Part I: Emergence and scaling laws for SGD learning

Yunwei Ren^{*1}, Eshaan Nichani^{*1}, Denny Wu²³, Jason D. Lee¹ ¹Princeton University ²NYU ³Flatiron Institute

April 3, 2025

Neural Scaling Laws [Kaplan et al. 20] [Hoffmann et al. 22]. Increasing compute and data leads to *power-law decay* in the loss.

Functional form: $\mathcal{L} \propto D^{-\alpha} + N^{-\beta} + C$.

- *D* number of data points.
- *N* number of parameters.



Neural Scaling Laws [Kaplan et al. 20] [Hoffmann et al. 22]. Increasing compute and data leads to *power-law decay* in the loss.

Functional form: $\mathcal{L} \propto D^{-\alpha} + N^{-\beta} + C$.

- *D* number of data points.
- N number of parameters.

Emergent Capabilities [Wei et al. 22] [Ganguli et al. 22]. Learning of individual tasks (skills) exhibits *sharp transition* with scale.

③ Unpredictable skill acquisition time.



Question: How do we reconcile the *emergent behavior* in skill acquisition and the *smooth power-law decay* in the cumulative loss?

Question: How do we reconcile the *emergent behavior* in skill acquisition and the *smooth power-law decay* in the cumulative loss?

Hypothesis: Additive Model [Michaud et al. 24] [Nam et al. 24]

- Cumulative objective can be decomposed into a large number of distinct "skills", the learning of each exhibits *abrupt emergence*.
- □ Juxtaposition of numerous emergent learning curves at different timescales results in a predictable *power-law* rate in the loss.



Question: How do we reconcile the *emergent behavior* in skill acquisition and the *smooth power-law decay* in the cumulative loss?

Hypothesis: Additive Model [Michaud et al. 24] [Nam et al. 24]

- Cumulative objective can be decomposed into a large number of distinct "skills", the learning of each exhibits *abrupt emergence*.
- Juxtaposition of numerous emergent learning curves at different timescales results in a predictable *power-law* rate in the loss.



This work: theoretical justification for the *additive model* hypothesis in SGD learning of shallow NN.

Gaussian single-index model: $f_*(\mathbf{x}) = \sigma_*(\langle \mathbf{x}, \boldsymbol{\theta} \rangle), \ \mathbf{x} \sim \mathcal{N}(0, \mathbf{I}_d).$

- **D** Requires learning the direction $\theta \in \mathbb{R}^d$ and link function $\sigma_* : \mathbb{R} \to \mathbb{R}$.
 - Learning algorithm should adapt to low-dimensional structure.

Gaussian single-index model: $f_*(\mathbf{x}) = \sigma_*(\langle \mathbf{x}, \boldsymbol{\theta} \rangle), \ \mathbf{x} \sim \mathcal{N}(0, \mathbf{I}_d).$

D Requires learning the <u>direction</u> $\theta \in \mathbb{R}^d$ and <u>link function</u> $\sigma_* : \mathbb{R} \to \mathbb{R}$.

• Learning algorithm should adapt to low-dimensional structure.

Hermite expansion: $\sigma_*(z) = \sum_{j=0}^{\infty} \alpha_j^* \operatorname{He}_j(z), \ \alpha_j^* = \mathbb{E}[\sigma_*(z) \operatorname{He}_j(z)].$

Information exponent [Ben Arous et al. 2021] [Dudeja & Hsu 2018]

The information exponent of σ_* is $k = \text{IE}(\sigma_*) = \min\{k \in \mathbb{N}_+ : \alpha_k^* \neq 0\}.$

Intuition: amount of information in the gradient at random initialization.

Hermite expansion: $\sigma_*(z) = \sum_{j=0}^{\infty} \alpha_j^* \operatorname{He}_j(z), \ \alpha_j^* = \mathbb{E}[\sigma_*(z) \operatorname{He}_j(z)].$

Information exponent [Ben Arous et al. 2021] [Dudeja & Hsu 2018]

The information exponent of σ_* is $k = \text{IE}(\sigma_*) = \min\{k \in \mathbb{N}_+ : \alpha_k^* \neq 0\}.$

$$\begin{aligned} - \mathbb{E}[\nabla_{\boldsymbol{w}}\mathcal{L}(f_{\mathsf{N}\mathsf{N}})] &\approx \mathbb{E}[\nabla_{\boldsymbol{w}}(f_{\mathsf{N}\mathsf{N}}(\boldsymbol{x})f_{*}(\boldsymbol{x}))] \\ &= \boldsymbol{\theta} \cdot \mathbb{E}[\sigma'_{*}(\langle \boldsymbol{x}, \boldsymbol{\theta} \rangle)\sigma'(\langle \boldsymbol{x}, \boldsymbol{w} \rangle)] + \boldsymbol{w} \cdot \mathbb{E}[...] \quad \text{Stein's lemma} \end{aligned}$$

Hermite expansion: $\sigma_*(z) = \sum_{j=0}^{\infty} \alpha_j^* \operatorname{He}_j(z), \ \alpha_j^* = \mathbb{E}[\sigma_*(z) \operatorname{He}_j(z)].$

Information exponent [Ben Arous et al. 2021] [Dudeja & Hsu 2018]

