

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 & Emergence

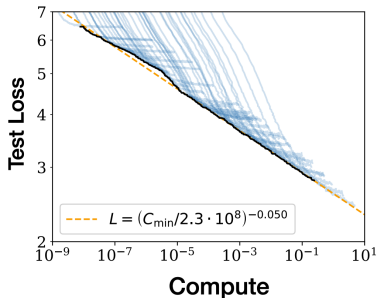
Neural Scaling Laws [Kaplan et al. 20] [Hofmann et al. 22].

Increasing compute and data leads to power-law decay in the loss.

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

D { number of data points.

N { number of parameters.



Neural Scaling Laws & Emergence

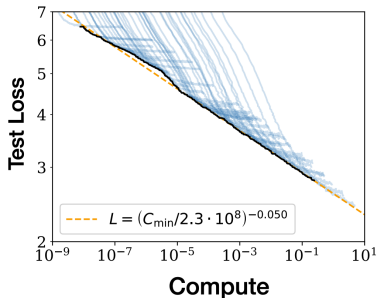
Neural Scaling Laws [Kaplan et al. 20] [Hochmann et al. 22].

Increasing compute and data leads to power-law decay in the loss.

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

D { number of data points.

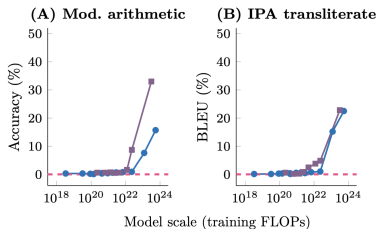
N { number of parameters.



Emergent Capabilities [Wei et al. 22] [Ganguli et al. 22].

Learning of individual tasks (skills) exhibits sharp transition with scale.

/ Unpredictable skill acquisition time.



Neural Scaling Laws & Emergence

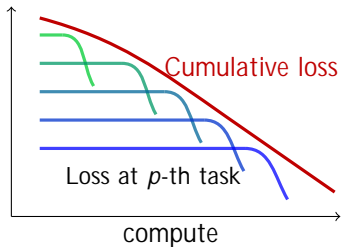
Question: How do we reconcile the *emergent behavior* in skill acquisition and the *smooth power-law decay* in the cumulative loss?

Neural Scaling Laws & Emergence

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.

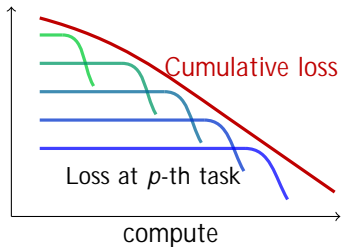


Neural Scaling Laws & Emergence

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.

Emergence in Gradient-based Feature Learning

Gaussian single-index model: $f(x) = (hx; \theta); x \sim N(0; I_d)$.

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

Emergence in Gradient-based Feature Learning

Gaussian single-index model: $f(x) = \langle hx; \theta \rangle; x \sim N(0; I_d)$.

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

Hermite expansion: $f(z) = \sum_{j=0}^{\infty} \hat{f}_j \text{He}_j(z), \hat{f}_j = \mathbb{E}[f(z)\text{He}_j(z)]$.

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

The information exponent of f is $k = \text{IE}(f) = \min\{k \in \mathbb{N}_+ : \hat{f}_k \neq 0\}$.

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

Emergence in Gradient-based Feature Learning

Hermite expansion: $f(z) = \sum_{j=0}^{\infty} \gamma_j \text{He}_j(z), \quad \gamma_j = \mathbb{E}[f(z)\text{He}_j(z)].$

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

The information exponent of f is $k = \text{IE}(f) = \min\{k \in \mathbb{N}_+ : \gamma_k \neq 0\}.$

$$\begin{aligned} \mathbb{E}[r_w L(f_{\text{NN}})] &= \mathbb{E}[r_w (f_{\text{NN}}(x) - f(x))] \\ &= \mathbb{E}[\partial^k (hx; i) \partial^k (hx; wi)] + w \mathbb{E}[\dots] \quad \text{Stein's lemma} \end{aligned}$$

