Asymptotics of SGD in Sequence-Single Index Models and Single-Layer Attention Networks

Luca Arnaboldi IdePhics Laboratory

EPFL Lausanne, Switzerland

Bruno Loureiro

Département d'Informatique École Normale Supérieure - PSL Paris, France **Ludovic Stephan**

ENS AI University Rennes Rennes, France

Florent Krzakala

IdePhics Laboratory EPFL Lausanne, Switzerland Lenka Zdeborová

SPOC Laboratory EPFL Lausanne, Switzerland

Abstract

We study the dynamics of stochastic gradient descent (SGD) for a class of sequence models termed Sequence Single-Index (SSI) models, where the target depends on a single direction in input space applied to a sequence of tokens. This setting generalizes classical single-index models to the sequential domain, encompassing simplified one-layer attention architectures. We derive a closed-form expression for the population loss in terms of a pair of sufficient statistics capturing semantic and positional alignment, and characterize the induced high-dimensional SGD dynamics for these coordinates. Our analysis reveals two distinct training phases: escape from uninformative initialization and alignment with the target subspace, and demonstrates how the sequence length and positional encoding influence convergence speed and learning trajectories. These results provide a rigorous and interpretable foundation for understanding how sequential structure in data can be beneficial for learning with attention-based models.

Stochastic Gradient Descent (SGD) is the core optimization tool driving modern machine learning. Recent years have seen substantial progress in understanding its dynamics, particularly in two-layer networks [Saad and Solla, 1995, Mei et al., 2018, Chizat and Bach, 2018, Rotskoff and Vanden-Eijnden, 2022, Sirignano and Spiliopoulos, 2020, Arnaboldi et al., 2023a]. While global convergence is qualitatively well-understood when the network is wide enough, quantitative results are scarcer. A particularly fruitful body of recent theoretical work addressing this gap has focused on deriving precise convergence rates for particular model classes on synthetic data, such as high-dimensional Gaussian single and multi-index models [Ben Arous et al., 2021, Abbe et al., 2022, 2023]. These advances have sparked a wave of follow-up studies [Damian et al., 2022, 2023, Dandi et al., 2024, Bietti et al., 2023, Ba et al., 2023, Moniri et al., 2023, Mousavi-Hosseini et al., 2023, Zweig and Bruna, 2024, Berthier et al., 2024, Arnaboldi et al., 2024a,b], deepening our understanding of what problems are hard to learn for neural networks trained under SGD.

While multi-index models have served as a cornerstone in theoretical analyses of learning, they remain far from the architectures driving recent breakthroughs in machine learning. Modern advances in learning from sequential data — particularly in natural language processing — are increasingly dominated by attention-based models such as the Transformer architecture [Vaswani et al., 2017]. These models introduce a paradigm shift through self-attention mechanisms, which dynamically reweight the influence of each input token based on its relevance to others. Through successive layers of attention, Transformers capture intricate dependencies across sequences, enabling state-of-the-art

performance in tasks ranging from machine translation to large-scale language modeling [Brown et al., 2020, Kenton and Toutanova, 2019].

The main goal of this work is to extend our theoretical understanding of learning with SGD on multi-index models to attention-based architectures and sequential data. Inspired by Cui et al. [2024], Troiani et al. [2025], our focus will be on the following class of single-layer, tied attention model:

$$f_{\boldsymbol{w}}(X) = R \left[\operatorname{softmax} \left(\left(X + \frac{P}{\sqrt{d}} \right) \boldsymbol{w} \boldsymbol{w}^{\top} \left(X + \frac{P}{\sqrt{d}} \right)^{\top} \right) \right].$$
 (1)

where $\boldsymbol{w} \in \mathbb{R}^d$ are the trainable weights and R is the *reduction map*, which allows passing from a $L \times L$ matrix to a k-dimensional vector. As detailed in Appendix A, this model is a reduction of the standard attention mechanism, where (i) key and query matrix are tied, and $d_{\text{head}} = 1$: $Q = K = X \cdot \boldsymbol{w} \in \mathbb{R}^{L \times 1}$; (ii) since we are considering a single layer attention, and we do not need to learn a new representation of the sequence, the value matrix is the identity: $V = I_L$. In this work, we will be interested in the optimization properties of the model in eq. (1) when trained under (spherical) one-pass stochastic gradient descent (SGD) from a random initial condition $\boldsymbol{w}^0 \sim \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d}))$:

$$\boldsymbol{w}^{\tau+1} = \frac{\boldsymbol{w}^{\tau} - \gamma \nabla_{\boldsymbol{w}} \ell(X^{\tau}, \boldsymbol{y}^{\tau}; f_{\boldsymbol{w}})}{\|\boldsymbol{w}^{\tau} - \gamma \nabla_{\boldsymbol{w}} \ell(X^{\tau}, \boldsymbol{y}^{\tau}; f_{\boldsymbol{w}})\|} \|\boldsymbol{w}^{\tau}\|.$$
(2)

with the squared loss:

$$\ell(X, \boldsymbol{y}; f_{\boldsymbol{w}}) = \|\boldsymbol{y} - f_{\boldsymbol{w}}(X)\|_{F}^{2} = \sum_{i=1}^{k} (y_{i} - f_{\boldsymbol{w}}(X)_{i})^{2}.$$
 (3)

Note that for one-pass (a.k.a. *online* or *streaming*) SGD each sample is only seen once, meaning that the sample complexity of the algorithm coincides with the convergence rate. The spherical constraint is considered to simplify the mathematical analysis, a common assumption in the analysis of SGD for single-index models [Ben Arous et al., 2021, Damian et al., 2022].

To derive a sharp characterization of the sample complexity and convergence rate of SGD for the single-layer attention mechanism in eq. (1), we assume that the sequence data (X, y) is generated from the following Gaussian *sequence single-index* (SSI) model:

Assumption 1 (Data distribution). We assume training data $(X, y) \in \mathbb{R}^{L \times d} \times \mathbb{R}^k$ is independently drawn from a Gaussian Sequence Single Index (SSI) model:

$$f_{\boldsymbol{w}_{\star}}^{\mathrm{SSI}}(X) = g(X \cdot \boldsymbol{w}_{\star})$$
 (4)

where $X \in \mathbb{R}^{L \times d}$ is a Gaussian matrix with entries $\mathcal{N}(0, 1/d)$, $\mathbf{g} \colon \mathbb{R}^L \to \mathbb{R}^k$ is a vector-valued link function that depends only on the scalar product of the input with a fixed vector $\mathbf{w}_{\star} \in \mathbb{S}^{d-1}(\sqrt{d})$ and k is a integer that does not depend on d nor L.

Recent work by [Cui et al., 2024, Troiani et al., 2025, Cui, 2025] has shown that the single-layer tied attention model in eq. (1) is a particular instance of a class of *sequence multi-index (SMI) models*, creating a bridge between single-index analysis and attention-based learning. This mapping implies that model eq. (1) can learn at best a predictor in this class, justifying the choice for training data. The SSI model class was introduced by Cui et al. [2024] in the context of studying phase transitions for the attention mechanism. These results position the SMI model as a signature synthetic model to study the interplay between attention mechanisms, data structure, and the dynamics of learning algorithms. Despite this progress, the learning dynamics of SMI models under SGD remain unexplored. Prior studies [Cui, 2025, Cui et al., 2024] analyzed the empirical risk minimizer through the heuristic replica method, while [Troiani et al., 2025] rigorously studied the Bayes-optimal estimator for SMI models. However, a theoretical understanding of the population landscape and SGD dynamics in this model is missing. Our work addresses this gap, providing the first rigorous characterization of SGD dynamics in SMI models by leveraging techniques developed for single and multi-index models.

Main results — Our main methodological contribution is the generalization of analytical tools to study multi-index models to variants that process sequences of tokens rather than simple vector inputs. Our contributions can be summarized as follows:

- We introduce the notion of *sequence information exponent* (SIE), as a generalization of the information exponent for single-index models Ben Arous et al. [2021]; the SIE has a direct correspondence with the sample complexity of SGD. We also discuss the implications of the positional encoding on the sample complexity, proving it could help SGD to learn faster.
- We analyze the speed-up introduced by the attention mechanism when learning sequential data, compared to models not adapted to treat sequence structures, e.g. fully connected networks. For many problems, we show that the gain is proportional to the sequence length L and in some cases even larger.
- We investigate the interplay between the positional and semantic structure of the data following the setting from Cui et al. [2024], showing that SGD dynamics is not always able to disentangle the two and that a rich phase diagram arises describing the structure of the corresponding population loss and the performance of SGD.

All our formal claim are supported by rigorous proofs, as well as numerical experiments; the code developed is available at https://github.com/IdePHICS/Sequence-Single-Index.

Further Related works — There have much activity discussing SGD with synthetic data on multi-index models over the last decades or so, see e.g. [Ben Arous et al., 2021, Veiga et al., 2022, Arnaboldi et al., 2023b, Collins-Woodfin et al., 2024, Marion and Berthier, 2023, Montanari and Urbani, 2025] and reference there in. Information and generative exponent have been the topic of many works over the last few years [Ben Arous et al., 2021, Abbe et al., 2022, Damian et al., 2024, Troiani et al., 2024]. Here we discuss and adapt these notions for sequence models Troiani et al. [2025].

On the topic of the **Theory of SGD in transformers**, Wu et al. [2023] convergence guarantees on SGD for the single layer transformer. Song et al. [2024] also study GD convergence in simple architectures and highlight the existence of suboptimal local solutions. Li et al. [2025] point out that rapid convergence does not guarantee meaningful learning. Li et al. [2023] give sample complexity bounds for a shallow vision transformer. Zhang et al. [2025] being overfitting in SGD trained transformer. Yüksel and Flammarion [2025] focus on gradient-based dynamics for next token prediction tasks. Compared to these work we move beyond convergence results to study how the data structure, e.g. sequence structure and positional encodings, affect sample complexity and behaviours of the the population loss and recovery dynamics in the high-dimensional limit.

Authors of Marion et al. [2024] introduce a model that can be seen as a sequence two-index model. While we assume the input data to be iid Gaussian they assume a spiked covariance structure in the data which makes their results not directly comparable to ours. We anticipate that our results can be generalized to their setting. Another recent work Mousavi-Hosseini et al. [2025] looked at sample complexity separation between attention-based networks and more traditional architectures, while their setting is different, this question is related to ours.

1 Setting and definitions

Let $(X^{\tau}, \boldsymbol{y}^{\tau}) \in \mathbb{R}^{L \times d} \times \mathbb{R}^{k}$ denote $i = 1, \ldots, n$ samples drawn from the Gaussian sequence single-index model with weights $\boldsymbol{w}_{\star} \in \mathbb{R}^{d}$ and link function \boldsymbol{g} , defined in eq. (4). As motivated in the introduction, our goal in this work is to characterize the sample complexity of learning the sequence task $(X^{\tau}, \boldsymbol{y}^{\tau})$ with a tied single-layer attention trained under one-pass (spherical) SGD defined in eq. (2). Note that while the model might appear as too simplified because of the lack of correlation between the tokens, we will show that it is sufficient to capture the main features of sequence models.

As shown in the classical result by Robbins and Monro [1951], one-pass SGD can be understood as a noisy discretization of gradient flow on the *population risk* (often refereed also as *population loss*):

$$R(\boldsymbol{w}) = \mathbb{E}_{X \sim \mathcal{N}(0, I_{\boldsymbol{d}/\boldsymbol{d}})} \left[\ell(X, \boldsymbol{f}_{\boldsymbol{w}_{\star}}^{\mathrm{SSI}}(X); \boldsymbol{f}_{\boldsymbol{w}}) \right].$$
 (5)

Therefore, in order to understand the dynamics in eq. (2) it is important to understand the landscape of the risk above. The key property of single-index models underlying the convergence rate analysis of Ben Arous et al. [2021] is that rotation invariance of the population risk implies that it only depends on a single parameter: the correlation between the target weights and the predictor weight, also

known as an *order parameter* or *sufficient statistic*. A similar property holds for the family of SSI models defined by Assumption 1. Indeed, conditionally on the weights, the following projections

$$z_{\star} = X \cdot w_{\star} \in \mathbb{R}^{L} \quad \text{and} \quad z = \left(X + \frac{P}{\sqrt{d}}\right) w \in \mathbb{R}^{L},$$
 (6)

define joint Gaussian variables, which can be fully characterized by their means and covariances:

$$\mathbb{E}\left[z_{\star,i}\right] = 0, \quad \mathbb{E}\left[z_{i}\right] = e_{i}$$

$$\operatorname{Cov}\left(z_{\star,i}, z_{\star,j}\right) = \delta_{ij}, \quad \operatorname{Cov}\left(z_{i}, z_{j}\right) = \delta_{ij}, \quad \operatorname{Cov}\left(z_{\star,i}, z_{j}\right) = \delta_{ij}m,$$
(7)

where we have introduced the sufficient statistics:

$$m = \frac{\boldsymbol{w}_{\star}^{\top} \boldsymbol{w}}{d}$$
 and $\boldsymbol{e} = \frac{P \boldsymbol{w}}{\sqrt{d}}$. (8)

These play exactly the same role as the overlap between target and predictor weights in the standard single-index model. Therefore, the population risk can be written as a function of these statistics:

$$R(\boldsymbol{w}) \equiv R(\boldsymbol{e}, m) = \mathbb{E}_{(\boldsymbol{z}_{\star}, \boldsymbol{z})} \left[\left\| \boldsymbol{g}(\boldsymbol{z}_{\star}) - \boldsymbol{R} \left[\operatorname{softmax} \left(\boldsymbol{z} \boldsymbol{z}^{\top} \right) \right] \right\|_{F}^{2} \right]. \tag{9}$$

Note that this formulation reduces the problem of understanding the landscape geometry of R in the $\boldsymbol{w} \in \mathbb{R}^d$ space to understanding it in $(\boldsymbol{e},m) \in \mathbb{R}^{L+1}$ — a significant simplification when the token size d is large with respect to the sequence length L, the regime we will focus in this work. Note that, given this direction constraint on \boldsymbol{w} , the sufficient statistics are constrained inside the unit ball of \mathbb{R}^{L+1} : $\|(\boldsymbol{e},m)\| \leq 1$.