The information exponent of σ_* is $k = \text{IE}(\sigma_*) = \min\{k \in \mathbb{N}_+ : \alpha_k^* \neq 0\}.$

$$\begin{aligned} -\mathbb{E}[\nabla_{\boldsymbol{w}}\mathcal{L}(f_{\mathsf{NN}})] &\approx \mathbb{E}[\nabla_{\boldsymbol{w}}(f_{\mathsf{NN}}(\boldsymbol{x})f_{*}(\boldsymbol{x}))] \\ &= \boldsymbol{\theta} \cdot \mathbb{E}[\sigma'_{*}(\langle \boldsymbol{x}, \boldsymbol{\theta} \rangle)\sigma'(\langle \boldsymbol{x}, \boldsymbol{w} \rangle)] + \boldsymbol{w} \cdot \mathbb{E}[...] \quad Stein's \ lemma \\ &= \boldsymbol{\theta} \cdot \sum_{j=0}^{\infty} (j+1)^{2} \alpha^{*}_{j+1}\beta_{j+1} \underbrace{\langle \boldsymbol{w}, \boldsymbol{\theta} \rangle^{j}}_{\boldsymbol{d}=j/2} + ... \quad Hermite \ expansion \end{aligned}$$

Hermite expansion: $\sigma_*(z) = \sum_{j=0}^{\infty} \alpha_j^* \operatorname{He}_j(z), \ \alpha_j^* = \mathbb{E}[\sigma_*(z) \operatorname{He}_j(z)].$

Information exponent [Ben Arous et al. 2021] [Dudeja & Hsu 2018]

The information exponent of σ_* is $k = \text{IE}(\sigma_*) = \min\{k \in \mathbb{N}_+ : \alpha_k^* \neq 0\}.$

$$\begin{aligned} -\mathbb{E}[\nabla_{\boldsymbol{w}}\mathcal{L}(f_{\mathsf{NN}})] &\approx \mathbb{E}[\nabla_{\boldsymbol{w}}(f_{\mathsf{NN}}(\boldsymbol{x})f_{*}(\boldsymbol{x}))] \\ &= \boldsymbol{\theta} \cdot \mathbb{E}[\sigma'_{*}(\langle \boldsymbol{x}, \boldsymbol{\theta} \rangle)\sigma'(\langle \boldsymbol{x}, \boldsymbol{w} \rangle)] + \boldsymbol{w} \cdot \mathbb{E}[...] \quad Stein's \ lemma \\ &= \boldsymbol{\theta} \cdot \sum_{j=0}^{\infty} (j+1)^{2} \alpha^{*}_{j+1} \beta_{j+1} \underbrace{\langle \boldsymbol{w}, \boldsymbol{\theta} \rangle^{j}}_{\boldsymbol{d}^{-j/2} \ \text{at initialization}} + ... \quad Hermite \ expansion \end{aligned}$$

Theorem ([Ben Arous et al. 21], [Bietti et al. 22], [Damian et al. 23]...) A two-layer neural network can learn single-index f_* with information exponent k using $n \simeq T \gtrsim d^{\Theta(k)}$ samples and SGD steps.

Hermite expansion: $\sigma_*(z) = \sum_{j=0}^{\infty} \alpha_j^* \operatorname{He}_j(z), \ \alpha_j^* = \mathbb{E}[\sigma_*(z) \operatorname{He}_j(z)].$

Information exponent [Ben Arous et al. 2021] [Dudeja & Hsu 2018] The information exponent of σ_* is $k = \text{IE}(\sigma_*) = \min\{k \in \mathbb{N}_+ : \alpha_k^* \neq 0\}$.

Most samples in SGD are used to escape from the high-entropy "equator".



Emergent learning curve

- **Search Phase.** Online SGD exhibits loss plateau up to $T \simeq d^{k-1}$ steps.

Learning Extensive-width Neural Network

<u>Target Function</u>: Width-*P* two-layer neural network $f_*(\mathbf{x}) = \sum_{m=1}^{P} a_p \cdot \sigma(\langle \mathbf{x}, \boldsymbol{\theta}_p \rangle), \quad \sum a_p^2 = 1, \ \{\boldsymbol{\theta}_p\}_{p=1}^{P} \text{ orthonormal.}$

Motivation: learning sum of (orthogonal) single-index tasks.

- σ with IE $k > 2 \Rightarrow$ emergent learning curve for each direction.
- Varying second-layer $\{a_p\}_{p=1}^P \Rightarrow$ separation in length of search phase.

Learning Extensive-width Neural Network

 $\underline{\text{Target Function: Width-}P \text{ two-layer neural network}}_{f_*}(\mathbf{x}) = \sum_{m=1}^{P} a_p \cdot \sigma(\langle \mathbf{x}, \boldsymbol{\theta}_p \rangle), \quad \sum a_p^2 = 1, \ \{\boldsymbol{\theta}_p\}_{p=1}^{P} \text{ orthonormal.}}$

- **D** Extensive width. $P \simeq d^{\gamma}$ for $\gamma \in [0, c], c > 0$.
 - Large number of tasks \Rightarrow infinite-dimensional *effective dynamics*.
- **□** Large condition number. $a_{\max}/a_{\min} = \Theta(\operatorname{Poly}(P))$.
 - Covers power-law decay in second-layer $a_p \simeq p^{-\beta}, \beta \in [0, \infty)$.
- **Single-phase SGD learning** with MSE loss.
 - Avoids layer-wise training which alters scaling law.