Emergence in Gradient-based Feature Learning

Hermite expansion: $f(z) = \sum_{j=0}^{\infty} \gamma_j \text{He}_j(z), \quad \gamma_j = \mathbb{E}[f(z)\text{He}_j(z)].$

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

The information exponent of f is $k = \text{IE}(f) = \min\{k \geq 0 : \gamma_k \neq 0\}.$

$$\begin{aligned}
 \mathbb{E}[r_w L(f_{\text{NN}})] &= \mathbb{E}[r_w (f_{\text{NN}}(x) - f(x))^2] \\
 &= \mathbb{E}[\sum_{j=0}^{\infty} \gamma_j^2 \text{He}_j(w x)^2] + w \mathbb{E}[\sum_{j=0}^{\infty} \gamma_j \text{He}_j(w x)] \quad \text{Stein's lemma} \\
 &= \sum_{j=0}^{\infty} (j+1)^2 \gamma_{j+1}^2 \underbrace{\mathbb{E}[\text{He}_j(w x)^2]}_{\substack{d \\ j=2 \text{ at initialization}}} + \dots \quad \text{Hermite expansion}
 \end{aligned}$$

Emergence in Gradient-based Feature Learning

Hermite expansion: $f(z) = \sum_{j=0}^{\infty} c_j \text{He}_j(z)$, $c_j = \mathbb{E}[f(z) \text{He}_j(z)]$.

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

The information exponent of f is $k = \text{IE}(f) = \min\{k \geq 0 : c_k \neq 0\}$.

$$\begin{aligned} \mathbb{E}[r_w L(f_{\text{NN}})] &= \mathbb{E}[r_w (f_{\text{NN}}(x) - f(x))] \\ &= \mathbb{E}[\sum_{j=0}^{\infty} c_j (\text{He}_j(hx; i) - \text{He}_j(hx; wi))] + w \mathbb{E}[\dots] \quad \text{Stein's lemma} \\ &= \sum_{j=0}^{\infty} c_j (j+1)^2 \frac{w^j}{j!} + \dots \quad \text{Hermite expansion} \\ &\quad \text{with } c_j \text{ at initialization} \end{aligned}$$

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 $n \geq T$ & $d^{\Theta(k)}$ samples and SGD steps.

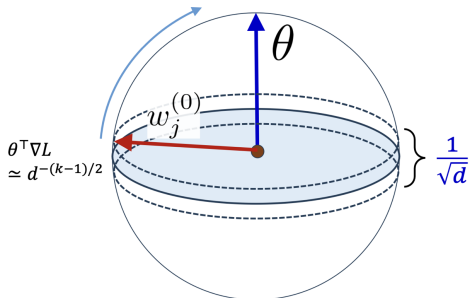
Emergence in Gradient-based Feature Learning

Hermite expansion: $f(z) = \sum_{j=0}^{\infty} \gamma_j \text{He}_j(z)$, $\gamma_j = \mathbb{E}[f(z) \text{He}_j(z)]$.

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

The information exponent of f is $k = \mathbb{E}(\gamma_j) = \min\{k \in \mathbb{N}_+ : \gamma_k \neq 0\}$.

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 $T \approx d^{k-1}$ steps.

⌘ **Descent Phase.** Loss sharply decreases in $T \approx d$ steps.

Learning Extensive-width Neural Network

Target Function: Width- P two-layer neural network

$$f(x) = \sum_{m=1}^P a_p (hx; p_i); \quad \sum_{p=1}^P a_p^2 = 1; \quad \{g_{p=1}^P\} \text{ orthonormal:}$$

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

with IE $k > 2$) *emergent* learning curve for each direction.

Varying second-layer $\{a_p g_{p=1}^P\}$) separation in length of search phase.