Escaping mediocrity — As previously discussed, studying the convergence rate of one-pass SGD is akin to studying the population risk landscape. In the standard single-index model, the picture arising from [Ben Arous et al., 2021, Arnaboldi et al., 2023c] is rather simple: the only critical points of the population risk are a single global minima at the target weights and (possibly) a strict saddle at zero correlation. Therefore, the convergence rate of one-pass SGD from random initialization is dominated by the time taken to escape this saddle-point, a scenario a scenario commonly referred to as *escaping mediocrity* [Arnaboldi et al., 2023c].

As we shall see, the risk landscape of sequence models is richer, with in particular the presence of local minima. Nevertheless, these models share the common property of mediocrity at initialization, with the convergence rate dominated by the flatness of the initial saddle-point. Therefore, we start our discussion by formalizing this notion in the context of SSI models. In the high-dimensional scenario where d is large, the initial weight \mathbf{w}^0 is approximately orthogonal to the target direction \mathbf{w}_{\star} , as well as the positional embedding P. Quantitatively, the sufficient statistics are distributed as

$$\lim_{d \to +\infty} \sqrt{d}(e^0, m^0) \sim \mathcal{N}(0, I_{L+1}). \tag{10}$$

Namely, the initial value of the sufficient statistic is $(e^0, m^0) \approx (\mathbf{0}, 0)$, with fluctuations of order $O(1/\sqrt{d})$. For the model in eq. (1), $(e, m) = (\mathbf{0}, 0)$ is a saddle-point of risk, and the dynamics is divided in two phases:

- The escape from the initial condition, where the model develops a weak correlation with the target direction w_{*} and/or the positional embedding P;
- Full recovery where it reaches a complete overlap with either the target direction w_{\star} and/or the positional embedding $P: ||(e,m)|| \approx 1$.

As previously discussed, the first phase is the one that requires the most number of gradient steps Ben Arous et al. [2021], Arnaboldi et al. [2023c]: the sample complexity required for the first phase is always greater or equal to the one required for the second phase; after having reached a small correlation with the target, the attention decay exponentially fast to a complete overlapped state.

Definition 1 (Weak recovery). Let $\eta \in (0,1)$ a parameter independent from d. We say that the model has weakly recovered the target when $\|(e,m)\| \ge \eta$. The weak recovery time is then

$$\tau_{\eta}^{\text{weak}} = \min \left\{ \tau \geq 0 \colon \left\| (\boldsymbol{e}^{\tau}, m^{\tau}) \right\| \geq \eta \right\}.$$

We use this definition of weak recovery as a proxy for identifying the learning has happened, since the subsequent *strong recovery* will be faster. Figure 1 shows some examples of population loss surface: apart from initialization, there are no other critical points where the dynamic can get slowed down.

Finally, for simplicity of the discussion we will make the following assumption on the spherical one-pass SGD dynamics in eq. (2).

Assumption 2 (Gradient flow approximation). We approximate the training dynamics of eq.(2) via the following ODE, which corresponds to an order-2 Taylor expansion in γ :

$$\frac{d\boldsymbol{w}}{dt} = \mathbb{E}\left[\nabla_{\boldsymbol{w}}^{\perp}\ell(X, y, f_{\boldsymbol{w}})\right] - \frac{\gamma}{2} \,\mathbb{E}\left[\left\|\nabla_{\boldsymbol{w}}^{\perp}\ell(X, y, f_{\boldsymbol{w}})\right\|^{2}\right]\boldsymbol{w},\tag{11}$$

where $\nabla_{\boldsymbol{w}}^{\perp} = (I - \boldsymbol{w} \boldsymbol{w}^{\top}) \nabla_{\boldsymbol{w}}$ is the spherical gradient. The time scaling corresponds to $t = \tau / \gamma$.

As shown in Ben Arous et al. [2021], Arnaboldi et al. [2024c], such an ODE captures both the right weak recovery time for a fixed γ , as well as the maximal value of γ for which the dynamics do not stay trapped in the uninformative region.

2 The sample complexity of SGD

In this section, we focus on understanding the complexity of the SGD algorithm, i.e., how the number of total gradient steps n scales with the dimension d of the token embeddings, in the high-dimensional limit $d \gg 1$; for simplicity of exposition, we focus on the case k=1, but the same arguments can be repeated for each of the components of the output.

No positional encoding — We first focus on the case without positional encoding: P = 0 implies that the only relevant sufficient statistic is m. An analogy with single-index models for networks can be done: For single-index models, the *information exponent* fully characterizes the sample complexity of the SGD algorithm Ben Arous et al. [2021]. The definition can be generalized for sequential data.

Definition 2 (Sequence Information Exponent (SIE)). Given $f_{w_{\star}}^{\mathrm{SSI}}$ a sequence single-index model, let g be the function that acts on the local field z_{\star} , we define the sequence information exponent as

$$\mathit{SIE}(oldsymbol{f}_{oldsymbol{w}_{\star}}^{\mathrm{SSI}}) \coloneqq \min \left\{ \sum_{l=1}^{L} k_l > 0 \colon oldsymbol{k} \in \mathbb{N}^L, \mathbb{E}_{oldsymbol{z} \sim \mathcal{N}(0,I_L)} \left[\left(\prod_{l=1}^{L} \mathrm{He}_{k_l}(z_l)
ight) g(oldsymbol{z})
ight]
eq 0
ight\},$$

where He_k is the k-th order Hermite polynomial.

In Appendix B we provide more details on Hermite polynomials. Let us give some examples of SIE for different SSI models:

- $g(z_{\star}) = z_{\star,1} + z_{\star,2} + \cdots + z_{\star,L}$ has SIE = 1;
- $q(z_{\star}) = z_{\star 1} z_{\star 2} = \text{He}_1(z_{\star 1}) \text{He}_1(z_{\star 2}) \text{ has SIE} = 2;$
- $g(z_{\star}) = \text{He}_1(z_{\star,1}) \text{He}_4(z_{\star,2}) + \text{He}_2(z_{\star,3}) \text{He}_2(z_{\star,4})$ has SIE = 4;
- $g(\mathbf{z}_{\star}) = \prod_{l=1}^{L} \operatorname{He}_{k_{l}}(z_{\star,l})$ has $\operatorname{SIE} = \sum_{l=1}^{L} k_{l}$.

The main feature of the information exponent is that it can be connected to the sample complexity of the SGD algorithm; we prove an equivalent result for the sequence information exponent.

Theorem 1 (Informal). Let $f_{w_{\star}}^{\mathrm{SSI}}(X)$ be a sequence single-index model, and let SIE be its sequence information exponent. If the model f_{w} has a rich enough Hermite expansion, then the sample complexity of the SGD algorithm is

$$t_{\eta}^{+} = \begin{cases} \mathcal{O}_{L}(d) & \text{if SIE} = 1\\ \mathcal{O}_{L}(d\log^{2}d) & \text{if SIE} = 2\\ \mathcal{O}_{L}(d^{SIE-1}) & \text{if SIE} \geq 3 \end{cases}$$

A formal statement of the Theorem and the proof are given in the Appendix C. It relies on the following connection between the *flatness* of the landscape near initialization and the sequence information exponent:

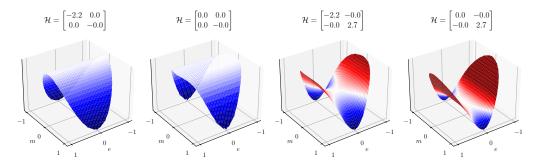


Figure 1: the landscape of the population risk, together with the hessian at initialization, for different values of the SIE and positional encoding. (left) $g(z_{\star}) = \text{He}_2(z_{\star,1}) + \text{He}_2(z_{\star,2})$: SIE=2, no positional encoding: null gradient, but non-null hessian; (center-left) SIE=4, no positional encoding: the first non-null term at initialization is at the 4th order; (center-right) $g(z_{\star}) = \text{He}_4(z_{\star,1}) + \text{He}_4(z_{\star,2})$: SIE=2, with positional encoding: again dynamic dominated by the hessian, but we have a positive curvature in the direction of e; (right) SIE=4, with positional encoding: hessian is positive semidefinite, and the dynamic is again at 4th order in direction of e. In all the examples L=2, $P_1=-P_2$, R=Tr.

Proposition 1. Define the flatness index $\kappa(f_{\boldsymbol{w}}, f_{\boldsymbol{w}_{\star}}^{SSI})$ of the model as

$$\kappa(f_{\boldsymbol{w}}, \boldsymbol{f}_{\boldsymbol{w}_{\star}}^{\mathrm{SSI}}) = \min \left\{ k > 0 : \nabla^{k} R(\boldsymbol{0}, 0) \neq 0 \right\},$$

where R is the reduced population loss defined in (9). Then, if the model f_w has a rich enough Hermite expansion, then

$$\kappa\left(f_{m{w}}, m{f}_{m{w}_{\star}}^{\mathrm{SSI}}\right) = \mathrm{SIE}\left(m{f}_{m{w}_{\star}}^{\mathrm{SSI}}\right)$$

Intuitively, a higher value of κ implies that the landscape of R is flatter around the initialization point $(\mathbf{0},0)$, and thus the number of gradient steps needed to build a weak correlation along either the e or m directions is higher.

In Figure 1, we show some examples of the population loss landscape for different values of the SIE. Figure 1 focuses on even SIE because the symmetry of Equation (1) restricts the possible targets to even functions; in the Appendix D we discuss how to surpass this limitation, and we present settings with odd SIE.

The effect of positional encoding — Despite the fact that positional encoding only acts on the trained model, and not the target function, it changes the population loss, potentially changing the dynamic at initialization. In particular, adding positional encoding increase the expressivity of the model, and can ultimately lead to faster weak-recovery of the SGD.

Lemma 1. Let $f_{\boldsymbol{w}}$ be a model with P=0 that learns a target $\boldsymbol{f}_{\boldsymbol{w}_{\star}}^{\mathrm{SSI}}(X)$ with a given SIE. If we add a positional encoding P to the model, and let $f_{\boldsymbol{w}}^{\mathrm{new}}$ be new model, then

$$\kappa\left(f_{m{w}}^{new}, m{f}_{m{w}_{\star}}^{\mathrm{SSI}}
ight) \leq \kappa\left(f_{m{w}}, m{f}_{m{w}_{\star}}^{\mathrm{SSI}}
ight) = \mathrm{SIE}\left(m{f}_{m{w}_{\star}}^{\mathrm{SSI}}
ight).$$

In other words, the positional encoding can only decrease the flatness of the loss landscape, thus the sample complexity of the SGD algorithm could be reduced. The proof of this lemma is given in the Appendix C. In the right part of Figure 2, we present an example where adding the positional encoding can improve the sample complexity of the SGD algorithm. The left part of Figure 2 shows that the population loss landscape at initialization is flat, and the hessian is null: the first non-null term in the Taylor expansion is at order 4, hence the SIE is 4. The right part of Figure 2 shows instead that when we add

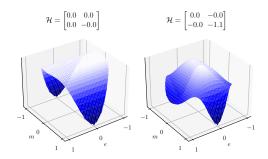


Figure 2: Population loss landscape for P=0 (left) and $P\neq 0$ (right). Example of a case where SIE = 4, while SIE_{positional} = 2. Target: $g(\boldsymbol{z}_{\star})={}^{4}/\!{}_{3}+\mathrm{He}_{4}(z_{1})+2\mathrm{He}_{4}(z_{\star,2}), P_{1}=-P_{2}, R=\mathrm{Tr}.$

the positional encoding, a non-null term appears at order 2, and $\mathrm{SIE}_{\mathrm{positional}} = 2$. Note that while the positions of the global minima are not affected by the positional encoding, SGD can converge to the new local minima instead; more discussion on this point is given in Section 4. In contrast, the right part of Figure 1 shows that the positional encoding is not necessarily beneficial: there are cases for which the loss landscape changes, but not the sample complexity.

3 The role of the sequence length

In this section, we focus a first new emerging characteristic of the *sequence single-index models* over the vanilla *single-index models*: the sequence length L. Therefore, we neglect the effect of the positional encoding, by setting P=0, in order to isolate just the effect of the sequence length. The goal of the section is to measure the speed-up that a model like Equation (1) can achieve over the plain *single-index models* when processing sequential data with length L.

Linear attention — Our first goal is to understand the dependence of the convergence rate of SGD on the sequence length. For that, consider the particular case of *linear attention*, given by the reduction map:

$$R[A] = \boldsymbol{a}_{\text{left}}^{\top} A \boldsymbol{a}_{\text{right}} \quad \text{with } \boldsymbol{a}_{\text{left}} = \boldsymbol{a}_{\text{right}} = \frac{1}{\sqrt{L}} (1, 1, \dots, 1)^{\top} \in \mathbb{R}^{L}.$$
 (12)

Rearranging the terms

$$f_{\boldsymbol{w}}(\boldsymbol{z}) = \boldsymbol{a}_{\text{left}}^{\top} \left(\boldsymbol{z} \boldsymbol{z}^{\top} \right) \boldsymbol{a}_{\text{right}} = \frac{1}{L} \sum_{i=1}^{L} \sum_{j=1}^{L} z_{i} z_{j} = \left(\sum_{i=1}^{L} \frac{z_{i}}{\sqrt{L}} \right)^{2} = \left(\frac{\text{flatten}(X) \cdot \boldsymbol{w}_{\text{tied}}}{\sqrt{L}} \right)^{2}, (13)$$

where $w_{\text{tied}} := \text{concat}(w, \dots, w) \in \mathbb{R}^{Ld}$ is the concatenation of L copies of w. This model is equivalent to a *generalized linear model* with *tied weights*, and activation function $\sigma(x) = x^2$. In terms of performance, taking a general activation σ will at worst be the same of the *attention mechanism* originally considered, if not better. More precisely, we consider:

$$f_{\boldsymbol{w}}(X) = \sigma \left(\frac{\text{flatten}(X) \cdot \boldsymbol{w}_{\text{tied}}}{\sqrt{L}} \right).$$
 (14)

This is the most generic model we study in this section. Further numerical experiments elucidating the equivalence of the speed-up for this *tied network* and for the attention models can be found in Appendix E. Note that *tied networks* are not restricted to learn even function only, differently from the model in Equation (1).

