

ON LEARNING LINEAR DYNAMICAL SYSTEMS IN CONTEXT WITH ATTENTION LAYERS

Maria-Luiza Vladarean* † ‡
maria.vladarean@tum.de

Xuhui Zhang* † ‡
xuhui.zhang@tum.de

Suvrit Sra † ‡
s.sra@tum.de

ABSTRACT

This paper studies the expressive power of linear attention layers for in-context learning (ICL) of linear dynamical systems (LDS). We consider training on sequences of inexact observations produced by noise-corrupted LDSs, with all perturbations being Gaussian. Importantly, this non-i.i.d. data setting is a significant step towards modeling real-world scenarios. We provide the optimal weight construction for a single linear-attention layer and show its equivalence to one step of Gradient Descent relative to an autoregression objective of window size one. Guided by experiments, we uncover a connection to a generalization of the Preconditioned Conjugate Gradient method for larger window sizes. We back our findings with numerical evidence. These results add to the existing understanding of transformers’ expressivity as in-context learners and offer plausible hypotheses for recent observations that place their performance on par with that of the Kalman Filter — the optimal model-dependent learner for this setting.

1 INTRODUCTION

This paper contributes towards understanding transformers’ expressive power when learning from *non-i.i.d.* data. In particular, we consider sequences produced by a time-invariant linear dynamical system (LDS) doubly-corrupted by Gaussian noise

$$\begin{cases} \mathbf{x}_{t+1} = \mathbf{A}\mathbf{x}_t + \mathbf{w}_t, \\ y_t = \mathbf{c}^\top \mathbf{x}_t + v_t, \end{cases} \quad (1)$$

where $\mathbf{w}_t \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \Sigma_{\mathbf{w}})$ and $v_t \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma_v^2)$ with mutually independent \mathbf{w}_t and v_t . Our starting point is the well-known ability of transformer-based Large Language Models to perform in-context learning (ICL) (Brown et al., 2020).

ICL boils down to accurately answering a query based on a set of examples given as a textual prefix (“in context”) (Brown et al., 2020). This behaviour is desirable, as it dampens the requirement for expensive data collection and fine-tuning stages (Liu et al., 2023). Present research on ICL spans the spectrum of directions between enhancing it through specialized training and prompt engineering, and building a mechanistic understanding of it — see the comprehensive review of Dong et al. (2022). Our work is aligned with the latter direction.

Broadly, there are two theoretical perspectives on ICL mechanics: a Bayesian view in which transformers recover latent concepts from prompts, thereby performing implicit Bayesian inference (Wang et al., 2023; Jiang, 2023; Wies et al., 2023; Xie et al., 2021, and subsequent works in this direction), and an algorithmic view in which they implement iterative optimization methods in context (Von Oswald et al., 2023a; Giannou et al., 2023; Akyürek et al., 2022; Garg et al., 2022; Ahn et al., 2023; Mahankali et al., 2023; Sander & Peyré, 2024; Von Oswald et al., 2023b; Sander et al., 2024). While the latter view does not account for the “emergent” aspect of “in-the-wild” ICL (Shen et al., 2023), it provides concrete expressions for transformers’ modelling power and identifies the minimal functional unit wielding it — a single, causally-masked, linear attention layer, without positional encoding.

*Equal contribution.

†School of Computation, Information and Technology, Technical University of Munich

‡Munich Center for Machine Learning (MCML)

Works adopting the algorithmic view investigate transformers’ ability to perform linear regression in the forward pass, confirming this through both empirical evidence and formal analysis. The i.i.d. setting is well understood, with results showing that for specific tokenizations, data distributions, and architectural choices, the optimal transformer weights implement gradient-based optimization relative to a context-dependent least-squares loss (Von Oswald et al., 2023a; Mahankali et al., 2023; Ahn et al., 2023; Von Oswald et al., 2023b; Sander et al., 2024). In contrast to the i.i.d. case, our theoretical grasp of the more realistic non-i.i.d. setting is weak. The main hurdle in analyzing this scenario is handling tokens’ statistical dependence on their preceding context. Our work takes the first steps towards unraveling this difficulty.

Specifically, we study the ability of a single linear attention layer to predict observation y_T based on a context of past observations $\{y_t\}_{t=1}^{T-1}$ generated by LDS (1). This setting is well-motivated: firstly, sequence $\{y_t\}_{t=1}^T$ is built on a temporal scaffold similar to that of language-induced tokens, which stands in stark contrast to the i.i.d. setup predominantly addressed by prior works (with few exceptions discussed shortly). Notably, dynamical systems have already been proposed as conveniently flexible models for grammatical sentence formation (Elman, 1995; Tabor et al., 1996; Beim Graben et al., 2004; Belanger & Kakade, 2015), thus making setting (1) particularly relevant. Secondly, considering LDS-produced data places the algorithmic and Bayesian views on closer footing, since LDSs are a subclass of the Hidden Markov Models (Minka, 1999) used in the Bayesian formulation of ICL (Xie et al., 2021). Finally, empirical observations highlight the strong performance of transformers relative to the Kalman Filter (KF) (Kalman, 1960) for predicting y_T , even in regimes where KF is provably optimal (Du et al., 2023). To our knowledge, the underlying mechanism is yet to be understood.

Our goal is to characterize the optimal single linear self-attention layer trained to predict y_T based on the context $\{y_t\}_{t=1}^{T-1}$. We begin by defining a context-dependent loss for time-series data via the improper learning approach to system identification, whereby processes of type (1) are well approximated by autoregressive models. We then identify and interpret the structure of optimally trained linear attention layers as algorithmic steps on the corresponding autoregressive loss. Our contributions are the following.