Learning Extensive-width Neural Network

Target Function: Width- P two-layer neural network

$$f(x) = \sum_{m=1}^P a_p (hx; p i); \quad \sum_{p=1}^P a_p^2 = 1; \quad \{g_{p=1}^P\} \text{ orthonormal.}$$

- p **Extensive width.** $P \gg d$ for $\lambda \in [0; c]; c > 0$.
 Large number of tasks λ in finite-dimensional *effective dynamics*.
- p **Large condition number.** $a_{\max} = a_{\min} = \text{Poly}(P)$.
 Covers **power-law decay** in second-layer $a_p \sim p^{-\lambda}; \lambda \in [0; 1)$.
- p **Single-phase SGD learning** with MSE loss.
 Avoids layer-wise training which alters scaling law.

Learning Extensive-width Neural Network

Target Function: Width- P two-layer neural network

$$f(x) = \sum_{m=1}^P a_p (hx; p_i); \quad \sum_{p=1}^P a_p^2 = 1; \quad \{g_{p=1}^P\} \text{ orthonormal.}$$

- p **Extensive width.** $P \gg d$ for $\lambda \in [0; c]; c > 0$.
 Large number of tasks λ in finite-dimensional *effective dynamics*.
- p **Large condition number.** $a_{\max} = a_{\min} = \text{Poly}(P)$.
 Covers **power-law decay** in second-layer $a_p \sim p^{-\lambda}; \lambda \in [0; 1)$.
- p **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 p , smooth scaling law for MSE.

Main Theorem: Complexity of SGD Learning

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

$$f(x) = \sum_{i=1}^m k w_i k_2^2 \quad (hx; w_i = k w_i k_i):$$

Main Theorem: Complexity of SGD Learning

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

$$f(x) = \sum_{i=1}^m k w_i k_2^2 (h x; w_i = k w_i k_i):$$

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

$$w_i(t+1) = w_i(t) - \eta \nabla_{w_i} (f(x_t) - f(x_t))^2; \quad x_t \sim N(0; I_d):$$

Main Theorem: Complexity of SGD Learning

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

$$f(x) = \sum_{i=1}^m k w_i k_2^2 (h x; w_i = k w_i k_i):$$

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

$$w_i(t+1) = w_i(t) - \eta \nabla_{w_i} (f(x_t) - f(x_t))^2; \quad x_t \sim N(0; I_d):$$

Theorem (Complexity of SGD learning)

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

- If $p < p$, alignment with p emerges at $T_p = a_p^{-1} d^{k=2-1}$.

Main Theorem: Complexity of SGD Learning

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

$$f(x) = \sum_{i=1}^m k w_i k_2^2 (h x; w_i = k w_i k_i):$$

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

$$w_i(t+1) = w_i(t) - \eta_{w_i} (f(x_t) - f(x_t))^2; \quad x_t \sim N(0; I_d):$$

Theorem (Complexity of SGD learning)

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

- If $p < p$, alignment with p emerges at $T_p = a_p^{-1} d^{k=2}^{-1}$.

- All directions up to p are learned at $n = T = a_p^{-2} d^{k=2}^{-1} \text{poly}(P)$.

Main Theorem: Complexity of SGD Learning

Theorem (Complexity of SGD learning)

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

• If $p < p$, alignment with p emerges at $T_p = a_p^{-1} d^{k=2} \text{poly}(P)$.

• All directions up to p are learned at $n = T = a_p^{-2} d^{k=2} \text{poly}(P)$.

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

Main Theorem: Complexity of SGD Learning

Theorem (Complexity of SGD learning)

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

• If $p < p$, alignment with p emerges at $T_p = a_p^{-1} d^{k=2} \text{poly}(P)$.

• All directions up to p are learned at $n = T_p a_p^{-2} d^{k=2} \text{poly}(P)$.

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

Prior works required sample size or compute *exponentially* large in the condition number: $n; P = \exp\left(\frac{a_{\max}}{a_{\min}}\right)$ [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].