The corresponding untied network — Given the model in Equation (14), a natural benchmark is the model with untied weights. Let $W \in \mathbb{R}^{L \times d}$ be the matrix of weights, whose rows are all updated with Equation (2), the untied network is given by

$$f_W(X) = \sigma\left(\frac{\text{flatten}(X) \cdot \text{flatten}(W)}{\sqrt{L}}\right).$$
 (15)

Since we have L independent weights (the rows of W) the sufficient statistic measuring the overlap between the model and the target SSI is not a scalar as for the tied network, but a vector of length L

$$m = \frac{W w_{\star}}{d} \in \mathbb{R}^{L}$$
 compacted to a scalar as $m_{\text{untied}} = \frac{\|m\|}{\sqrt{L}}$. (16)

Measuring the speedup — The learning rate plays an important role in determining the number of gradient steps needed to reach weak recovery: the larger the learning rate $\gamma(L)$ is, the faster the model learns. However, if it becomes too large, SGD will fail to converge, never achieving weak recovery. The gradient-flow approximation in Eq. (11) exhibits this effect: when γ becomes too large the dynamic of the system is not attracted by w_{\star} or P anymore, and there is no learning. In order to have a faithful measure of the speed-up, we will assume that the learning rate is taken to be the largest possible that guarantees weak recovery; we discuss this upper bound on the learning rate in App. E.

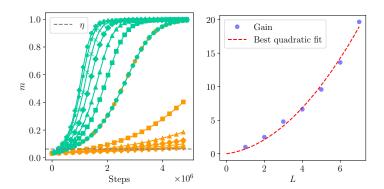


Figure 3: Left: overlap m for tied (green) and untied (orange) networks as a function of the number of gradient steps; different symbols represent different values of L. Right: measured gain as a function of the sequence length L, with the best fit line showing its scaling as L^2 . $g(z_*) = \sum_{i=1}^L \text{He}_2(z_{*,i}), d = 1000, \sigma = \text{ReLU}$.

We measure the speed-up of tied networks with respect to untied networks, as measured in terms of number of gradient steps needed to reach weak recovery, by the ratio of the weak recovery times in the two cases

$$gain(L) := \frac{\tau_{\eta, \text{untied}}^{\text{weak}}}{\tau_{\eta}^{\text{weak}}}.$$
(17)

where $\tau_{\eta, \text{untied}}^{\text{weak}}$ is given by Definition 1 where m is replaced by m_{untied} ; the dependence of gain on η is subleading, and we will neglect it in the following.

Theorem 2. Let $C_{\text{SIE}} \in \mathbb{R}^{L^{\text{SIE}}}$ be the first non-zero tensor in the Hermite expansion of g (see Appendix B). Then the gain satisfies with high probability

$$\operatorname{gain} \gtrsim \left(\frac{C_{\operatorname{SIE}} \times (\mathbf{1}, \dots, \mathbf{1})}{\|C_{\operatorname{SIE}}\|_{\operatorname{op}}}\right)^2 \cdot \begin{cases} L & \textit{if } \operatorname{SIE} = 1\\ 1 & \textit{otherwise} \end{cases}.$$

If the tensor $C_{\rm SIE}$ is orthogonally decomposable, in particular in the cases where ${\rm SIE} \leq 2$ or g is separable, then

$$\operatorname{gain} symp \left(\frac{C_{\operatorname{SIE}} \times (\mathbf{1}, \dots, \mathbf{1})}{\|C_{\operatorname{SIE}}\|_{\operatorname{op}}} \right)^2 \cdot \begin{cases} L & \textit{if } \operatorname{SIE} = 1\\ 1 & \textit{otherwise} \end{cases}.$$

By definition of the operator norm, we have

$$C_{\rm SIE} \times (\mathbf{1}, \dots, \mathbf{1}) \leq \|C_{\rm SIE}\|_{\rm op} \, L^{\rm SIE/2}, \quad \text{and hence} \quad 0 \leq {\rm gain} \lesssim L^{\rm SIE \vee 2}.$$

Since the untied network has L times the number of parameters compared to the tied one, a naive parameter counting argument would yield a gain $=L^{(SIE-1)\vee 1}$ expected gain. Counter-intuitively, the actual gain of using a tied network can either exceed or fall short of this naive value, depending on the function g. In pathological cases (see Appendix E), the tied network can even either fail to learn the target, or do so slower than its untied counterpart.

Example for SIE = 2 — Let's assume to have a target function $g(z_{\star}) = \sum_{i=1}^{L} \text{He}_2(z_{\star,i})$. In this case, the SIE is 2 and the tensor C_2 is simply the identity matrix I_L . The gain is by

$$\operatorname{gain} \asymp \left(\frac{I_L \times (\mathbf{1}, \dots, \mathbf{1})}{\|I_L\|_{\operatorname{op}}}\right)^2 \cdot 1 = \left(\frac{L}{1}\right)^2 \cdot 1 = L^2.$$

Fig. 3 show a numerical experiment, with a ReLU activation, confirming the result of Th. 2.

4 Positional encoding and training dynamics

We now turn our attention to the role played by positional encoding in the attention layer when trained under SGD. Since the focus is on the effect of positional encoding, we stick with a class of target

functions that can exhibit either a *semantic* (label mostly depends on tokens value, but not the order) or a *positional* (where tokens most important feature is their position in the sequence, rather then their embedding). Consider a SSI target function of the form

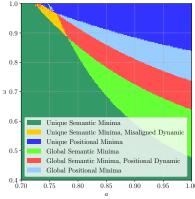
$$f_{\boldsymbol{w}_{\star}}^{\mathrm{SSI}}(X) = (1 - \omega) \operatorname{softmax} \left(X \boldsymbol{w}_{\star} \boldsymbol{w}_{\star}^{\top} X^{\top} \right) + \omega \operatorname{softmax} \left[\begin{pmatrix} a^2 & -a^2 \\ -a^2 & a^2 \end{pmatrix} \right] \in \mathbb{R}^{2 \times 2}, \quad (18)$$

where he parameter $\omega \in [0,1]$ allows the target to switch from a semantic to a positional behavior, while the parameter $a \in (0,1]$ controls the alignment of the target with its positional part.

We shall train the model in Equation (1) with positional encoding $P_1 = -P_2$ and reduction map R the identity function. We focus on the *gradient flow* limit where η is sufficiently small [Robbins and Monro, 1951]. Our analysis will focus on the population loss given by Equation (9), since it completely characterize the behavior of SGD in this regime; the sufficient statistics (e, m) are the only free variables of the setting.

The sequential information exponent of this setting is SIE = 2 (see Appendix F for the explicit derivation), thus the sample complexity for escaping the initialization is $\mathcal{O}(d\log d)$. After the initial phase, SGD fast converges to the minimum of the population loss that is fully aligned with either the semantic or the positional sufficient statistic, i.e. ||(e,m)|| = 1, but is not guaranteed to be the global minimum. Figure 5 shows an example where the population loss has 2 minimums, one semantic with (e,m)=(1,0) and one positional with (e,m)=(0,1), and the steepest direction at initialization, namely the eigenvector associated with the lowest eigenvalue of the Hessian, points towards the local one; in this case, the gradient flow will converge to the wrong minima.

phase diagram with all the possible behaviors:



Varying the parameters, ω and a the high-dimensional SGD dynamics from random initialization exhibits a rich phenomenology. In Figure 4 we show the

- Unique Positional Minima: the population loss has unique positional minima, and the SGD converges to it. This is the case for $\omega = 1$ and a = 1.
- Global Positional Minima: the population loss has both a semantic and a positional minimum, and SGD converges to the global positional one.
- Global Semantic Minima, Positional Dynamic: the population loss has both a semantic and a positional minimum, and SGD converges to the local positional one. This is the case where SGD does not converge to the global minima.
- Unique Semantic Minima, Misaligned Dynamics: the population loss has unique semantic minima, and the SGD converges to it, even though the steepest direction at initialization points orthogonal to it.
- Global Semantic Minima: the population loss has both a semantic and positional minima, and SGD converges to the global semantic one.
- Unique Semantic Minima: the population loss has unique semantic minima, and the SGD converges to it. This is the case for $\omega = 0$ and a = 1.

We verify this by simulating many runs of SGD with different initializations and data samples. In Figure 5 we compute the empirical probability of convergence to the semantic minima, for a=1 and varying ω . The theoretical value of ω where we have a transition from sematic to positional dynamics is $\omega_{\rm trans}=0.64$, which is in good agreement with the transition observed in the measured probabilities; the transition becomes sharper as d increases: ideally, in the limit $d\to\infty$ we expect a step function. The simulations in Figure 5 are performed with d=1000, and some finite size effects are still present. In Appendix F we show present a more detailed analysis, including different values of d.

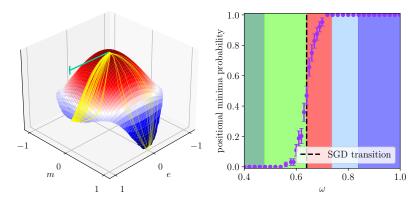


Figure 5: (left) surface of the population loss for $\omega=0.67$ and a=1. The steepest direction at initialization (green vector) points towards the *positional* local minimum, while the global minima is semantic. Some examples of SGD trajectories are shown in yellow: most of them fall into the semantic local minimum, while some others manage to fully-recover the global minimum due to finite size effects (d=1000). (right) empirical probability of convergence to the semantic minima as a function of ω for a=1. The probability is computed over 64 SGD runs with different initializations and data samples. The theoretical prediction of the transition from semantic to positional minima is at $\omega_{\rm trans}\approx 0.64$.

Conclusion

In this paper, we introduced the Sequence-Single Index model as a new, high-dimensional theoretical framework for analyzing single-layer attention architectures. The most significant contribution of this study is the definition of the Sequence Information Exponent. This exponent, which serves as a rigorous tool to quantify the inherent hardness and predict the sample complexity of Stochastic Gradient Descent, is defined in direct analogy with classical single-index models. This analysis transcends qualitative understanding, providing precise scaling laws for the required number of gradient steps, denoted by n, relative to the token dimension, denoted by d. We demonstrated that the sequential setting is significantly richer, showing how the sequence length L accelerates convergence and how positional encoding can proactively reduce the SIE, thereby lowering the computational barrier to learning.

This work establishes a foundational line of research essential for a principled understanding of modern sequential data models, particularly those based on the Transformer architecture. The intricate dynamics manifesting in the population loss landscape, exemplified by the identification of phase transitions and the convergence to suboptimal local minima, unveil a wealth of avenues for future investigation. Subsequent research should aim to fully map these complex high-dimensional dynamics, develop techniques for robustly breaking inherent symmetries, and extend the SIE framework to more complex multi-index and multi-layer attention systems. We hope that this quantitative framework will stimulate further theoretical explorations, leading to the development of a robust, generalizable understanding of learning on structured sequential data.

Limitations – The theoretical analysis of SGD in Sequence Single-Index models is subject to several simplifications for tractability. Specifically, the model under scrutiny is a simplified, single-layer attention architecture where the key and query matrices are tied, the attention head dimension is one $(d_{head} = 1)$, and the value matrix is the identity $(V = I_L)$. It should be noted that the scope of this architecture is restricted to a specific class of target functions, including even functions, although Appendix D provides an extension to this class. Additionally, the analysis operates under the assumption of token independence across the sequence, a crucial simplification that facilitates the execution of numerous sequence modeling tasks. The analysis of SGD dynamics is founded on the approximation of the discrete dynamics via a second-order ODE (Gradient Flow approximation).

Acknowledgement— We would like to thank Luca Pesce, Luca Biggio, and Yatin Dandi for their insightful discussions. We acknowledge funding from the Swiss National Science Foundation grants SNSF SMArtNet (grant number 212049), OperaGOST (grant number 200021 200390), DSGIANGO (grant number 225837) and by the French government, managed by the National Research Agency (ANR), under the France 2030 program with the reference "ANR-23-IACL-0008" and the Choose France - CNRS AI Rising Talents program.