Learning Extensive-width Neural Network

<u>Target Function</u>: Width-*P* two-layer neural network $f_*(\mathbf{x}) = \sum_{m=1}^{P} a_p \cdot \sigma(\langle \mathbf{x}, \boldsymbol{\theta}_p \rangle), \quad \sum a_p^2 = 1, \ \{\boldsymbol{\theta}_p\}_{p=1}^{P} \text{ orthonormal}.$

- **D** Extensive width. $P \simeq d^{\gamma}$ for $\gamma \in [0, c], c > 0$.
 - Large number of tasks \Rightarrow infinite-dimensional *effective dynamics*.
- **□** Large condition number. $a_{\max}/a_{\min} = \Theta(\operatorname{Poly}(P))$.
 - Covers power-law decay in second-layer $a_p \simeq p^{-\beta}, \beta \in [0, \infty)$.
- **Single-phase SGD learning** with MSE loss.
 - Avoids layer-wise training which alters scaling law.

This work: optimization and sample complexity for SGD training: • sharp emergence time for each θ_p , • smooth scaling law for MSE.

Student model: 2-homogeneous two-layer NN (matching σ)

$$f(\boldsymbol{x}) = \sum_{i=1}^{m} \|\boldsymbol{w}_i\|_2^2 \cdot \sigma(\langle \boldsymbol{x}, \boldsymbol{w}_i / \|\boldsymbol{w}_i\|\rangle).$$

Student model: 2-homogeneous two-layer NN (matching σ)

$$f(\mathbf{x}) = \sum_{i=1}^{m} \|\mathbf{w}_i\|_2^2 \cdot \sigma(\langle \mathbf{x}, \mathbf{w}_i / \|\mathbf{w}_i\|\rangle).$$

Online SGD training: at time *t*, compute MSE gradient update

$$\boldsymbol{w}_i(t+1) = \boldsymbol{w}_i(t) - \eta \nabla_{\boldsymbol{w}_i}(f_*(\boldsymbol{x}_t) - f(\boldsymbol{x}_t))^2, \quad \boldsymbol{x}_t \sim \mathcal{N}(0, \boldsymbol{I}_d).$$

Student model: 2-homogeneous two-layer NN (matching σ)

$$f(\mathbf{x}) = \sum_{i=1}^{m} \|\mathbf{w}_i\|_2^2 \cdot \sigma(\langle \mathbf{x}, \mathbf{w}_i / \|\mathbf{w}_i\|\rangle).$$

Online SGD training: at time *t*, compute MSE gradient update

$$\boldsymbol{w}_i(t+1) = \boldsymbol{w}_i(t) - \eta \nabla_{\boldsymbol{w}_i} (f_*(\boldsymbol{x}_t) - f(\boldsymbol{x}_t))^2, \quad \boldsymbol{x}_t \sim \mathcal{N}(0, \boldsymbol{I}_d).$$

Theorem (Complexity of SGD learning)

Assume $P \ll d^{0.1}$, $\bar{p} < P$. If we train a student NN with $m = \tilde{\Theta}(\bar{p})$ neurons using online SGD with $\eta \simeq \frac{a_{\bar{p}}}{d^{k/2} \text{poly}(P)}$, then w.h.p.,

 $\Box \text{ If } p < \bar{p}, \text{ alignment with } \theta_p \text{ emerges at } T_p \asymp a_p^{-1} \cdot \eta^{-1} d^{k/2-1}.$

J Student model: 2-homogeneous two-layer NN (matching σ)

$$f(\mathbf{x}) = \sum_{i=1}^{m} \|\mathbf{w}_i\|_2^2 \cdot \sigma(\langle \mathbf{x}, \mathbf{w}_i / \|\mathbf{w}_i\|\rangle).$$

Online SGD training: at time *t*, compute MSE gradient update

$$\boldsymbol{w}_i(t+1) = \boldsymbol{w}_i(t) - \eta \nabla_{\boldsymbol{w}_i} (f_*(\boldsymbol{x}_t) - f(\boldsymbol{x}_t))^2, \quad \boldsymbol{x}_t \sim \mathcal{N}(0, \boldsymbol{I}_d).$$

Theorem (Complexity of SGD learning)

Assume $P \ll d^{0.1}$, $\bar{p} < P$. If we train a student NN with $m = \tilde{\Theta}(\bar{p})$ neurons using online SGD with $\eta \simeq \frac{a_{\bar{p}}}{d^{k/2} \text{poly}(P)}$, then w.h.p.,

 $\Box \text{ If } p < \bar{p}, \text{ alignment with } \theta_p \text{ emerges at } T_p \asymp a_p^{-1} \cdot \eta^{-1} d^{k/2-1}.$

 \square All directions up to \bar{p} are learned at $n \asymp T \asymp a_{\bar{p}}^{-2} d^{k-1} \mathrm{poly}(P)$.