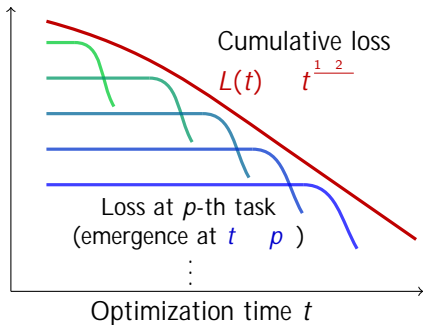
Neural Scaling Laws for Feature Learning

Corollary (Emergence & Scaling Laws)

Assume $a_p = p$ for $p > 1=2$ and fixed learning rate, then

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

$$L(t) \sim m^{1/2} - t d^{1/2} \sim t^{-1/2}.$$



Neural Scaling Laws for Feature Learning

Corollary (Emergence & Scaling Laws)

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

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

$$L(t) \propto m^{-1/2} t^{-1/2}.$$

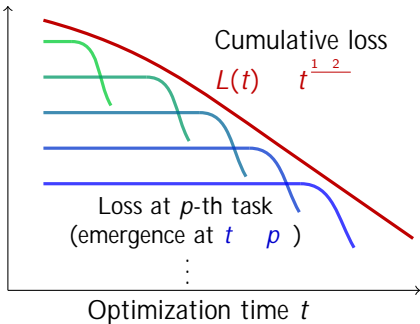
Intuition: Assume *decoupled* learning of different teacher components,

Direction p learned at $T_p \propto p^{-1} d^{k=2}$.

$$L(t) = \sum_{p=1}^P a_p^2 \mathbb{1}(ft < T_p g)$$

$$L(T_p) \propto \int_p^P s^{-2} ds \propto p^{-1/2}.$$

$m < P$ student neurons reaches $s > m$ $a_s^2 \propto m^{-1/2}$ error.



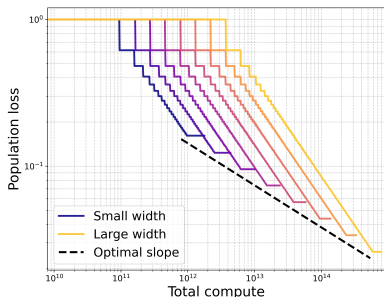
Neural Scaling Laws for Feature Learning

$\rho \propto m^{1/2}$ { approximation barrier.

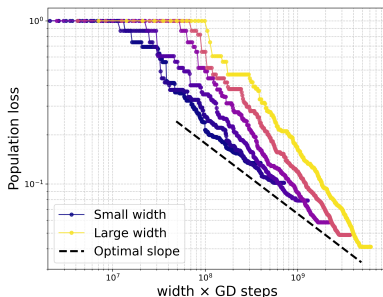
Determined by the student network width.

$\rho \propto t d^{1/k} \propto t^{1/2}$ { optimization error.

Determined by the number of online SGD steps.



(a) Theoretical scaling law.



(b) Empirical scaling law.

Remark on Discretization

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

/ Exponent does not match known rates for weak ρ ball [Johnstone 17].

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

\Optimal" step size for the ρ -th task: $\frac{a_\rho}{d^{k-2}\text{poly}(P)}, a_\rho \propto \rho^{-1}$.

Direction ρ now learned at
 $n = T_\rho \propto \rho^{-1} d^{k-2} = \rho^{-2} d^{k-1} \text{poly}(P)$.

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

$$L(n; m) \approx \frac{n}{d^{k-1} \text{poly}(P)} + m^{1/2}$$

, This now matches the known rates for weak ρ ball.

Proof Sketch

Decoupled gradient dynamics.

Let $w_{(p)}$ denote the neuron that eventually converged to direction p , and $m_{p;q}(t) = \frac{\langle w_p(t), w_q(t) \rangle}{\|w_p(t)\| \|w_q(t)\|}$ measures the overlap at time t .