References

- David Saad and Sara A Solla. On-line learning in soft committee machines. *Physical Review E*, 52 (4):4225, 1995.
- Song Mei, Andrea Montanari, and Phan-Minh Nguyen. A mean field view of the landscape of two-layer neural networks. *Proceedings of the National Academy of Sciences*, 115(33):E7665–E7671, 2018.
- Lenaic Chizat and Francis Bach. On the global convergence of gradient descent for over-parameterized models using optimal transport. *Advances in neural information processing systems*, 31, 2018.
- Grant Rotskoff and Eric Vanden-Eijnden. Trainability and accuracy of artificial neural networks: An interacting particle system approach. *Communications on Pure and Applied Mathematics*, 75(9): 1889–1935, 2022. doi: https://doi.org/10.1002/cpa.22074.
- Justin Sirignano and Konstantinos Spiliopoulos. Mean field analysis of neural networks: A central limit theorem. *Stochastic Processes and their Applications*, 130(3):1820–1852, 2020.
- Luca Arnaboldi, Ludovic Stephan, Florent Krzakala, and Bruno Loureiro. From high-dimensional & mean-field dynamics to dimensionless odes: A unifying approach to sgd in two-layers networks. In *The Thirty Sixth Annual Conference on Learning Theory*, pages 1199–1227. PMLR, 2023a.
- Gerard Ben Arous, Reza Gheissari, and Aukosh Jagannath. Online stochastic gradient descent on non-convex losses from high-dimensional inference. *Journal of Machine Learning Research*, 22 (106):1–51, 2021. URL http://jmlr.org/papers/v22/20-1288.html.
- Emmanuel Abbe, Enric Boix-Adsera, and Theodor Misiakiewicz. The merged-staircase property: a necessary and nearly sufficient condition for sgd learning of sparse functions on two-layer neural networks. In *Conference on Learning Theory*, pages 4782–4887. PMLR, 2022.
- Emmanuel Abbe, Enric Boix Adserà, and Theodor Misiakiewicz. Sgd learning on neural networks: leap complexity and saddle-to-saddle dynamics. In Gergely Neu and Lorenzo Rosasco, editors, *Proceedings of Thirty Sixth Conference on Learning Theory*, volume 195 of *Proceedings of Machine Learning Research*, pages 2552–2623. PMLR, 12–15 Jul 2023. URL https://proceedings.mlr.press/v195/abbe23a.html.
- Alexandru Damian, Jason Lee, and Mahdi Soltanolkotabi. Neural networks can learn representations with gradient descent. In Po-Ling Loh and Maxim Raginsky, editors, *Proceedings of Thirty Fifth Conference on Learning Theory*, volume 178 of *Proceedings of Machine Learning Research*, pages 5413–5452. PMLR, 02–05 Jul 2022.
- Alex Damian, Eshaan Nichani, Rong Ge, and Jason D. Lee. Smoothing the landscape boosts the signal for SGD: optimal sample complexity for learning single index models. In Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine, editors, Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 16, 2023, 2023. URL http://papers.nips.cc/paper_files/paper/2023/hash/02763667a5761ff92bb15d8751bcd223-Abstract-Conference.html.
- Yatin Dandi, Florent Krzakala, Bruno Loureiro, Luca Pesce, and Ludovic Stephan. How two-layer neural networks learn, one (giant) step at a time. *Journal of Machine Learning Research*, 25(349): 1–65, 2024. URL http://jmlr.org/papers/v25/23-1543.html.
- Alberto Bietti, Joan Bruna, and Loucas Pillaud-Vivien. On learning gaussian multi-index models with gradient flow. *arXiv preprint arXiv:2310.19793*, 2023.
- Jimmy Ba, Murat A Erdogdu, Taiji Suzuki, Zhichao Wang, and Denny Wu. Learning in the presence of low-dimensional structure: A spiked random matrix perspective. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine, editors, *Advances in Neural Information Processing Systems*, volume 36, pages 17420–17449. Curran Associates, Inc., 2023. URL https://proceedings.neurips.cc/paper_files/paper/2023/file/38a1671ab0747b6ffe4d1c6ef117a3a9-Paper-Conference.pdf.

- Behrad Moniri, Donghwan Lee, Hamed Hassani, and Edgar Dobriban. A theory of non-linear feature learning with one gradient step in two-layer neural networks. *arXiv preprint arXiv:2310.07891*, 2023.
- Alireza Mousavi-Hosseini, Denny Wu, Taiji Suzuki, and Murat A Erdogdu. Gradient-based feature learning under structured data. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine, editors, *Advances in Neural Information Processing Systems*, volume 36, pages 71449–71485. Curran Associates, Inc., 2023. URL https://proceedings.neurips.cc/paper_files/paper/2023/file/e21955c93dede886af1d0d362c756757-Paper-Conference.pdf.
- Aaron Zweig and Joan Bruna. Symmetric single index learning. In *The Twelfth International Conference on Learning Representations*, 2024.
- Raphaël Berthier, Andrea Montanari, and Kangjie Zhou. Learning time-scales in two-layers neural networks. *Foundations of Computational Mathematics*, pages 1–84, 2024.
- Luca Arnaboldi, Florent Krzakala, Bruno Loureiro, and Ludovic Stephan. Escaping mediocrity: how two-layer networks learn hard generalized linear models with sgd. *arXiv preprint arXiv:2305.18502*, 2024a.
- Luca Arnaboldi, Yatin Dandi, Florent Krzakala, Bruno Loureiro, Luca Pesce, and Ludovic Stephan. Online learning and information exponents: The importance of batch size &; Time/Complexity tradeoffs. In Ruslan Salakhutdinov, Zico Kolter, Katherine Heller, Adrian Weller, Nuria Oliver, Jonathan Scarlett, and Felix Berkenkamp, editors, *Proceedings of the 41st International Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning Research*, pages 1730–1762. PMLR, 21–27 Jul 2024b.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in Neural Information Processing Systems*, 31, 2017.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners, 2020.
- Jacob Devlin Ming-Wei Chang Kenton and Lee Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding, 2019.
- Hugo Cui, Freya Behrens, Florent Krzakala, and Lenka Zdeborová. A phase transition between positional and semantic learning in a solvable model of dot-product attention. In A. Globerson, L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang, editors, *Advances in Neural Information Processing Systems*, volume 37, pages 36342–36389. Curran Associates, Inc., 2024. URL https://proceedings.neurips.cc/paper_files/paper/2024/file/3fefebc2d4e3c1c6ee9b892bd293117d-Paper-Conference.pdf.
- Emanuele Troiani, Hugo Cui, Yatin Dandi, Florent Krzakala, and Lenka Zdeborová. Fundamental limits of learning in sequence multi-index models and deep attention networks: High-dimensional asymptotics and sharp thresholds, 2025. URL https://arxiv.org/abs/2502.00901.
- Hugo Cui. High-dimensional learning of narrow neural networks. *Journal of Statistical Mechanics: Theory and Experiment*, 2025(2):023402, 2025.
- Rodrigo Veiga, Ludovic Stephan, Bruno Loureiro, Florent Krzakala, and Lenka Zdeborová. Phase diagram of stochastic gradient descent in high-dimensional two-layer neural networks. Advances in Neural Information Processing Systems, 35:23244–23255, 2022.
- Luca Arnaboldi, Ludovic Stephan, Florent Krzakala, and Bruno Loureiro. From high-dimensional & mean-field dynamics to dimensionless odes: A unifying approach to sgd in two-layers networks. In Gergely Neu and Lorenzo Rosasco, editors, *Proceedings of Thirty Sixth Conference on Learning Theory*, volume 195 of *Proceedings of Machine Learning Research*, pages 1199–1227. PMLR, 12–15 Jul 2023b. URL https://proceedings.mlr.press/v195/arnaboldi23a.html.

- Elizabeth Collins-Woodfin, Courtney Paquette, Elliot Paquette, and Inbar Seroussi. Hitting the high-dimensional notes: an ode for sgd learning dynamics on glms and multi-index models. *Information and Inference: A Journal of the IMA*, 2024. URL https://api.semanticscholar.org/CorpusID:273517012.
- Pierre Marion and Raphaël Berthier. Leveraging the two-timescale regime to demonstrate convergence of neural networks. *Advances in Neural Information Processing Systems*, 36:64996–65029, 2023.
- Andrea Montanari and Pierfrancesco Urbani. Dynamical decoupling of generalization and overfitting in large two-layer networks. *arXiv preprint arXiv:2502.21269*, 2025.
- Alex Damian, Loucas Pillaud-Vivien, Jason D Lee, and Joan Bruna. Computational-statistical gaps in gaussian single-index models. *arXiv preprint arXiv:2403.05529*, 2024.
- Emanuele Troiani, Yatin Dandi, Leonardo Defilippis, Lenka Zdeborová, Bruno Loureiro, and Florent Krzakala. Fundamental limits of weak learnability in high-dimensional multi-index models. *arXiv* preprint arXiv:2405.15480, 2024.
- Yongtao Wu, Fanghui Liu, Grigorios Chrysos, and Volkan Cevher. On the convergence of encoderonly shallow transformers. Advances in Neural Information Processing Systems, 36:52197–52237, 2023.
- Bingqing Song, Boran Han, Shuai Zhang, Jie Ding, and Mingyi Hong. Unraveling the gradient descent dynamics of transformers. *Advances in Neural Information Processing Systems*, 37: 92317–92351, 2024.
- Bingrui Li, Wei Huang, Andi Han, Zhanpeng Zhou, Taiji Suzuki, Jun Zhu, and Jianfei Chen. On the optimization and generalization of two-layer transformers with sign gradient descent. In *The Thirteenth International Conference on Learning Representations*, 2025.
- Hongkang Li, Meng Wang, Sijia Liu, and Pin-Yu Chen. A theoretical understanding of shallow vision transformers: Learning, generalization, and sample complexity. arXiv preprint arXiv:2302.06015, 2023.
- Yingying Zhang, Zhenyu Wu, Jian Li, and Yong Liu. Understanding generalization in transformers: Error bounds and training dynamics under benign and harmful overfitting. *arXiv* preprint *arXiv*:2502.12508, 2025.
- Oğuz Kaan Yüksel and Nicolas Flammarion. On the sample complexity of next-token prediction. In *The 28th International Conference on Artificial Intelligence and Statistics*, 2025.
- Pierre Marion, Raphaël Berthier, Gérard Biau, and Claire Boyer. Attention layers provably solve single-location regression. *arXiv preprint arXiv:2410.01537*, 2024.
- Alireza Mousavi-Hosseini, Clayton Sanford, Denny Wu, and Murat A Erdogdu. When do transformers outperform feedforward and recurrent networks? a statistical perspective. *arXiv preprint arXiv:2503.11272*, 2025.
- Herbert Robbins and Sutton Monro. A Stochastic Approximation Method. *The Annals of Mathematical Statistics*, 22(3):400 407, 1951. doi: 10.1214/aoms/1177729586. URL https://doi.org/10.1214/aoms/1177729586.
- Luca Arnaboldi, Florent Krzakala, Bruno Loureiro, and Ludovic Stephan. Escaping mediocrity: how two-layer networks learn hard generalized linear models. In *OPT 2023: Optimization for Machine Learning*, 2023c.
- Luca Arnaboldi, Yatin Dandi, Florent Krzakala, Luca Pesce, and Ludovic Stephan. Repetita iuvant: Data repetition allows sgd to learn high-dimensional multi-index functions. *arXiv preprint arXiv:2405.15459*, 2024c.
- Elina Robeva. Orthogonal decomposition of symmetric tensors. *SIAM Journal on Matrix Analysis and Applications*, 37(1):86–102, 2016. doi: 10.1137/140989340. URL https://doi.org/10.1137/140989340.
- Harold Grad. Note on N-dimensional hermite polynomials. *Communications on Pure and Applied Mathematics*, 2(4):325–330, 1949. ISSN 1097-0312. doi: 10.1002/cpa.3160020402.

A Reduction from attention to sequence single-index

Let $X \in \mathbb{R}^{L \times d}$ denote a sequence of length L of d-dimensional tokens, and consider the standard dot-product attention:

Attention(X) = softmax
$$\left(\frac{QK^{\top}}{\sqrt{d_{\text{head}}}}\right)V$$
 (19)

where $Q=(X+P)W_Q, K=(X+P)W_K, V=(X+P)W_V\in\mathbb{R}^{L\times d_{\mathrm{head}}}$ are trainable weights known as the *query*, *key* and *value* matrices, respectively. The matrix $P\in\mathbb{R}^{L\times d}$ is the *positional encoding*, a fixed matrix needed to inject a representation of the position of the tokens in the sequence. To make the analysis tractable, [Troiani et al., 2025] considered the following simplifying assumptions:

- Key and query matrix are tied, and $d_{\text{head}} = 1$: $Q = K = X \cdot w \in \mathbb{R}^{L \times 1}$;
- Identity value matrix $V = I_L$.

Note the second assumption is mild for single-layer attention, since we do not need to learn a new representation of the sequence. Under these assumptions, eq. (19) reduces to:

$$TiedAttention(X) = softmax \left(X w w^{\top} X^{\top} \right)$$
(20)

This is a map from sequences $X \in \mathbb{R}^{L \times d}$ to $L \times L$ matrices. Further adding the reduction map $R : \mathbb{R}^{L \times L} \to \mathbb{R}^k$, we get the model in eq. (1). Finally, to get the reduction to a sequence single-index model, it suffices to consider P = 0 and the map on real-valued sequences $s \in \mathbb{R}^L$:

$$g(\mathbf{s}) = R \left[\text{softmax} \left(\mathbf{s} \mathbf{s}^{\top} \right) \right]$$
 (21)

B Mathematical preliminaries and notations

B.1 Tensors

We consider tensors as multidimensional arrays: a tensor T of order k and dimensions (d_1, \ldots, d_k) is simply an element of $\mathbb{R}^{d_1 \times \cdots \times d_k}$. Its elements are denoted by $T_{i_1 \dots i_k}$, where $i_\ell \in [d_\ell]$. The scalar product between two tensors with same dimensions is defined as

$$\langle T, T \rangle = \sum_{i_1, \dots, i_k} T_{i_1 \dots i_k} T'_{i_1 \dots i_k}.$$

We say that a tensor is *symmetric* if all its dimensions are equal and for any index (i_1, \ldots, i_k) and permutation $\sigma \in \mathfrak{S}_k$,

$$T_{i_1...i_k} = T_{i_{\sigma(1)}...i_{\sigma(k)}}.$$

We shall need two operations on tensors: the first is the *tensor product*, that turns two tensors of order k, ℓ into a tensor of order $k + \ell$ defined as

$$(T \otimes T')_{i_1...i_{k+\ell}} = T_{i_1...i_k} T'_{i_{k+1}...i_{k+\ell}}.$$

The second is the *tensor-matrix* contraction: given a tensor T of order k, an index ℓ and a matrix M of size $d_{\ell} \times d'\ell$, the tensor $T \times_{\ell} M$ is defined as

$$(T \times_{\ell} M)_{i_1 \dots i'_{\ell} \dots i_k} = \sum_{i_{\ell}} T_{i_1 \dots i_{\ell} \dots i_k} M_{i_{\ell} i'_{\ell}}$$

Given k matrices $M^{(1)}, \ldots, M^{(k)}$, we will use the shorthand

$$T \times (M^{(1)}, \dots, M^{(k)}) = T \times_1 M^{(1)} \dots \times_k M^{(k)}$$

Immediate properties of those operations are gathered in the following lemma:

Lemma 2. The operation \times is associative: if T is a tensor and $(M^{(1)}, \ldots, M^{(k)}), (N^{(1)}, \ldots, N^{(k)})$ are matrices with compatible dimensions,

$$\left(T \times (M^{(1)}, \dots, M^{(k)})\right) \times (N^{(1)}, \dots, N^{(k)}) = T \times (M^{(1)}N^{(1)}, \dots, M^{(k)}M^{(k)})$$
Let $T \in \mathbb{R}^{d_1 \times \dots \times d_k}$ and $(\boldsymbol{x}_1, \dots, \boldsymbol{x}_k) \in \mathbb{R}^{d_1} \times \dots \times \mathbb{R}^{d_k}$. Then
$$\langle T, \boldsymbol{x}_1 \otimes \dots \otimes \boldsymbol{x}_k \rangle = T \times (\boldsymbol{x}_1, \dots, \boldsymbol{x}_k).$$

Odeco tensors Since tensors of order $k \geq 3$ are sometimes hard to handle, we work with a restricted class. We say that a symmetric tensor T is *odeco* (short for *orthogonally decomposable*, see Robeva [2016]) if there exist real numbers $\lambda_1, \ldots, \lambda_r$ and orthogonal vectors v_1, \ldots, v_r such that

$$T = \sum_{i=1}^{r} \lambda_i v_i^{\otimes k}.$$

In particular, all tensors of order 1 (with r=1) and 2 (with r equal to the rank of T) are odeco.