Theorem (Complexity of SGD learning)

Assume $P \ll d^{0.1}$, $\bar{p} < P$. If we train a student NN with $m = \tilde{\Theta}(\bar{p})$ neurons using online SGD with $\eta \simeq \frac{a_{\bar{p}}}{d^{k/2} \text{poly}(P)}$, then w.h.p.,

 $\Box \text{ If } p < \bar{p}, \text{ alignment with } \theta_p \text{ emerges at } T_p \asymp a_p^{-1} \cdot \eta^{-1} d^{k/2-1}.$

 $\square \text{ All directions up to } \bar{p} \text{ are learned at } n \asymp T \asymp a_{\bar{p}}^{-2} d^{k-1} \mathrm{poly}(P) \text{ .}$

Corollary: ignoring logarithmic factors, achieving small population loss (i.e., learning *all* tasks) requires $m \asymp P$ neurons and $n \asymp T \asymp a_{\min}^{-2} d^{k-1} \text{poly}(P)$.

Theorem (Complexity of SGD learning)

Assume $P \ll d^{0.1}$, $\bar{p} < P$. If we train a student NN with $m = \tilde{\Theta}(\bar{p})$ neurons using online SGD with $\eta \simeq \frac{a_{\bar{p}}}{d^{k/2} \text{poly}(P)}$, then w.h.p.,

 $\Box \text{ If } p < \bar{p}, \text{ alignment with } \theta_p \text{ emerges at } T_p \asymp a_p^{-1} \cdot \eta^{-1} d^{k/2-1}.$

 $\square \text{ All directions up to } \bar{p} \text{ are learned at } n \asymp T \asymp a_{\bar{p}}^{-2} d^{k-1} \mathrm{poly}(P) \text{ .}$

Corollary: ignoring logarithmic factors, achieving small population loss (i.e., learning *all* tasks) requires $m \simeq P$ neurons and $n \simeq T \simeq a_{\min}^{-2} d^{k-1} \operatorname{poly}(P)$.

- Prior works required sample size or compute *exponentially* large in the condition number: n, P ≍ exp (^{amax}/_{amin}) [Li et al. 22] [Oko et al. 24].
- Our learning procedure does not involve reinitialization [Ge et al. 21] or Stiefel constraint [Ben Arous et al. 24].

Neural Scaling Laws for Feature Learning

Corollary (Emergence & Scaling Laws)

Assume $a_p \asymp p^{-\beta}$ for $\beta > 1/2$ and fixed learning rate η , then

- 1. Emergence: the p-th task is learned at time $\eta t \sim d^{k/2-1}p^{\beta}$.
- 2. Scaling law: population MSE exhibits power-law decay

 $\mathcal{L}(t) \sim m^{1-2\beta} \vee \left(\eta t d^{1-k/2}\right)^{rac{1-2\beta}{eta}}.$



Optimization time t

Neural Scaling Laws for Feature Learning

Corollary (Emergence & Scaling Laws)

Assume $a_p \asymp p^{-\beta}$ for $\beta > 1/2$ and fixed learning rate η , then

- 1. Emergence: the p-th task is learned at time $\eta t \sim d^{k/2-1}p^{\beta}$.
- 2. Scaling law: population MSE exhibits power-law decay

$$\mathcal{L}(t) \sim m^{1-2\beta} \vee \left(\eta t d^{1-k/2}\right)^{\frac{1-2\beta}{\beta}}$$

Intuition: Assume *decoupled* learning of different teacher components,

- Direction θ_p learned at $T_p \propto p^{\beta} \eta^{-1} d^{k/2-1}$.
- $\mathcal{L}(t) \approx \sum_{p=1}^{P} a_p^2 \mathbb{I}\{t < T_p\} \Rightarrow$ $\mathcal{L}(T_p) \approx \int_p^{\infty} s^{-2\beta} \mathrm{d}s \sim p^{1-2\beta}.$
- m < P student neurons reaches $\sum_{s>m} a_s^2 \asymp m^{1-2\beta}$ error.



Neural Scaling Laws for Feature Learning

- $\square m^{1-2\beta}$ approximation barrier.
 - Determined by the student network width.
- $\Box \left(\eta t d^{1-k/2}\right)^{\frac{1-2\beta}{\beta}} \text{optimization error.}$
 - Determined by the number of online SGD steps.



Remark on Discretization

Note: the continuous-time (or constant- η) rate can be misleading!

 \odot Exponent does not match known rates for weak ℓ_p ball [Johnstone 17].

Adaptive learning rate for *p***-th task.** Consider only learning the top-*p* neurons.

• "Optimal" step size for the *p*-th task: $\eta \simeq \frac{a_p}{d^{k/2} \text{poly}(P)}$, $a_p \simeq p^{-\beta}$.

• Direction
$$\theta_p$$
 now learned at
 $n = T_p \sim p^{\beta} \eta^{-1} d^{k/2-1} = p^{2\beta} d^{k-1} \text{poly}(P).$

Sample complexity: under this learning rate for the first p neurons,

$$\mathcal{L}(n,m) \sim \left(\frac{n}{d^{k-1}\mathrm{poly}(P)}\right)^{\frac{1-2\beta}{2\beta}} + m^{1-2\beta}$$

 \bigcirc This now matches the known rates for weak ℓ_p ball.