Claim 1 - decoupling. When all $m_{(p);q}^2$ ($q \neq p$) 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 $m_{(p);q}^2$ ($q \neq p$) remain small when the q -th direction is recovered by $w_{(q)}$ (after which $w_{(q)}$ no longer affects the learning dynamics ("automatic deflation").

Proof Sketch

Decoupled gradient dynamics.

Let $w_{(p)}$ denote the neuron that eventually converged to direction p , and $m_{p;q}(t) = \langle w_p(t), w_q(t) \rangle$ measures the overlap at time t .

Claim 1 - decoupling. When all $m_{(p);q}^2$ ($q \neq p$) 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 $m_{(p);q}^2$ ($q \neq p$) remain small when the q -th direction is recovered by $w_{(q)}$ (after which 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).
 - / 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).
 - / 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

- ⌘ Anisotropic input: $x \sim N(0; \Sigma)$, where $\Sigma = \sum_{i=1}^d v_i v_i^T$; $v_i \sim \dots$.
 - 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

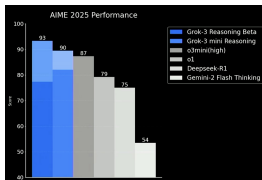
April 3, 2025

Motivation: multi-step reasoning tasks

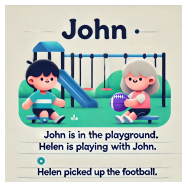
Large language models exhibits remarkable reasoning capabilities.

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

Mathematics



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."

Function composition

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} f_{with} $f_{\text{hold}}(\text{football})$ and get playground.

Note: (location, with, hold) are also function inputs.

Function composition

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.

Abstractly, the model can do the following function composition:

Inputs: (x_1, x_2, \dots, x_k)

Outputs: $f_k^{(k)} \circ f_{k-1}^{(k-1)} \circ \dots \circ f_1^{(1)}(x)$ 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?

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?

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?

Transformers and self-attention

Definition (Self-attention)

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

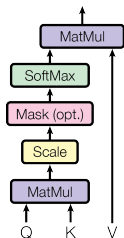
$$\text{attn}(X; W_{KQ}; W_{OV}) = X + W_{OV} X S(X^T W_{KQ} X);$$

where the softmax function S is applied column-wise.

Convex combination of all values $W_{OV} X$ with coefficient $S(X^T W_{KQ} X)_{(i,j)}$, 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.

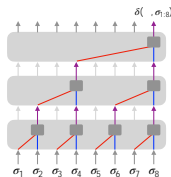
Composing k functions only requires $O(\log k)$ layers.

Examples (Efficiently parallelizable tasks)

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

baebc**ab**ebdea.

(2-hop Induction head)



(Finite state automata)

Our task: k -fold composition

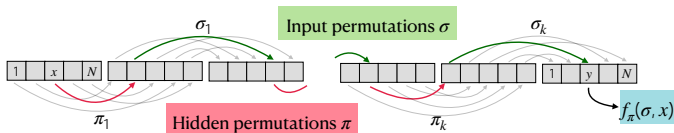
Definition. (k -fold composition)

$\sigma := (\sigma_1; \dots; \sigma_k) \in (S_N)^k$ is a tuple of k hidden permutations.

Input: $(\sigma; x) \in X := (S_N)^k \times [N]$, where $\sigma = (\sigma_1; \dots; \sigma_k)$, x is an index in $[N]$.

Output $f : X \rightarrow [N]$ is defined as

$$f(\sigma; x) := (\sigma_k \circ \sigma_{k-1} \circ \dots \circ \sigma_1)(x):$$



Remark: $\sigma := (\sigma_1; \dots; \sigma_k) \in (S_N)^k$ induces a function class F .