B.2 Hermite Polynomials

In this section we provide our definition of Hermite polynomials, which are used in the construction of the Hermite basis, both for the one-dimensional and the multidimensional case.

Gaussian measure and Gaussian ℓ^2 space We define the Gaussian density in p dimensions

$$\omega_p(\boldsymbol{x}) = \frac{1}{(2\pi)^{d/2}} \exp\left(-\frac{\|\boldsymbol{x}\|^2}{2}\right),$$

and $d\omega_p(x) = \omega_p(x)dx$. This measure defines a space $\ell^2(\omega_p)$ of functions f satisfying

$$||f||_{\omega} := \int f(\boldsymbol{x})^2 d\omega_p(\boldsymbol{x}) < \infty;$$

it is a Hilbert space w.r.t the scalar product

$$\langle f, g \rangle_{\omega} = \int f(\boldsymbol{x}) g(\boldsymbol{x}) d\omega_p(\boldsymbol{x}).$$

Hermite polynomials and tensors We follow the conventions of Grad [1949]. Define the k-th Hermite tensor \mathcal{H}_k as

$$\mathcal{H}_k = \frac{(-1)^k}{\omega_p} \nabla^k \omega_p,$$

where ∇^k is the k-th order derivative. This results in a k-th order symmetric tensor of size $p \times \cdots \times p$. The Hermite tensors are orthogonal, in the sense that

 $\langle (\mathcal{H}_k)_{i_1...i_k}, (\mathcal{H}_\ell)_{j_1...j_\ell} \rangle_{\omega} \neq 0$ if and only if $k = \ell$ and (i_1, \ldots, i_k) is a permutation of (j_1, \ldots, j_ℓ) .

When p = 1, all Hermite tensors are scalars, and we get the usual Hermite polynomials:

$$He_0(x) = 1, (22)$$

$$He_1(x) = x, (23)$$

$$He_2(x) = x^2 - 1,$$
 (24)

$$He_3(x) = x^3 - 3x, (25)$$

$$He_4(x) = x^4 - 6x^2 + 3. (26)$$

Hermite expansion The orthogonality properties of the Hermite tensors imply the following theorem:

Theorem 3. Let $f \in \ell^2(\omega_p)$. There exist a unique sequence of coefficients $(C_k(f))_{k\geq 0}$ such that $C_k(f)$ is a tensor of order k and

$$f = \sum_{k>0} \langle C_k(f), \mathcal{H}_k \rangle. \tag{27}$$

Those coefficients are given by the following formula:

$$C_k(f) = \frac{1}{k!} \int f(\boldsymbol{x}) \mathcal{H}_k(\boldsymbol{x}) d\omega_p(\boldsymbol{x}).$$

Further, the scalar product $\langle \cdot, \cdot \rangle_{\omega}$ can be written as

$$\langle f, g \rangle_{\omega} = \sum_{k>0} \frac{1}{k!} \langle C_k(f), C_k(g) \rangle$$

The proof of this theorem can be found in Grad [1949]. The identity (27) is called the *Hermite* expansion of f, and the $C_k(f)$ are its *Hermite coefficients*.

Finally, by invariance of the Gaussian distribution through orthogonal transformation, the following holds:

Lemma 3. Let $g: \mathbb{R}^p \to \mathbb{R}$, and $W \in \mathbb{R}^{p \times q}$ be a matrix satisfying $WW^{\top} = I_p$. Let f(x) = g(Wx). Then

$$C_k(f) = C_k(g) \times (W, \dots, W).$$

When p = 1 and w is a single vector, we get

$$C_k(f) = c_k(g) \boldsymbol{w}^{\otimes k}.$$

This gives rise to a link between odeco tensors and separable functions:

Lemma 4. Let $g: \mathbb{R}^{\ell} \to \mathbb{R}$ be a separable function, such that

$$g(\boldsymbol{z}) = \sum_{i} g_i(z_i),$$

 $W \in \mathbb{R}^{\ell imes q}$ an orthogonal matrix, and let $f({m x}) = g(W{m x}).$ Then

$$C_k(f) = \sum_{i=1}^{\ell} c_k(g_i) \boldsymbol{w}_i^{\otimes k},$$

and in particular every Hermite coefficient of f is odeco.

C Formalization and proofs

C.1 Preliminaries

We consider the following approximation of the SGD dynamics:

$$\frac{d\boldsymbol{w}}{dt} = -\mathbb{E}\left[\nabla_{\boldsymbol{w}}^{\perp}\mathcal{L}(X, y, f_{\boldsymbol{w}})\right] - \gamma \,\mathbb{E}\left[\left\|\nabla_{\boldsymbol{w}}^{\perp}\mathcal{L}(X, y, f_{\boldsymbol{w}})\right\|^{2}\right]\boldsymbol{w}$$
(28)

The results of Ben Arous et al. [2021] (when m is a scalar) and Arnaboldi et al. [2024c] (when m is a vector) imply the following:

Theorem 4. Let τ_{η} be the weak recovery time of the ODE (28), with the same initial conditions as the process (2). Then for small enough η and any $\delta > 0$, there exist constants $c(\delta)$, $C(\delta)$ such that if

$$\gamma = c(\delta)(d\tau_n)^{-1},\tag{29}$$

then with probability at least $1 - \delta$

$$t_{\eta}^{+} \leq C(\delta)d\tau_{\eta}^{2}.$$

On the other hand, for any $t \leq C(\delta)d\tau_n^2$, if $\gamma \leq c(\delta)(dt)^{-1/2}$, then with probability $1-\delta$

$$t_n^+ \geq c(\delta)t$$

When γ does not satisfy the bound (29), we cannot show a strong enough concentration around the deterministic ODE dynamics, and hence directly showing non-convergence of (2) is difficult. However, as we shall see in the proof, above this value of γ the inhibitive term in (28) dominates at initialization, and hence the ODE dynamics stay trapped around zero overlap with the target subspace. For this reason, we shall consider (in line with Ben Arous et al. [2021]) that the sample complexity cannot be improved by increasing γ above the bound (29).

Remark. When γ is instead fixed below the value (29), the hitting time of the dynamics 2 is instead given by

$$t_{\eta}^{+}(\gamma) \asymp \gamma^{-1} \tau_{\eta}.$$

We shall also assume that $m_0>0$, and that the coefficients appearing in the gain expression in Theorem 2 are all non-negative. As mentioned in Ben Arous et al. [2021], Arnaboldi et al. [2024b,c], this condition can be ensured with probability 1/2 by randomly setting the learning rate to $\pm \gamma$ with equal probability.

C.2 An ODE for overlap evolution

We begin by computing the expectation of the gradient term:

Lemma 5. In the tied case, we have when $\|\mathbf{w}\| = 1$

$$\mathbb{E}\big[\mathcal{L}(X,y,f_{\textit{tied}})\big] = \mathbb{E}[y^2] - 2\sum_{k>0} c_k(\sigma)C_k(g) \times (\mathbf{1}_L,\ldots,\mathbf{1}_L)m^k + \|\sigma\|_{\omega}.$$

In the untied case, we have instead

$$\mathbb{E}\big[\mathcal{L}(X,y,f_{\textit{untied}})\big] = \mathbb{E}[y^2] - 2\sum_{k \geq 0} c_k(\sigma) C_k(g) \times (\boldsymbol{m},\dots,\boldsymbol{m}) + \|\sigma\|_{\omega} \,.$$

Proof. Recall that $\mathcal{L}(X,y,f)=(y-f(X))^2=y^2-2yf(X)+f(X)^2$. For simplicity, define

$$ilde{oldsymbol{w}}_{ ext{tied}} = (oldsymbol{w} \, oldsymbol{w} \dots oldsymbol{w}) \in \mathbb{R}^{dL}, \quad ilde{oldsymbol{w}}_{ ext{untied}} = (oldsymbol{w}_1 \dots oldsymbol{w}_L) \in \mathbb{R}^{dL} \quad ext{ and } \quad ilde{W}^\star = egin{pmatrix} oldsymbol{w}_1^\star & oldsymbol{0} & \cdots & oldsymbol{0} \\ oldsymbol{0} & oldsymbol{w}_2^\star & \cdots & oldsymbol{0} \\ \vdots & \vdots & \ddots & \vdots \\ oldsymbol{0} & oldsymbol{0} & \cdots & oldsymbol{w}_L^\star \end{pmatrix} \in \mathbb{R}^{L \times dL}$$

Then, for $\square \in \{\text{tied}, \text{untied}\}\$, we have

$$f_{\square}(X) = \sigma\left(\frac{\langle \tilde{\boldsymbol{w}}_{\square}, \operatorname{flatten}(X) \rangle}{\sqrt{L}}\right) \quad \text{ and } \quad y(X) = g(\tilde{W}^{\star} \cdot \operatorname{flatten}(X)).$$

Since $\frac{\mathrm{flatten}(X)}{\sqrt{L}}$ is a standard normal vector, and both $\frac{\tilde{w}_{\square}}{\sqrt{L}}$ and $\frac{\tilde{W}^{\star}}{\sqrt{L}}$ are orthogonal matrices, we can use the Hermite expansion properties to find

$$\mathbb{E}[f_{\square}(X)y(X)] = \sum_{k\geq 0} \langle C_k(f_{\square}), C_k(y) \rangle$$

$$= \sum_{k\geq 0} \langle c_k(\sigma) \tilde{\boldsymbol{w}}_{\square}^{\otimes k}, C_k(g) \times (\tilde{W}^*, \dots, W^*) \rangle$$

$$= \sum_{k\geq 0} c_k(\sigma) C_k(g) \times (\tilde{W}^* \tilde{\boldsymbol{w}}_{\square}, \dots, \tilde{W}^* \tilde{\boldsymbol{w}}_{\square}).$$

In the tied case, we have $\tilde{W}^{\star}\tilde{w}_{\text{tied}} = m\mathbf{1}_{L}$, while in the untied case $\tilde{W}^{\star}\tilde{w}_{\text{untied}} = m$. For the last term, we have in both cases $\frac{\langle \tilde{w}_{\square}, \operatorname{flatten}(X) \rangle}{\sqrt{L}} \sim \mathcal{N}(0,1)$ whenever $\|\tilde{w}_{\square}\|^2 = L$, and hence

$$\mathbb{E}[f_{\square}(X)^2] = \|\sigma\|_{\omega}$$

This ends the proof.

When $\|\boldsymbol{w}\|$ (resp. $\|\boldsymbol{w}_i\|$ is different from one, the expressions of Lemma 5 depend on $q = \|\boldsymbol{w}\|^2$ (resp. $q_i = \|\boldsymbol{w}_i\|^2$). However, we can write

$$\nabla_{\boldsymbol{w}} \mathbb{E}[\mathcal{L}(X, y, f_{\mathsf{tied}})] = \frac{\partial}{\partial m} \mathbb{E}[\mathcal{L}(X, y, f_{\mathsf{tied}})] \boldsymbol{w}^{\star} + 2 \frac{\partial}{\partial g} \mathbb{E}[\mathcal{L}(X, y, f_{\mathsf{tied}})] \boldsymbol{w}.$$

As a result, we have

$$abla_{m{w}}^{\perp} \mathbb{E}[\mathcal{L}(X, y, f_{\mathsf{tied}})] = \left(\frac{\partial}{\partial m} \mathbb{E}[\mathcal{L}(X, y, f_{\mathsf{tied}})] \right) (m{w}^{\star} - mm{w}).$$

The same holds for the untied case:

$$\nabla_{\boldsymbol{w}}^{\perp} \mathbb{E}[\mathcal{L}(X, y, f_{\text{tied}})] = \left(\frac{\partial}{\partial m_i} \mathbb{E}[\mathcal{L}(X, y, f_{\text{tied}})]\right) (\boldsymbol{w}^{\star} - m_i \boldsymbol{w}_i)$$

We arrive at the following result:

Proposition 2. Define the drift functions

$$\phi_{tied}(m) = 2(1 - m^2) \sum_{k \ge 0} kc_k(\sigma) C_k(g) \times (\mathbf{1}_L, \dots, \mathbf{1}_L) m^{k-1}$$

$$\phi_{\mathit{untied}}(m{m}) = 2(1 - m{m} \circ m{m}) \circ \sum_{k \geq 0} c_k(\sigma) C_k(g) \times (I_L, m{m}, \dots, m{m}).$$