Proof Sketch

• Decoupled gradient dynamics.

Let $\boldsymbol{w}_{\iota(p)}$ denote the neuron that eventually converged to direction $\boldsymbol{\theta}_p$, and $m_{p,q}(t) = \langle \boldsymbol{w}_p(t) / \| \boldsymbol{w}_p(t) \|, \boldsymbol{\theta}_q \rangle$ measures the overlap at time t.

- Claim 1 decoupling. When all $m_{\iota(p),q}^2$ $(q \neq p)$ are small, the learning of different directions can be approximately decoupled.
- Claim 2 sharp transitions. Since the norm of w_p grows rapidly after alignment is achieved, all $m_{\iota(p),q}^2$ $(q \neq p)$ remain small when the *q*-th direction is recovered by $w_{\iota(q)}$ after which θ_q no longer affects the learning dynamics ("automatic deflation").

Proof Sketch

• Decoupled gradient dynamics.

Let $\boldsymbol{w}_{\iota(p)}$ denote the neuron that eventually converged to direction $\boldsymbol{\theta}_p$, and $m_{p,q}(t) = \langle \boldsymbol{w}_p(t) / \| \boldsymbol{w}_p(t) \|, \boldsymbol{\theta}_q \rangle$ measures the overlap at time t.

- Claim 1 decoupling. When all $m_{\iota(p),q}^2$ $(q \neq p)$ are small, the learning of different directions can be approximately decoupled.
- Claim 2 sharp transitions. Since the norm of w_p grows rapidly after alignment is achieved, all $m_{\iota(p),q}^2$ $(q \neq p)$ remain small when the *q*-th direction is recovered by $w_{\iota(q)}$ after which θ_q no longer affects the learning dynamics ("automatic deflation").

• From gradient flow to online SGD.

- Martingale decomposition in [Ben Arous et al. 21] + a refined stochastic induction argument from [Ren & Lee 24].
- "Unstable" discretization that couples the online SGD dynamics for learning the top-*p* teacher neurons.

Remarks and Future Directions

- **Edge Case:** information exponent k = 2 (e.g., quadratic σ).
 - © Dynamics of different directions cannot be decoupled.
- □ *Companion work:* introduce algorithmic modification (Stiefel SGD) and handle the special case of quadratic activation.
 - Gerard Ben Arous, Murat A. Erdogdu, N. Mert Vural, Denny Wu. "Scaling laws for learning quadratic two-layer neural networks with SGD in high dimensions".

Remarks and Future Directions

- **Edge Case:** information exponent k = 2 (e.g., quadratic σ).
 - © Dynamics of different directions cannot be decoupled.
- □ *Companion work:* introduce algorithmic modification (Stiefel SGD) and handle the special case of quadratic activation.
 - Gerard Ben Arous, Murat A. Erdogdu, N. Mert Vural, Denny Wu. "Scaling laws for learning quadratic two-layer neural networks with SGD in high dimensions".

Future Directions

- **D** Anisotropic input: $\mathbf{x} \sim \mathcal{N}(0, \mathbf{\Sigma})$, where $\mathbf{\Sigma} = \sum_{i=1}^{d} \lambda_i \mathbf{v}_i \mathbf{v}_i^{\top}, \lambda_i \asymp i^{-\alpha}$.
 - Two-parameter scaling law? (source & capacity conditions)
- □ Beyond additive structure ⇒ scaling laws for compositional generalization? (More on this in Part II)

Part II: Learning Compositional Functions with Transformers

Zixuan Wang*¹, Eshaan Nichani*¹, Alberto Bietti³, Alex Damian¹, Daniel Hsu⁴, Jason D. Lee¹, Denny Wu²³

¹Princeton University ²NYU ³Flatiron Institute ⁴Columbia University

April 3, 2025

Motivation: multi-step reasoning tasks

Large language models exhibits remarkable reasoning capabilities.

• They can solve complex tasks that require combining multiple reasoning steps. [Wei et al., 2022]



Multi-hop QA



Coding



Examples (Multi-hop QA)

Prompt: "John is in the playground. [...] Helen is playing with John. [...] Helen picked up the football. [...] Where is the football?" **Answer:** "In the playground."

Multi-step reasoning tasks can be viewed as function composition.

- **Prompt:** "John is in the playground. [...] Helen is playing with John. [...] Helen picked up the football. [...] Where is the football?" **Answer:** "In the playground."
- $f_{\text{location}}(\text{John}) = \text{playground}, f_{\text{with}}(\text{Helen}) = \text{John}, f_{\text{hold}}(\text{football}) = \text{Helen}$
- The model needs to compute $f_{\text{location}} \circ f_{\text{with}} \circ f_{\text{hold}}(\text{football})$ and get playground.
- Note: (location, with, hold) are also function inputs.

Multi-step reasoning tasks can be viewed as function composition.

- **Prompt:** "John is in the playground. [...] Helen is playing with John. [...] Helen picked up the football. [...] Where is the football?" **Answer:** "In the playground."
- $f_{\text{location}}(\text{John}) = \text{playground}, f_{\text{with}}(\text{Helen}) = \text{John}, f_{\text{hold}}(\text{football}) = \text{Helen}$
- The model needs to compute f_{location} o f_{with} o f_{hold} (football) and get playground.
- Note: (location, with, hold) are also function inputs.