Construction for k -fold composition

Theorem (Expressivity of Transformers)

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

Construction for k -fold composition

Theorem (Expressivity of Transformers)

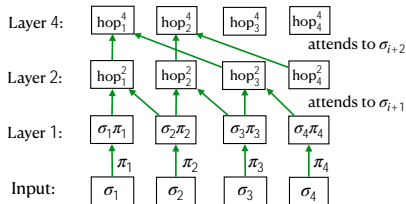
Assume k is a power of 2. For any $\epsilon > 0$, there exists an $L = \log_2 k + 1$ layer transformer with embedding dimension $\Theta(kN)$ that expresses k -fold composition.

Layer 1 encodes hidden σ_i 's.

Layer 2 composes σ_{i+1} and σ_i to get

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

Layer 3 composes $\text{hop}_i^{2^l}$ and $\text{hop}_{i+2^l}^{2^l}$ and get $\text{hop}_i^{2^{l+1}}$.



The parallel solution (Layer 1)

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

Input $X^{(0)}$:

σ_1			
(1,1)	(1,2)		(1,N)
$\sigma_1(1)$	$\sigma_1(2)$		$\sigma_1(N)$
$\mathbf{0}_{kNL}$	$\mathbf{0}_{kNL}$		$\mathbf{0}_{kNL}$

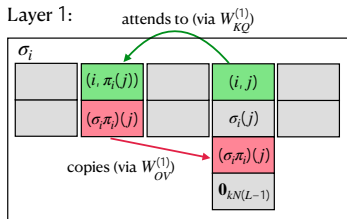
σ_i			
(i,1)	(i,2)		(i,N)
$\sigma_i(1)$	$\sigma_i(2)$		$\sigma_i(N)$
$\mathbf{0}_{kNL}$	$\mathbf{0}_{kNL}$		$\mathbf{0}_{kNL}$

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

The key-query matrix makes position $(i; j)$ attend to position $(i; i(j))$.

The value matrix copies $\text{hop}_i^1(j) := (i; i(j))$ into the residual stream of $(i; j)$

The first layer computes 1-hop!



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

Assume that the previous layer computes $\text{hop}_i^{2^{\ell-1}}$ (i.e residual stream of $(i; j)$ contains $\text{hop}_i^{2^{\ell-1}}(j)$).

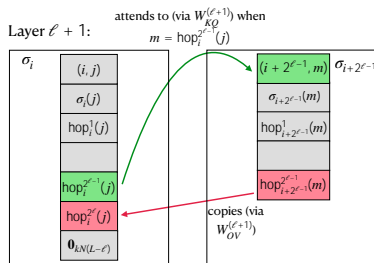
The key-query matrix makes position $(i; j)$ attend to position $(i + 2^{\ell-1}; \text{hop}_{i+2^{\ell-1}}(j))$.

Value matrix copies $\text{hop}_{i+2^{\ell-1}}^{2^{\ell-1}}(\text{hop}_i^{2^{\ell-1}}(j))$ into residual stream of $(i; j)$.

Thus the $(\ell + 1)$ th layer computes

$$\text{hop}_i^{2^{\ell}} = \text{hop}_{i+2^{\ell-1}}^{2^{\ell-1}} \circ \text{hop}_i^{2^{\ell-1}} :$$

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 $q : X \rightarrow [N] \times \mathbb{R}$ and tolerance ϵ . SQ oracle responds with \hat{q}_j , a noisy estimate of q

$$j\hat{q}_j = \mathbb{E}_{x \sim \mathcal{D}} [q(\hat{q}_j; x; f(\hat{q}_j; x))] + \epsilon$$

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 $q : X \rightarrow [N] \times \mathbb{R}$ and tolerance ϵ . SQ oracle responds with \hat{q} , a noisy estimate of q

$$j\hat{q} = \mathbb{E}_{j,x}[q(j;x;f(\cdot;x))] + \epsilon$$

Theorem (SQ lower bound for k -fold composition)

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

- $n^{-1/2}$ by i.i.d. concentration heuristic, where n is the sample size.
) Either runtime (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 $\text{hop}_i^1 = |i - i|$ for $i \in [k]$.