Then the tied and untied overlaps satisfy the following ODEs:

$$\frac{dm}{dt} = \phi_{tied}(m) - \gamma \mathbb{E} \left[\left\| \nabla_{\boldsymbol{w}}^{\perp} \mathcal{L}(X, y, f_{\boldsymbol{w}}) \right\|^{2} \right] m$$
(30)

$$\frac{dm}{dt} = \phi_{untied}(\boldsymbol{m}) - \gamma \mathbb{E} \left[\left\| \nabla_{\boldsymbol{w}}^{\perp} \mathcal{L}(X, y, f_{\boldsymbol{w}}) \right\|^{2} \right] \boldsymbol{m}$$
(31)

C.3 Controlling the gradient norm

We turn our attention to the inhibitive terms in Proposition 2. We show the following:

Lemma 6. There exist an $\eta > 0$ and two constants c, C such that if $m < \eta$ (resp. $||m|| \le \eta$),

$$cd \leq \mathbb{E} \bigg[\bigg\| \nabla_{\boldsymbol{w}}^{\perp} \mathcal{L}(X, y, f_{\boldsymbol{w}}) \bigg\|^2 \bigg] \leq Cd$$

Proof. We only treat the tied case; the untied one is done similarly. Differentiating the loss w.r.t w, we find

$$\nabla_{\boldsymbol{w}} \mathcal{L}(X, y, f) = 2\sigma' \left(\frac{\langle \boldsymbol{w}_{\text{tied}}, \text{flatten}(X) \rangle}{\sqrt{L}} \right) \left(\sigma' \left(\frac{\langle \boldsymbol{w}_{\text{tied}}, \text{flatten}(X) \rangle}{\sqrt{L}} \right) - y(X) \right) \cdot \frac{\sum_{i} \boldsymbol{x}_{i}}{\sqrt{L}}$$

We write $x_i = x_i^{\parallel} + x_i^{\perp}$, where x_i^{\parallel} is the projection on x_i on the subspace spanned by w and w^* . Then

$$\nabla_{\boldsymbol{w}}^{\perp} \mathcal{L}(X, y, f) = 2\sigma' \left(\frac{\langle \boldsymbol{w}_{\text{tied}}, \text{flatten}(X) \rangle}{\sqrt{L}} \right) \left(\sigma' \left(\frac{\langle \boldsymbol{w}_{\text{tied}}, \text{flatten}(X) \rangle}{\sqrt{L}} \right) - y(X) \right) \cdot \frac{\sum_{i} (I - \boldsymbol{w} \boldsymbol{w}^{\perp}) \boldsymbol{x}_{i}^{\parallel} + \sum_{i} \boldsymbol{x}_{i}^{\perp}}{\sqrt{L}}$$

Importantly, the vector x_i^{\perp} is independent from any of the prefactors, and has norm d-2, while the first vector is a fixed-dimensional Gaussian. As a result, we have

$$\mathbb{E}\left[\left\|\nabla_{\boldsymbol{w}}^{\perp}\mathcal{L}(X,y,f)\right\|^{2}\right] = \mathbb{E}\left[f_{\boldsymbol{w}}(X)^{2}(y(X) - f_{\boldsymbol{w}}(X))^{2}\right](d-2) + O(1)$$

It remains to show that the expectation above is bounded away from zero. Assuming that the labels are centered for simplicity, when m=0 the expectation simplifies to

$$L_0 = \mathbb{E}\left[f_{\boldsymbol{w}}(X)^2\right] \mathbb{E}\left[y(X)^2\right] + \mathbb{E}\left[f_{\boldsymbol{w}}(X)^4\right] > 0.$$

By continuity, we can choose $\eta>0$ such that if $|m|\leq\eta$

$$\frac{L_0}{2} \le \mathbb{E}\Big[f_{\boldsymbol{w}}(X)^2(y(X) - f_{\boldsymbol{w}}(X))^2\Big] \ge 2L_0,$$

which ends the proof.

C.4 Hitting time for the tied dynamics

We are now ready to prove Theorem 1 (as well as part of Theorem 2). In light of Theorem 4, it suffices to compute the hitting time τ_{η} of the tied dynamics (30). Since $\phi_{\text{tied}}(m) = 2C_{\text{SIE}} \times (\mathbf{1}_L, \ldots, \mathbf{1}_L) m^{\text{SIE}-1} + O(m^{\text{SIE}})$, there exists an $\eta_{\text{tied}} > 0$ such that if $m < \eta_{\text{tied}}$,

$$c_{\mathrm{SIE}}(\sigma)C_{\mathrm{SIE}}\times(\mathbf{1}_L,\ldots,\mathbf{1}_L)m^{\mathrm{SIE}-1}<\phi_{\mathrm{tied}}(m)<3c_{\mathrm{SIE}}(\sigma)C_{\mathrm{SIE}}\times(\mathbf{1}_L,\ldots,\mathbf{1}_L)m^{\mathrm{SIE}-1}$$

Combining this with the bound of Lemma 6, we find that whenever $m(t) \leq \eta$ for some constant c > 0,

$$\frac{dm}{dt} \ge c \cdot C_{\text{SIE}} \times (\mathbf{1}_L, \dots, \mathbf{1}_L) m^{\text{SIE}-1} - C\gamma dm$$

For any $\delta > 0$, there exists a constant $\kappa(\delta)$ such that with probability at least $1 - \delta$,

$$m(0) \ge \frac{\kappa(\delta)}{\sqrt{d}}.$$

As a result, when $\gamma \le \kappa(\delta)^{{
m SIE}-1}C^{-1}d^{-\frac{{
m SIE}}{2}+1}$, then for $0\le t\le au_\eta$

$$\frac{dm}{dt} \ge c' \cdot C_{\text{SIE}} \times (\mathbf{1}_L, \dots, \mathbf{1}_L) m^{\text{SIE}-1}$$

This implies:

• when SIE = 1,

$$m(t) \ge m_0 + \frac{C_{\text{SIE}} \times (\mathbf{1}_L, \dots, \mathbf{1}_L)}{2} t$$
, hence $\tau_{\eta} \le \frac{C' \eta}{C_{\text{SIE}} \times (\mathbf{1}_L, \dots, \mathbf{1}_L)}$;

• when SIE = 2,

$$m(t) \geq m_0 \exp\left(\frac{C_{\mathrm{SIE}} \times (\mathbf{1}_L, \dots, \mathbf{1}_L)}{2} t\right), \quad \text{hence} \quad \tau_\eta \leq \frac{C' \ln(d)}{C_{\mathrm{SIE}} \times (\mathbf{1}_L, \dots, \mathbf{1}_L)};$$

• when SIE > 2,

$$m(t) \geq \left(m(0)^{2-\mathrm{SIE}} - \frac{C_{\mathrm{SIE}} \times (\mathbf{1}_L, \dots, \mathbf{1}_L)}{2} t^{\mathrm{SIE}-2}\right)^{-\frac{1}{\mathrm{SIE}-2}}, \quad \text{hence} \quad \tau_{\eta} \leq \frac{C' d^{\frac{\mathrm{SIE}}{2}-1}}{C_{\mathrm{SIE}} \times (\mathbf{1}_L, \dots, \mathbf{1}_L)}.$$

The bound on γ above is always weaker (up to constant factors) than the bound $\gamma \leq C(d\tau_d)^{-1}$ of Theorem 4, and hence we always have

$$t_n^+ \le C d\tau_n^2$$

which corresponds to the statement of Theorem 1.

On the other hand, for any $t \ge 0$ we have

$$\frac{dm}{dt} \le 3C_{\text{SIE}} \times (\mathbf{1}_L, \dots, \mathbf{1}_L) m^{\text{SIE}-1},$$

which by the same reasoning implies that the bounds on τ_{η} that we obtained are sharp up to constants.

C.5 Hitting time for untied dynamics

We are now ready to finish the proof of Theorem 2. Since

$$\phi_{\text{untied}}(\boldsymbol{m}) = 2c_{\text{SIE}}(\sigma)C_{\text{SIE}} \times (I_L, \boldsymbol{m}, \dots, \boldsymbol{m}) + O(\|\boldsymbol{m}\|^{\text{SIE}}),$$

for small enough $\eta > 0$ we have for $t \leq \tau_n$

$$\frac{d\|\boldsymbol{m}\|}{dt} \leq C \cdot C_{\text{SIE}} \times (\boldsymbol{m}, \dots, \boldsymbol{m}) + \|C_{\text{SIE}}\| \|\boldsymbol{m}\|^{\text{SIE}-1} \leq (C+1) \|C_{\text{SIE}}\| \cdot \|\boldsymbol{m}\|^{\text{SIE}-1}.$$

Recall that the hitting time $\tau_{\eta,\text{untied}}$ corresponds to ||m|| hitting the threshold $\eta\sqrt{L}$. As a result, since the dependency in η is of leading order only for SIE=1, we find

$$\tau_{\eta, \text{untied}} \geq \frac{c}{\|C_{\text{SIE}}\|} \begin{cases} \eta \sqrt{L} & \text{if SIE} = 1\\ \log(d) & \text{if SIE} = 2\\ d^{\frac{\text{SIE}}{2} - 1} & \text{if SIE} \geq 3 \end{cases}$$

Since the gain satisfies

$$\operatorname{gain} symp \left(rac{ au_{\eta, ext{untied}}}{ au_{\eta}}
ight)^2,$$

this proves the first part of Theorem 2.

For the second part, assume that C_{SIE} is *odeco*, hence there exists $\lambda_1, \ldots, \lambda_L$ and orthogonal vectors v_1, \ldots, v_L such that

$$C_{ ext{SIE}} = \sum_{i=1}^L \lambda_i oldsymbol{v}_i^{\otimes ext{SIE}}.$$

We assume that the λ_i are ordered by absolute value, so that $||C_{\text{SIE}}|| = |\lambda_1|$ Letting $m^{(1)} = \langle \boldsymbol{m}, \boldsymbol{v_1} \rangle$, we have

$$\langle \boldsymbol{v}_1, C_{\text{SIE}} \times (I_L, \boldsymbol{m}, \dots, \boldsymbol{m}) \rangle = C_{\text{SIE}} \times (\boldsymbol{v}_1, \boldsymbol{m}, \dots, \boldsymbol{m}) = \lambda_1(m^{(1)})^{\text{SIE}-1}$$

For small enough η , this implies that

$$\frac{dm^{(1)}}{dt} \ge c \|C_{\text{SIE}}\| (m^{(1)})^{\text{SIE}-1}$$

Since $||m|| \ge m^{(1)}$ by the Cauchy-Schwarz inequality, the hitting time $\tau_{\eta, \text{untied}}$ is at most that of $m^{(1)}$, and hence

$$\tau_{\eta, \text{untied}} \leq \frac{C}{\|C_{\text{SIE}}\|} \begin{cases} \eta \sqrt{L} & \text{if SIE} = 1\\ \log(d) & \text{if SIE} = 2\\ d^{\frac{\text{SIE}}{2} - 1} & \text{if SIE} \geq 3 \end{cases}$$

This closes the upper bound of Theorem 2.

C.6 Proof of 1

Let $R^{\text{new}}(e, m)$, R(m) be the reduced population losses for the model with and without positional encoding, respectively, so that $R(m) = R^{\text{new}}(\mathbf{0}, m)$. Then

$$\frac{\partial^k R^{\text{new}}}{\partial m^k}(\mathbf{0}, 0) = \frac{d^k R}{dm^k}(0),$$

and hence $\nabla^k R(0) \neq 0$ implies that $\nabla^k R^{\text{new}} \neq 0$.

D Sequence Information Exponent Beyond attention

In the main paper, we discussed the learning of a generic Sequence Single Index model with a parametrized model like Equation (1), with the goal of modelling the attention mechanism learning. The theory of Sequence Information Exponent goes beyond this, and can be used to understand the sample complexity of the learning of a generic sequence single-index model both as a target and as a trained model. In this Appendix, we focus on a particular choice of positional encoding that breaks the even symmetry of the model, allowing attention to weakly recover odd SIE targets. In particular, we consider a dynamical positional encoding of the form:

$$P_i = \frac{c_i}{\sqrt{d}} \boldsymbol{w} + \tilde{P}_i \quad \text{with } \tilde{P}_i \in \mathbb{R}^{L \times d} \text{ a fixed vector.}$$
 (32)

 $c \in \mathbb{R}^L$ is a fixed vector of coefficients. We call this special version of positional encoding *injected* positional encoding. The trained model becomes:

$$f_{\boldsymbol{w}}(X) = R \left[\operatorname{softmax} \left(\left(X + \frac{\tilde{P}}{\sqrt{d}} \right) \boldsymbol{w} \boldsymbol{w}^{\top} \left(X + \frac{\tilde{P}}{\sqrt{d}} \right)^{\top} + \boldsymbol{c} \boldsymbol{c}^{\top} \right) \right],$$
 (33)

and it has now a non-zero odd Hermite expansion. In Figure 6 we show examples of population losses with odd targets, learned with the model in Equation (33). The injected positional encoding breaks the symmetry, as the population risk plots show.

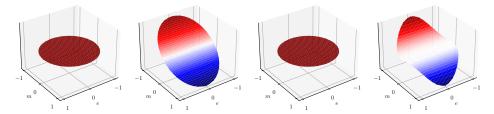


Figure 6: Population loss of the model in Equation (33) with odd targets. (left) $g(\boldsymbol{z}_{\star}) = \operatorname{He}_{1}(\boldsymbol{z}_{\star,1}) + \operatorname{He}_{1}(\boldsymbol{z}_{\star,2})$, no positional encoding; (center-left) $g(\boldsymbol{z}_{\star}) = \operatorname{He}_{1}(\boldsymbol{z}_{\star,1}) + \operatorname{He}_{1}(\boldsymbol{z}_{\star,2})$, injected positional encoding; (center-right) $g(\boldsymbol{z}_{\star}) = \operatorname{He}_{3}(\boldsymbol{z}_{\star,1}) + \operatorname{He}_{3}(\boldsymbol{z}_{\star,2})$, no positional encoding; (right) $g(\boldsymbol{z}_{\star}) = \operatorname{He}_{3}(\boldsymbol{z}_{\star,1}) + \operatorname{He}_{3}(\boldsymbol{z}_{\star,2})$, injected positional encoding. L = 2, $\tilde{P} = 0$ in order to isolate the effect of the injection.