Abstractly, the model can do the following function composition:

- Inputs: $(\sigma_1, \sigma_2, ..., \sigma_k, x)$
- Outputs: $f_{\sigma_k}^{(k)} \circ f_{\sigma_{k-1}}^{(k-1)} \cdots f_{\sigma_1}^{(1)}(x)$ where the hidden function is position dependent.

How does transformer compose functions?

With Chain-of-Thought (CoT):

- Expressivity [Li et al., 2024]
- Require explicit *k* inference steps to compose *k* functions.
- Inefficient for simple tasks?

With Chain-of-Thought (CoT):

- Expressivity [Li et al., 2024]
- Require explicit *k* inference steps to compose *k* functions.
- Inefficient for simple tasks?

Without Chain-of-Thought:

- Limited expressivity [Merrill and Sabharwal, 2023]
- Only a single inference step.
- Potentially faster?

How does transformer compose functions?

With Chain-of-Thought (CoT):

- Expressivity [Li et al., 2024]
- Require explicit *k* inference steps to compose *k* functions.
- Inefficient for simple tasks?

Without Chain-of-Thought:

- Limited expressivity [Merrill and Sabharwal, 2023]
- Only a single inference step.
- Potentially faster?

Focus of this talk:

- Can transformers compose (some) functions internally?
- Can they learn composition efficiently via gradient-based training?

Definition (Self-attention)

The self-attention head attn(·; W_{KQ}, W_{OV}) : $\mathbb{R}^{d \times T} \to \mathbb{R}^{d \times T}$ is, with $W_{KQ}, W_{OV} \in \mathbb{R}^{d \times d}$

$$\mathsf{attn}(X; W_{\mathcal{K}\mathcal{Q}}, W_{\mathcal{O}\mathcal{V}}) = X + W_{\mathcal{O}\mathcal{V}}X\mathcal{S}(X^ op W_{\mathcal{K}\mathcal{Q}}X),$$

where the softmax function ${\mathcal S}$ is applied column-wise.

- Convex combination of all values W_{OV}X with coefficient S(X[⊤]W_{KQ}X)_(i,j), i.e. attention score.
- A multi-layer transformer composes multiple self-attention layers.
- We focus on attention-only transformers now.



Transformers parallelize composition

Can transformers compose (some) functions internally?

- Yes! Enabled by parallelism.
- Composing k functions only requires $O(\log k)$ layers.

Transformers parallelize composition

Can transformers compose (some) functions internally?

- Yes! Enabled by parallelism.
- Composing k functions only requires $O(\log k)$ layers.

Examples (Efficiently parallelizable tasks)

• *k*-hop induction head [Sanford et al., 2024]; Finite state automata [Liu et al., 2023].

(2-hop Induction head)



(Finite state automata)

Definition. (*k*-fold composition)

- $\pi := (\pi_1, \ldots, \pi_k) \in (S_N)^k$ is a tuple of k hidden permutations.
- Input: $(\sigma, x) \in \mathcal{X} := (S_N)^k \times [N]$, where $\sigma = (\sigma_1, \dots, \sigma_k)$, x is an index in [N].
- Output $f_{\pi}: \mathcal{X} \rightarrow [N]$ is defined as

$$f_{\pi}(\sigma, x) := (\sigma_k \circ \pi_k \circ \sigma_{k-1} \circ \pi_{k-1} \circ \cdots \circ \sigma_1 \circ \pi_1)(x).$$



Remark: $\pi := (\pi_1, \ldots, \pi_k) \in (S_N)^k$ induces a function class \mathcal{F} .

Theorem (Expressivity of Transformers)

Assume k is a power of 2. For any $\pi \in (S_N)^k$, there exists an $L = \log_2 k + 1$ layer transformer with embedding dimension $\tilde{O}(kN)$ that expresses k-fold composition.

Theorem (Expressivity of Transformers)

Assume k is a power of 2. For any $\pi \in (S_N)^k$, there exists an $L = \log_2 k + 1$ layer transformer with embedding dimension $\tilde{O}(kN)$ that expresses k-fold composition.

- Layer 1 encodes hidden π_i 's.
- Layer 2 composes $\sigma_{i+1} \circ \pi_{i+1}$ and $\sigma_i \circ \pi_i$ to get

$$\mathsf{hop}_i^2 := \sigma_{i+1} \circ \pi_{i+1} \circ \sigma_i \circ \pi_i$$

• Layer ℓ composes $hop_i^{2^{l-2}}$ and $hop_{i+2^{l-2}}^{2^{l-2}}$ and get $hop_i^{2^{l-1}}$.