Stage ℓ : GD on the data with label $\text{hop}_i^{2^{\ell-1}}$.

We can learn the k -fold function with $\text{poly}(N; k)$ 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 $\text{hop}_i^1 = i - i$ for $i \in [k]$.

Stage ℓ : GD on the data with label $\text{hop}_i^{\ell-1}$.

We can learn the k -fold function with $\text{poly}(N; k)$ samples!

Theorem (GD upper bound with curriculum)

Assume $k = 2^{L-1}$, if the sample size $n \sim (k^4 N^6)$, 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 $\text{hop}_i^1 = i$ for $i \in [k]$.

Stage ℓ : GD on the data with label $\text{hop}_i^{\ell-1}$.

We can learn the k -fold function with $\text{poly}(N; k)$ samples!

Theorem (GD upper bound with curriculum)

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

Takeaway: $n \sim N^{\Omega(k)}$ (**Exponential**) \rightarrow $n \sim \text{poly}(N; k)$ (**Polynomial**).

Easy examples/intermediate supervision is essential for complex tasks!

Learning with curriculum

Stage 1:

In the construction, the first layer makes the $(i; j)$ position attend to $(i; \pi(j))$ 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 's.

Learning with curriculum

Stage 1:

In the construction, the first layer makes the $(i; j)$ position attend to $(i; \sigma_i(j))$ 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 σ_i 's.

Stage ℓ ($\ell \geq 2$):

ℓ th layer learns to compose $\text{hop}_i^{2^{\ell-2}}$ and $\text{hop}_{i+2^{\ell-2}}^{2^{\ell-2}}$ computed in the previous layer.

With $2^{\ell-1}$ -hop data, GD correctly makes $(i; j)$ position attend to $(i + 2^{\ell-2}; \text{hop}_i^{2^{\ell-2}}(j))$ 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 $n \sim (k^4 N^6)$, the output of gradient descent training on the mixture of f_1, f_2, \dots, f_k -fold data learns the k -fold composition function.

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 $n \sim (k^4 N^6)$, the output of gradient descent training on the mixture of $f_1, 2, 4, \dots, k$ -fold data learns the k -fold composition function.

Intuition: if layer ℓ (the $2^{\ell-1}$ -fold function) is not learned, GD cannot learn later layers with $\ell' > \ell$ | It is essential to learn from easy to hard!

Data mixture **implicitly induces curriculum learning**.

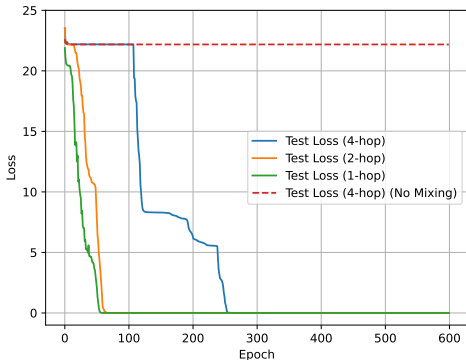
Experiments

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 $n = N^{\Omega(k)}$.

Mirrors k -parity in d dimensions (requires d^k samples/compute).

Remarks

Connection to k -sparse parity

k -fold composition exhibits a *statistical-computational gap*.

Information theoretic sample complexity is $n \sim kN \log N$, while a computationally efficient SQ learner requires $n \sim N^{\Omega(k)}$.

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*.

Remarks

Connection to k -sparse parity

k -fold composition exhibits a *statistical-computational gap*.

Information theoretic sample complexity is $n \sim kN \log N$, while a computationally efficient SQ learner requires $n \sim N^{\Omega(k)}$.

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 Lot , 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.
-  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 tradeo : Limitations of log-precision transformers.
Transactions of the Association for Computational Linguistics, 11:531{545.
-  Sanford, C., Hsu, D., and Telgarsky, M. (2024).