E Further analysis on the effect of sequence length

The aim of this Appendix is to clarify and give further examples on Theorem 2. Here we provide an intuitive explanation of the result, while the mathematical details can be found in Appendix C.

Let's start from the main result on the gain, that can be broken down in three parts:

gain
$$\sim$$
 (gain at costant γ) $\cdot \left(\frac{\gamma_{\text{tied}}}{\gamma_{\text{untied}}}\right) \cdot (\text{special factor for SIE} = 1)$ (34)

The gain constant γ This speedup comes from the different structure of the two networks, the tied one can built up correlation faster than the tied one because it compose L different signals. The strength of this effect is strongly dependent on the target function, the overall result is

gain at constant
$$\gamma \sim \frac{C_{\text{SIE}} \times (\mathbf{1}, \dots, \mathbf{1})}{\|C_{\text{SIE}}\|_{\text{op}}}.$$
 (35)

The ratio of the learning rates In Appendix \mathbb{C} , we showed that the learning rate can grow with the sequence length L at most as

for the tied network

$$\gamma_{\rm tied} \lesssim C_{\rm SIE} \times (\mathbf{1}, \dots, \mathbf{1}),$$

• for the untied network

$$\gamma_{\text{untied}} \lesssim ||C_{\text{SIE}}||_{\text{op}}$$
.

It is clear that saturating the bounds above leads to a factor

$$\frac{\gamma_{\text{tied}}}{\gamma_{\text{untied}}} \lesssim \frac{\max \gamma_{\text{tied}}}{\max \gamma_{\text{untied}}} \sim \frac{C_{\text{SIE}} \times (\mathbf{1}, \dots, \mathbf{1})}{\|C_{\text{SIE}}\|_{\text{op}}},$$
 (36)

that is exactly the same as the gain at constant γ . Obviously the gain measure is fair only if the learning rates bounds are saturated, but our result predict the speed-up factor even in the case where the optimal learning rates are not used.

The special factor for SIE = 1 The case SIE = 1 is special because the dependence of the weak recovery time on the constant η is not negligible, differently from the cases $SIE \geq 2$. This slows down further the learning of the untied network, by a factor \sqrt{L} , leading to a factor

special factor for SIE =
$$\begin{cases} 1 & \text{SIE} \ge 2\\ \sqrt{L} & \text{SIE} = 1 \end{cases}$$
 (37)

E.1 Example at not optimal learning rate

In Figure 3 we presented the speed up in the case of a target function with SIE = 2, and optimal learning rate for both the tied and untied network. Here we want to show that our result predicts the

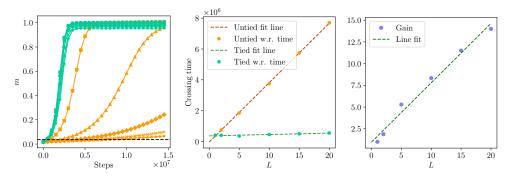


Figure 7: The gain in the case of a target function with SIE = 2, where both cases use the same learning rate $\gamma_0 = 0.005$. (left) The evolution of the overlap, (middle) the weak recovery time and (right) the gain. The gain is proportional to L, as predicted.

gain well even when the learning rate are not scaled up optimally. In particular, consider the same setting as in Figure 3, but with the learning rates

$$\gamma_{\rm tied} = \gamma_{\rm untied} = {\rm costant \ with } L = \gamma_0.$$
 (38)

Using Equation (34) we can predict the gain as

$$gain \sim L \cdot 1 \cdot 1 = L. \tag{39}$$

In Figure 7 we show the result of the simulation, where we can see that the gain is indeed proportional to L, as predicted.

E.2 The upper bound of learning rate scaling

In this subsection we want to show an example proving that not all scalings of the learning rate are allowed: if it grows too fast with the sequence length, the network will not be able to learn.

Let's take as example the SIE = 1 target function

$$g(\boldsymbol{z}_{\star}) = \frac{1}{\sqrt{L}} \sum_{i=1}^{L} z_{\star,i},\tag{40}$$

with the corresponding leading Hermite tensor

$$C_1 = \begin{pmatrix} 1/\sqrt{L} & \dots & 1/\sqrt{L} \end{pmatrix} \in \mathbb{R}^L. \tag{41}$$

We know that the maximum scaling for the learning rate of the untied network is

$$\gamma_{\text{untied}} \lesssim ||C_1||_{\text{op}} = 1,$$

thus we stick with a constant learning rate γ_0 for both networks

$$\gamma_{\rm tied} = \gamma_{\rm untied} = \gamma_0 \sim 1.$$
 (42)

We can use Equation (34) to have the theoretical prediction of the gain in this case:

$$C_1 \times (1, \dots, 1) = \sqrt{L} \implies \text{gain } \sim \sqrt{L} \cdot 1 \cdot \sqrt{L} = L.$$
 (43)

In Figure 8 we show the result of the simulation, where we can see that the gain is indeed proportional to L, as predicted. We can now ask what happens if we push the learning rate scaling of the untied network beyond the limit of Theorem 2. If we set

$$\gamma_{\text{untied}} = \gamma_{\text{tied}} = \gamma_0 \cdot L,$$
 (44)

the ratio between the learning rates becomes now $\gamma_{\text{tied}}/\gamma_{\text{untied}} = 1/L$, and the gain

$$gain \sim \sqrt{L} \cdot \frac{1}{L} \cdot \sqrt{L} = 1. \tag{45}$$

Apparently, the untied network performance is matching the one of the tied network. Figure 9 shows the result of the simulation of such a case: the learning rate of the untied network is too high, and the network is not able for large L, and thus the gain diverges with a maximum value of . This plot proves our bounds on the learning rate scaling are indeed correct.

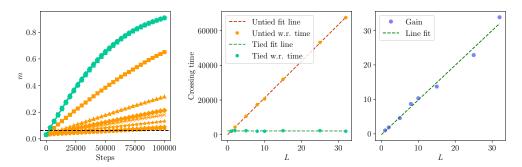


Figure 8: The gain in the case of a target function with SIE = 1, where both cases use the same learning rate $\gamma_0 = 0.005$. (left) The evolution of the overlap, (middle) the weak recovery time. The gain is proportional to L, as predicted. d = 1000, $\sigma = \mathrm{ReLU}$

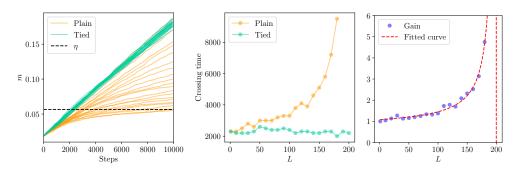


Figure 9: The gain in the case of a target function with SIE = 1, where both cases use the same learning rate $\gamma_0 = 0.005$. (left) The evolution of the overlap, (middle) the weak recovery time and (right) the gain. The gain diverges, as predicted. d = 1000, $\sigma = \mathrm{ReLU}$

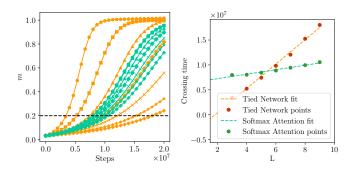


Figure 10: comparison between the learning of the tied network with squared activation (orange), namely linear attention, and the single-layer softmax attention(green) (left) The evolution of the overlap; each symbol is a different value of L. (right) the weak recovery time. The two networks learn at the same rate, but with different costants. $g(z_{\star}) = \sum_{i=1}^{L} \operatorname{He}_2(z_{\star,i}), d = 1000, \sigma = \operatorname{ReLU}$

E.3 Equivalence between Single-Layer attention and Tied network

In the main paper we proved that the tied network model is a generalization of the single-layer attention model. In particular, if the activation function is $\sigma(x)=x^2$, the tied network is equivalent to the single-layer *linear* attention. We also claimed that adding a non-linearity to the attention can only improve the performance of the network, although not affecting the scaling with the sequence length. In Figure 10 we show the comparison between the learning of the tied network with squared activation (orange), namely linear attention, and the single-layer softmax attention(green). The two networks learn at the same linear rate, but with different growth constants. Softmax attention performs better than the tied network with squared activation for sufficiently large sequence lengths.

E.4 Pathological cases

Theorem 2 shows that there could be cases where the gain from learning is 0, meaning that the tied network cannot learn anything, while the untied one possibly can. In particular the degenerate condition happens when

$$C_{\text{SIE}} \times (\mathbf{1}, \dots, \mathbf{1}) = 0, \tag{46}$$

where Theorem 2 guarantees that the gain is 0. Examples of such pathological cases are:

• SIE = 1: a possible pathological target could be

$$g(\boldsymbol{z}_{\star}) = z_{\star,1} - z_{\star,2}.$$

In this case the leading term in Hermite expansion is

$$C_1 = \begin{pmatrix} 1 \\ -1 \end{pmatrix}$$
 and consequently $C_1 \times \mathbf{1} = 0$.

• SIE = 2: anlogously, a possible pathological target could be

$$g(\mathbf{z}_{\star}) = z_{\star,1}^2 - z_{\star,2}^2.$$

In this case the leading term in Hermite expansion is

$$C_2 = \begin{pmatrix} 1 & 0 \\ 0 & -1 \end{pmatrix}$$
 and consequently $C_2 \times \begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix} = 0$.

In Figure 11 we show the pathological case for SIE = 1 and SIE = 2. The untied network is not able to learn the target function because of the symmetry: the tied network is by design symmetric, while the target is constructed to be as antisymmetric as possible in the sequence length. The untied network instead is able to learn the target function, because the weights are not constrained to be equal, thus the symmetry is broken by the initialization.

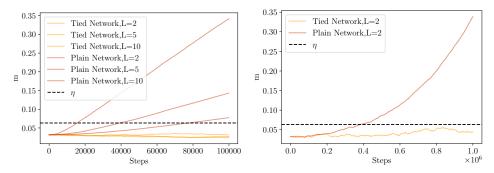


Figure 11: Pathological cases for ${\rm SIE}=1$ (left) and ${\rm SIE}=2$ (right). The untied network is able to learn the target function, while the tied one cannot.

The degeneracy of the tied network could be solved by weighting randomly the neurons. In that case, almost surely we have

$$C_{\text{SIE}} \times (\mathbf{1}, \dots, \mathbf{1}) \neq 0. \tag{47}$$

Although this is a possible solution, we leave the study of the effect of random weights fro symmetry breaking for future work.

F Further discussion on the positional-semantic transition

F.1 SIE of model (18)

Definition 2 of the sequence information exponent does not apply directly to the model (18) because the model has a non-scalar output. However, we can still extend the concept of SIE to this case by looking at the update rule of sufficient statistics around initialization.

The population loss function is given by

$$R(e,m) = \mathbb{E}_{\boldsymbol{z},\boldsymbol{z}_{\star} \sim \mathcal{P}(e,m)} \left[\left\| (1-\omega)\operatorname{softmax} \left(\boldsymbol{z}_{\star} \boldsymbol{z}_{\star}^{\top} \right) + \omega \operatorname{softmax} \left[\begin{pmatrix} a^{2} & -a^{2} \\ -a^{2} & a^{2} \end{pmatrix} \right] - \operatorname{softmax} \left(\boldsymbol{z} \boldsymbol{z}^{\top} \right) \right\|_{F} \right]$$

$$= \mathbb{E}_{\boldsymbol{z},\boldsymbol{z}_{\star} \sim \mathcal{P}(e,m)} \left[\mathcal{L}(\boldsymbol{z},\boldsymbol{z}_{\star}) \right]$$

$$(48)$$

with

$$\mathcal{P}(e,m) \equiv \mathcal{N} \left(\begin{pmatrix} 0 \\ 0 \\ e \\ -e \end{pmatrix}, \begin{pmatrix} 1 & 0 & m & 0 \\ 0 & 1 & 0 & m \\ m & 0 & 1 & 0 \\ 0 & m & 0 & 1 \end{pmatrix} \right). \tag{49}$$

Using parity arguments, we can easily show that the gradient of the population loss is 0 at the initialization point (e, m) = (0, 0). Let $p(z_{\star}, z; e, m)$ be the probability density function associated with the distribution $\mathcal{P}(e, m)$. Then the gradient at initialization is

$$\nabla_{(e,m)} R(e,m) \Big|_{e=m=0} = \mathbb{E}_{\boldsymbol{z},\boldsymbol{z}_{\star} \sim \mathcal{P}(e,m)} \left[\frac{\nabla_{(e,m)} p(\boldsymbol{z}_{\star}, \boldsymbol{z}; e, m)}{p(\boldsymbol{z}_{\star}, \boldsymbol{z}; 0, 0)} \mathcal{L}(\boldsymbol{z}, \boldsymbol{z}_{\star}) \right]$$
(50)

 $\mathcal{L}(z_{\star},z)$ is an even function of (z_{\star},z) , while the gradient of the probability density function is an odd function of (z_{\star},z) . Therefore, the product is an odd function of (z_{\star},z) and the expectation is 0.

$$\nabla_{(e,m)}R(e,m)\Big|_{e=m=0} = (0,0).$$
 (51)

We can use the same computation technique to compute the Hessian of the population loss at initialization. This time the parity arguments does not apply, and we can show numerically that the Hessian is non-null

$$\frac{\partial^2 R}{\partial e^2}\bigg|_{e=m=0}, \frac{\partial^2 R}{\partial m^2}\bigg|_{e=m=0}, \frac{\partial^2 R}{\partial e \partial m}\bigg|_{e=m=0} \neq 0.$$
 (52)

The information exponent We can also compute the *information exponent* by looking at what rate m grows around initialization. We use *spherical gradient descent* because we assumed the norm of the vector w to be costant (as Ben Arous also does in his paper on Information Exponent):