The parallel solution (Layer 1)

• Embedding for each $(i,j) \in [N] \times [k]$ stores one-hot embedding of (i,j) and $(i, \sigma_i(j))$

Input $X^{(0)}$:

σ_1						
	(1,1)	(1,2)		(1,N)		
Ī	$\sigma_{l}(1)$	$\sigma_1(2)$		$\sigma_1(N)$		
Ī	0 _{kNL}	0 _{kNL}		0 _{kNL}		

	σ_i			
	(<i>i</i> ,1)	(i,2)		(<i>i</i> , <i>N</i>)
•••	$\sigma_i(1)$	σ _i (2)		$\sigma_i(N)$
	0 _{kNL}	0 _{kNL}		0 _{kNL}
	KNL	KINL		KNL

	σ_k				
	(k,1)	(k,2)		(k, N)	
	$\sigma_k(1)$	$\sigma_k(2)$		$\sigma_k(N)$	
	0 _{kNL}	0 _{kNL}		0 _{kNL}	

- The key-query matrix makes position (i, j) attend to position $(i, \pi(j))$.
- The value matrix copies hop_i¹(j) := (σ_iπ_i)(j) into the residual stream of (i, j)
- The first layer computes 1-hop!



The parallel solution (Layer $\ell + 1$)

- Assume that the previous layer ℓ computes hop_i^{2^{ℓ-1}} (i.e residual stream of (i, j) contains hop_i^{2^{ℓ-1}}(j)).
- The key-query matrix makes position (i, j) attend to position (i + 2^{l-1}, hop<sub>i<sup>2^{l-1}</sub>(j)).
 </sub></sup>
- Value matrix copies $\operatorname{hop}_{i+2^{\ell-1}}^{2^{\ell-1}}(\operatorname{hop}_{i}^{2^{\ell-1}}(j))$ into residual stream of (i,j).
- Thus the $(\ell + 1)$ th layer computes

 $\operatorname{hop}_{i}^{2^{\ell}} = \operatorname{hop}_{i+2^{\ell-1}}^{2^{\ell-1}} \circ \operatorname{hop}_{i}^{2^{\ell-1}}.$

• Finally, the *L*-th layer outputs *k*-hop.



Computational hardness of *k*-fold composition

Is this parallel construction for the k-fold composition efficiently learnable by GD?

Computational hardness of *k*-fold composition

Is this parallel construction for the k-fold composition efficiently learnable by GD?

- Negative result: Statistical Query lower bound.
- SQ framework: Learner specifies query q : X × [N] → ℝ and tolerance τ. SQ oracle responds with ĝ, a noisy estimate of q

$$|\hat{q} - \mathbb{E}_{\sigma,x}[q(\sigma,x,f_{\pi}(\sigma,x))]| \leq au$$

Computational hardness of *k*-fold composition

Is this parallel construction for the k-fold composition efficiently learnable by GD?

- Negative result: Statistical Query lower bound.
- SQ framework: Learner specifies query q : X × [N] → ℝ and tolerance τ. SQ oracle responds with ĝ, a noisy estimate of q

$$|\hat{q} - \mathbb{E}_{\sigma,x}[q(\sigma, x, f_{\pi}(\sigma, x))]| \leq au$$

Theorem (SQ lower bound for *k*-fold composition)

Any SQ learner for the k-fold composition function class \mathcal{F} must either make $q \geq N^{\Omega(k)}$ queries or use a tolerance $\tau \leq N^{-\Omega(k)}$ to output a predictor \hat{f} with loss $L(\hat{f}) \lesssim 1$.

• $\tau \approx n^{-1/2}$ by i.i.d. concentration heuristic, where *n* is the sample size. \Rightarrow Either runtime (\geq number of queries) or sample size must be $N^{\Omega(k)}$.

Easy-to-hard: efficiency of curriculum learning

However, a curriculum can guide the transformer to learn the solution efficiently!

- Stage 1: GD on the data with label hop¹_i = $\sigma_i \pi_i$ for $i \in [k]$.
- Stage ℓ : GD on the data with label hop_i^{2^{ℓ-1}}.
- We can learn the k-fold function with poly(N, k) samples!

Easy-to-hard: efficiency of curriculum learning

However, a curriculum can guide the transformer to learn the solution efficiently!

- Stage 1: GD on the data with label hop¹_i = $\sigma_i \pi_i$ for $i \in [k]$.
- Stage ℓ : GD on the data with label hop_i^{2^{ℓ-1}}.
- We can learn the k-fold function with poly(N, k) samples!

Theorem (GD upper bound with curriculum)

Assume $k = 2^{L-1}$, if the sample size $n \ge \tilde{\Omega}(k^4 N^6)$, the output of gradient descent with curriculum learning learns the k-fold composition function.

Easy-to-hard: efficiency of curriculum learning

However, a curriculum can guide the transformer to learn the solution efficiently!

- Stage 1: GD on the data with label hop¹_i = $\sigma_i \pi_i$ for $i \in [k]$.
- Stage ℓ : GD on the data with label hop_i^{2^{ℓ-1}}.
- We can learn the k-fold function with poly(N, k) samples!

Theorem (GD upper bound with curriculum)

Assume $k = 2^{L-1}$, if the sample size $n \ge \tilde{\Omega}(k^4 N^6)$, the output of gradient descent with curriculum learning learns the k-fold composition function.

<u>Takeaway:</u> $n \asymp N^{\Omega(k)}$ (Exponential) $\Rightarrow n \asymp poly(N, k)$ (Polynomial).

 $\ensuremath{\textcircled{}}$ Easy examples/intermediate supervision is essential for complex tasks!

Stage 1:

- In the construction, the first layer makes the (i, j) position attend to $(i, \pi_i(j))$ position.
- With 1-hop data, learning this first layer is easy!
- Gradient descent can correctly learn the constructed attention pattern, i.e. recover the hidden permutations π_i 's.

Stage 1:

- In the construction, the first layer makes the (i, j) position attend to $(i, \pi_i(j))$ position.
- With 1-hop data, learning this first layer is easy!
- Gradient descent can correctly learn the constructed attention pattern, i.e. recover the hidden permutations π_i 's.

Stage ℓ ($\ell \ge 2$):

- ℓ th layer learns to compose hop_i^{2^{ℓ-2}} and hop_{i+2^{ℓ-2}}^{2^{ℓ-2}} computed in the previous layer.
- With $2^{\ell-1}$ -hop data, GD correctly makes (i, j) position attend to $(i + 2^{\ell-2}, \operatorname{hop}_{i}^{2^{\ell-2}}(j))$ position, thus matching the construction.

Learning with data mixture

In practice, designing explicit curriculum is challenging:

- Hard to ensure increasing difficulty; facing catastrophic forgetting
- Instead, practitioners use data mixture with different difficulties

Learning with data mixture

In practice, designing explicit curriculum is challenging:

- Hard to ensure increasing difficulty; facing catastrophic forgetting
- Instead, practitioners use data mixture with different difficulties

Training on data mixture also works!

Theorem (Upper bound on mixed data, informal)

Assume $k = 2^{L-1}$, if the sample size $n \ge \tilde{\Omega}(k^4 N^6)$, the output of gradient descent training on the mixture of $\{1, 2, 4, ..., k\}$ -fold data learns the k-fold composition function.

In practice, designing explicit curriculum is challenging:

- Hard to ensure increasing difficulty; facing catastrophic forgetting
- Instead, practitioners use data mixture with different difficulties

Training on data mixture also works!

Theorem (Upper bound on mixed data, informal)

Assume $k = 2^{L-1}$, if the sample size $n \ge \tilde{\Omega}(k^4 N^6)$, the output of gradient descent training on the mixture of $\{1, 2, 4, ..., k\}$ -fold data learns the k-fold composition function.

- Intuition: if layer ℓ (the 2^{ℓ-1}-fold function) is not learned, GD cannot learn later layers with ℓ' > ℓ — It is essential to learn from easy to hard!
- Data mixture implicitly induces curriculum learning.

Experiments

Simulation on normal transformers with MLP:

- Without data mixture: Transformer cannot learn the 4-hop task.
- With mixture: Transformer learns from easy to hard!
- Learning sequentially from easy (1-hop) to hard tasks (4-hop).



Remarks

Connection to *k*-sparse parity

- k-fold composition exhibits a statistical-computational gap. Information theoretic sample complexity is n ≍ kN log N, while a computationally efficient SQ learner requires n ≥ N^{Ω(k)}.
- Mirrors k-parity in d dimensions (requires d^k samples/compute).

Connection to *k*-sparse parity

- k-fold composition exhibits a statistical-computational gap. Information theoretic sample complexity is n ≍ kN log N, while a computationally efficient SQ learner requires n ≥ N^{Ω(k)}.
- Mirrors k-parity in d dimensions (requires d^k samples/compute).
- In both settings, efficient gradient-based learning requires curriculum/intermediate supervision [Abbe et al., 2023, Wen et al., 2024, Kim and Suzuki, 2024].
- Unlike *k*-parity, *k*-fold composition *cannot be expressed by a shallow neural network*.

Connection to *k*-sparse parity

- k-fold composition exhibits a statistical-computational gap. Information theoretic sample complexity is n ≍ kN log N, while a computationally efficient SQ learner requires n ≥ N^{Ω(k)}.
- Mirrors k-parity in d dimensions (requires d^k samples/compute).
- In both settings, efficient gradient-based learning requires curriculum/intermediate supervision [Abbe et al., 2023, Wen et al., 2024, Kim and Suzuki, 2024].
- Unlike *k*-parity, *k*-fold composition *cannot be expressed by a shallow neural network*.

Future directions

- How to learn general compositions?
- Learning a more efficient embedding.

References

Abbe, E., Cornacchia, E., and Lotfi, A. (2023).

Provable advantage of curriculum learning on parity targets with mixed inputs.

Advances in Neural Information Processing Systems, 36:24291–24321.

Kim, J. and Suzuki, T. (2024).

Transformers provably solve parity efficiently with chain of thought. arXiv preprint arXiv:2410.08633.

- Li, Z., Liu, H., Zhou, D., and Ma, T. (2024). Chain of thought empowers transformers to solve inherently serial problems. *arXiv preprint arXiv:2402.12875.*

5

Liu, B., Ash, J. T., Goel, S., Krishnamurthy, A., and Zhang, C. (2023). Transformers learn shortcuts to automata.

In The Eleventh International Conference on Learning Representations.

- Merrill, W. and Sabharwal, A. (2023).
 The parallelism tradeoff: Limitations of log-precision transformers.
 Transactions of the Association for Computational Linguistics, 11:531–545.
 - Sanford, C., Hsu, D., and Telgarsky, M. (2024).