The update rule of w is

$$\boldsymbol{w}^{\tau+1} = \frac{\boldsymbol{w}^{\tau} - \gamma \nabla_{\boldsymbol{w}} \mathcal{L}}{\|\boldsymbol{w}^{\tau} - \gamma \nabla_{\boldsymbol{w}} \mathcal{L}\|_{2}} \sqrt{d},$$

multiplying both sides by w_{\star}/d we get

$$m^{\tau+1} = \frac{m^{\tau} - \gamma \frac{\boldsymbol{w}_{\star} \cdot \nabla_{\boldsymbol{w}} \mathcal{L}}{d}}{\|\boldsymbol{w}^{\tau} - \gamma \nabla_{\boldsymbol{w}} \mathcal{L}\|_{2}} \sqrt{d}.$$

We can assume to be in the gradient flow regime, where the learning rate $\gamma \ll 1$

$$\|\boldsymbol{w} - \gamma \nabla_{\boldsymbol{w}} \mathcal{L}\|_{2} = \sqrt{\|\boldsymbol{w}\|^{2} - 2\gamma \boldsymbol{w} \cdot \nabla_{\boldsymbol{w}} \mathcal{L} + \gamma^{2} \|\nabla_{\boldsymbol{w}} \mathcal{L}\|^{2}} \approx \sqrt{d} \sqrt{1 - 2\gamma \frac{\boldsymbol{w} \cdot \nabla_{\boldsymbol{w}} \mathcal{L}}{d}},$$

that finally lead us to

$$m^{\tau+1} = \left(m^{\tau} - \gamma \frac{\boldsymbol{w}_{\star} \cdot \nabla_{\boldsymbol{w}} \mathcal{L}}{d}\right) \left(1 + \gamma \frac{\boldsymbol{w} \cdot \nabla_{\boldsymbol{w}} \mathcal{L}}{d}\right) \approx m^{\tau} - \gamma \frac{\boldsymbol{w}_{\star} \cdot \nabla_{\boldsymbol{w}} \mathcal{L}}{d} + \gamma \frac{\boldsymbol{w} \cdot \nabla_{\boldsymbol{w}} \mathcal{L}}{d}, \quad (53)$$

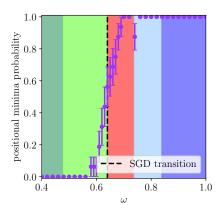


Figure 12: reproduction of Figure 5 for d = 100. The transition is smoother than in the d = 1000. Averaged over 25 runs.

where in the last step we used again the small learning rate limit.

Now we can use the chain rule for getting the gradient in terms of the derivatives we already have

$$\nabla_{w} \mathcal{L} = \frac{\partial \mathcal{L}}{\partial m} \cdot \nabla_{w} m + \frac{\partial \mathcal{L}}{\partial e} \cdot \nabla_{w} e = \frac{\partial \mathcal{L}}{\partial m} \cdot \frac{w_{\star}}{d} + \frac{\partial \mathcal{L}}{\partial e} \cdot \frac{p}{\sqrt{d}}$$

We can rearrange the equation above as

$$\frac{m^{\tau+1}-m^{\tau}}{^{\gamma}\!/\!d} = -\boldsymbol{w}_{\star} \cdot \left(\frac{\partial \mathcal{L}}{\partial m} \cdot \frac{\boldsymbol{w}_{\star}}{d} + \frac{\partial \mathcal{L}}{\partial e} \cdot \frac{\boldsymbol{p}}{\sqrt{d}}\right) + \boldsymbol{w} \cdot \left(\frac{\partial \mathcal{L}}{\partial m} \cdot \frac{\boldsymbol{w}_{\star}}{d} + \frac{\partial \mathcal{L}}{\partial e} \cdot \frac{\boldsymbol{p}}{\sqrt{d}}\right) = -\frac{\partial \mathcal{L}}{\partial m} + m\frac{\partial \mathcal{L}}{\partial m} + e\frac{\partial \mathcal{L}}{\partial e} \cdot \frac{\partial \mathcal{L}}{\partial e} + e\frac{\partial \mathcal{L}}{\partial e} + e\frac{\partial \mathcal{L}}{\partial e} \cdot \frac{\partial \mathcal{L}}{\partial e} + e\frac{\partial \mathcal{L}}{\partial e} + e\frac{\partial \mathcal{L}}{\partial e} \cdot \frac{\partial \mathcal{L}}{\partial e} + e\frac{\partial \mathcal{L}}{\partial e} \cdot \frac{\partial \mathcal{L}}{\partial e} + e\frac{\partial \mathcal$$

where we used the fact $\mathbf{p} \cdot \mathbf{w}_{\star} \approx 0$ in high dimension (as already assumed above), and $\|\mathbf{w}_{\star}\| = d$.

We are interested in what is happening at initialization, therefore all the derivatives should be evaluated around (e,m)=0. We already know that $\nabla_{e,m}\mathcal{L}\big|_{e=m=0}=0$, so we expand the at the next order in m

$$\frac{m^{\tau+1} - m^{\tau}}{\gamma/d} = -m \left. \frac{\partial^2 \mathcal{L}}{\partial m^2} \right|_{m,e=0} + 2me \left. \frac{\partial^2 \mathcal{L}}{\partial m \partial e} \right|_{m,e=0}.$$

We can repeat the same derivation for finding the analogous equation for e:

$$\frac{e^{\tau+1}-e^{\tau}}{\gamma/d}=-e\left.\frac{\partial^2\mathcal{L}}{\partial e^2}\right|_{m,e=0}+2me\left.\frac{\partial^2\mathcal{L}}{\partial m\partial e}\right|_{m,e=0}.$$

Since the hessian of he loss has always a negative eigenvalue, both these equations need $\tau = O(d \log d)$ to escape a neighborhood of initialization, leading to IE = 2.

F.2 Numerical experiments on the transition

In this Appendix we would like to provide more details on the phase diagram of Figure 4.

In Figure 12 we reproduce Figure 5 for d=100, clarifying the *finite size* effects we claimed in the main text. Smaller values of d lead to a more smooth transition, since the gradient flow assumption is less valid. In the limit $d \to \infty$ we expect to see a step function.

The yellow region in the phase diagram of Figure 4 is predicted to have a unique global semantic minima, while SGD is aligned with the positional direction at initialization. The resulting trajectories are shown in Figure 13: the dynamics moves towards the minima direction, get stuck in the local flat (but not critical) region around (e,m)=(1,0), and then it moves towards the global minima. In this region the convergence is very slow. The turning point of the dynamics is a bit misplaced because of the numerical errors given by the loss integration; details on this in Appendix G.

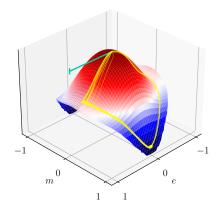


Figure 13: The loss surface in the yellow region of the phase diagram. The SGD trajectory is shown in yellow. The training moves towards the direction orthogonal to the one of global minima, ultimately slowing down the convergence.

F.3 An alternative model without phase transition

In order to highlight the peculiarity of the model (18), we present an example of a model that does not exhibit a phase transition between the positional and semantic regimes. Let's assume that our target is a softmax attention matrix *with* positional encoding

$$y(X) = \operatorname{softmax} \left[(X + P_* / \sqrt{d}) \boldsymbol{w}_* \boldsymbol{w}_*^\top (X + P_* / \sqrt{d})^\top \right]$$

where

$$P_{\star} = \begin{pmatrix} + \boldsymbol{p}_{\star} \\ -\boldsymbol{p}_{\star} \end{pmatrix}, \quad \boldsymbol{w}_{\star} = \sqrt{1 - \omega^2} \boldsymbol{w}_s + \omega \boldsymbol{p}_{\star} \sqrt{d} \quad \text{with} \quad \boldsymbol{w}_s \perp \boldsymbol{p}_{\star}.$$

 ω plays the same role as in the model (18), and we can set $\omega=0$ to get the semantic model, or $\omega=1$ to get the positional model. The difference is that in this case the positional encoding is added to the input of the softmax function, while in the previous model it was added to the output. This change has a significant impact on the behavior of the model.

The gloabl minima of the population loss function does not transition from a positional to a semantic regime, but rather it smoothly move from the semantic to the positional regime as ω increases. We leave the details and numerical experiments for future work.

G Numerical experiments details

All the codes used to run the experiments are available at https://github.com/IdePHICS/Sequence-Single-Index; where details for reproducing figures are not available in the paper, we provide the code to reproduce them in the repository.

The experiments run on a Mac Studio M2 Ultra, within at most few hours for the largest ones. The code is written in Python, using the libraries numpy, scipy, torch and matplotlib. hydra is used to manage the configuration files.

G.1 Details on the integration method of squared loss

All the plot of population loss we showed are a numerical integration of the loss function. As showed in Section 1, the population loss is given by a multivariate Gaussian integral of 2L dimensions, where the mean and the covariance are determined by the sufficient statistics. The integral can't be computed analytically, so we use a custom numerical procedure based on the Gauss-Hermite quadrature.

Let $f: \mathbb{R}^{2L} \to \mathbb{R}$ be a function of 2L variables to be integrated, and let $\mu \in \mathbb{R}^{2L}$ and $\Sigma \in \mathbb{R}^{2L \times 2L}$ be the mean and covariance of the multivariate Gaussian distribution. The integral we want to compute

is

$$I = \int_{\mathbb{R}^{2L}} f(x) \frac{1}{(2\pi)^L} \exp\left(-\frac{1}{2}(x-\mu)^{\top} \Sigma^{-1}(x-\mu)\right) dx$$
 (54)

The numerical procedure is as follows:

(i) Compute the 1D Gauss-Hermite nodes and weights for $N_{\rm int}$ points

$$\{x_i\}_{i=1}^{N_{\text{int}}}$$
 and $\{w_i\}_{i=1}^{N_{\text{int}}}$

where x_i are the nodes given by the roots of the Hermite polynomial $H_{N_{\rm int}}(x)$ and w_i are the corresponding weights, computed as

$$w_i = \frac{2^{N_{\text{int}}}\sqrt{\pi}}{N_{\text{int}}!} \frac{1}{H'_{N_{\text{int}}}(x_i)^2}$$

(ii) Compute the 2L dimensional nodes and weights by taking the Kronecker product of the 1D nodes and weights

$$\{X_i\}_{i=1}^{N_{\mathrm{int}}^{2L}} = \bigotimes_{l=1}^{2L} \{x_i\}_{i=1}^{N_{\mathrm{int}}} \text{ and } \{W_i\}_{i=1}^{N_{\mathrm{int}}^{2L}} = \bigotimes_{l=1}^{2L} \{w_i\}_{i=1}^{N_{\mathrm{int}}}$$

(iii) Let T be the Cholesky decomposition of Σ , i.e. $\Sigma = T^{\top}T$. We can then change the variable to $y = T^{-1}(x - \mu)$ and compute the integral as

$$I = \int_{\mathbb{R}^{2L}} f(x) \frac{1}{(2\pi)^L} \exp\left(-\frac{1}{2}(x-\mu)^\top \Sigma^{-1}(x-\mu)\right) dx = \int_{\mathbb{R}^{2L}} f(Ty+\mu) \frac{1}{(2\pi)^L} \exp\left(-\frac{1}{2}y^\top y\right) dy$$

(iv) The integral can then be approximated, as

$$I \approx \sum_{i=1}^{N_{\text{init}}^{2L}} \left(\prod_{l=1}^{2L} W_{i,l} \right) f(TX_i + \mu)$$

The precision of the integral is obviously regulated by the number of points $N_{\rm int}$ we use. In our experiments, we used $N_{\rm int}=17$, while for the phase diagram in Figure 4 we used $N_{\rm int}=19$. Some effects of this integration error are visible in Figure 13.

NeurIPS Paper Checklist

(i) Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [Yes]

Justification: We state all the result clearly in the Main Result section.

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

(ii) Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: We state all the assumptions made in the paper and discuss their limitations.

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

(iii) Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [Yes]

Justification: All hypothesis and proof of the theoretical results are reported either in the main paper or in the appendix.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

(iv) Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: We release the code used for the experiments, as well as an appendix with details on numerical simulations.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived
 well by the reviewers: Making the paper reproducible is important, regardless of
 whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
- (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

(v) Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: See answer above.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- The authors should provide instructions on data access and preparation, including how
 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

(vi) Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: The released code contains all the details needed to reproduce the experiments, including the hyperparameters. Sometimes the hyperparameters are also specified in the text.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

(vii) Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: When applicable, we report the standard deviation of the results as symmetric error bars.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.

- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error
 of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

(viii) Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: In the Appendix G, we briefly describe our hardware and give upper bounds on running times.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

(ix) Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines?

Answer: [Yes]

Justification: We verified our Reserach Process is in accordance with the NeurIPS Code of Ethics. Concering Societal Impact, our result do not work with any real dataset, nor we release any trained model that could be used for malicious purposes.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

(x) Broader impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [NA]
Justification:

Guidelines: Our result do not work with any real dataset, nor we release any trained model that could be used for malicious purposes.

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

(xi) Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: our result do not work with any real dataset, nor we release any trained model that could be used for malicious purposes.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with
 necessary safeguards to allow for controlled use of the model, for example by requiring
 that users adhere to usage guidelines or restrictions to access the model or implementing
 safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

(xii) Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: Authors developed all the code and data used in the paper, based on open source library like PyTorch, NumPy, and SciPy.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.

- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the
 package should be provided. For popular datasets, paperswithcode.com/datasets
 has curated licenses for some datasets. Their licensing guide can help determine the
 license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

(xiii) New assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [NA]

Justification: The code released is not meant to be used outside this paper. The code should be simple enough to be understood without any documentation.

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

(xiv) Crowdsourcing and research with human subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: The paper does not involve crowdsourcing nor research with human subjects. Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

(xv) Institutional review board (IRB) approvals or equivalent for research with human subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: The paper does not involve crowdsourcing nor research with human subjects. Guidelines:

 The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.

- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

(xvi) Declaration of LLM usage

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does not impact the core methodology, scientific rigorousness, or originality of the research, declaration is not required.

Answer: [NA]

Justification: The core method development in this research does not involve LLMs as any important, original, or non-standard components.

Guidelines:

- The answer NA means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy (https://neurips.cc/Conferences/2025/LLM) for what should or should not be described.