- C1.** In Theorem 4.1, we prove that for an order-one autoregressive approximation of (1), the optimal linear attention layer implements a step of Gradient Descent on the associated least-squares loss. To our knowledge, this is the first optimality result for LDS data.
- C2.** In Lemma 4.1, we identify a salient banded pattern of the matrices involved in the stationarity condition for generic order- s approximations of (1). We further define a class of parameters that satisfy this structural constraint and empirically observe that minimizers obey it, thereby narrowing the search for the provably optimal linear attention layer when $s \geq 2$.
- C3.** In Section 5, we provide numerical experiments verifying our theory for order-one autoregressive approximations. Furthermore, we connect the tiling pattern of empirically determined minimizers for order- s approximations, $s \geq 2$, with a generalization of the Preconditioned Conjugate Gradient method, thus further highlighting the view of ICL as on-the-fly optimization. To our knowledge, this is the first interpretation of the in-context algorithm for general order- s autoregression.
- C4.** Conceptually, we provide theoretical grounding for interpreting ICL as implicit optimization in the LDS setting, bridging system identification theory and empirical observations of transformer-KF parity.

2 RELATED LITERATURE

We review the studies treating ICL as in-context optimization, together with works on filtering and system identification. Further comparisons are discussed in Section 4.1.

ICL for linear regression with i.i.d data. This line of work studies whether transformers trained on a few-shot learning objective can perform linear regression in-context, and how. Garg et al. (2022); Akyürek et al. (2022); Von Oswald et al. (2023a) provide empirical results in the affirmative, along with possible architecture constructions implementing Gradient Descent (GD) steps relative to a context-induced least squares loss. Through this lens, ICL reduces to on-the-fly optimization

executed in the transformer’s forward pass. Mahankali et al. (2023); Zhang et al. (2024); Ahn et al. (2023) complement these findings by proving that the one-layer linear self-attention implementing a (preconditioned) GD step is a global minimizer of the pretraining loss when covariates are Gaussian and i.i.d. Finally, Zhang et al. (2024) complete the picture by proving that Gradient Flow converges to these global minimizers. Our results extend this line of work to the non-i.i.d. setting.

ICL and system identification. Different from the prior section, the following papers use the standard causal pretraining objective (minimizing prediction error over all sequence positions) and, unless stated otherwise, the results concern a single layer of linear self-attention. Von Oswald et al. (2023b) give a construction implementing a GD step on $\mathcal{L}(\mathbf{W}) := \sum_{i=1}^{t-1} \|\mathbf{W}\mathbf{y}_i - \mathbf{y}_{i+1}\|^2$ in parallel for all positions t , for a specific token augmentation. Sander et al. (2024), further characterize the global minimizers of this causal pretraining loss relative to the noiseless data $\mathbf{y}_{t+1} = \mathbf{A}\mathbf{y}_t$, with \mathbf{A} sampled from the set of commuting orthogonal matrices. Notably, under the same token augmentation, this characterization coincides with the construction of Von Oswald et al. (2023b). Sander et al. (2024) also describe minimizers for architectures using positional encoding, albeit under the restriction of diagonal weight matrices. Zheng et al. (2024) complement these results by showing that, for a diagonal weight initialization and a controlled distribution of \mathbf{y}_0 , Gradient Flow (GF) recovers the aforementioned GD-implementing optimum. Finally, Sander & Peyré (2024) extend these results to arbitrary orthogonal \mathbf{A} s via an infinite-depth attention-only transformer that correctly predicts \mathbf{y}_T in the limit $T \rightarrow \infty$. This result also extends to the standard softmax activation.

Moving away from the noiseless settings above, Cole et al. (2025) establish approximation theoretic results for deep attention-only transformers predicting the sequence $\mathbf{y}_{t+1} = \mathbf{A}\mathbf{y}_t + \mathbf{w}_t$, with $\mathbf{w}_t \sim \mathcal{N}(\mathbf{0}, \sigma_w^2 \mathbf{I})$ and $\mathbf{A} \in \mathbb{S}_{++}^d$. They establish the existence of a $\log(T)$ -depth transformer that achieves uniform-over- \mathbf{A} error of order $\frac{\log(T)}{T}$ for predicting $\mathbb{E}[\mathbf{x}_{T+1} | \mathbf{x}_t, \mathbf{A}]$, and derive a lower bound on the prediction accuracy attainable by a single linear attention layer. Regarding transformers’ capacity, Ziemann et al. (2024) establish that achieving a uniform-in-time error bound for next-step prediction requires a number of parameters at least quadratic in the algebraic multiplicities of \mathbf{A} ’s unstable eigenvalues, as well as a context length at least logarithmic in T .

In summary, these works either study transformers’ ICL performance relative to simplified LDSs or do not address the question of weight optimality. In contrast, we study fully-fledged systems (1) with the aim of characterizing the pretraining loss minimizers in the few-shot training setting.

Transformers and linear filtering. The classical model-based prediction tool for systems of type (1) is the Kalman Filter (KF) (Kalman, 1960). Using knowledge of system parameters, the KF gives the minimum expected squared error estimates $\hat{\mathbf{x}}_i$ of the hidden states \mathbf{x}_i as linear combinations of the past observations \mathbf{y}_i . Transformers as potential implementers of KF were studied by Goel & Bartlett (2024), who prove that a softmax causal attention layer is an arbitrarily good approximator for a given type (1) system. Akram & Vikalo (2024) further construct a transformer emulating the KF. Finally, Du et al. (2023) provide empirical evidence that a GPT-2 architecture (Radford et al., 2019) competes in accuracy with the KF for predicting the next observation in-context, though the mechanism remains unstudied. We partially fill this gap with our present work.

3 PROBLEM FORMULATION & ASSUMPTIONS

Notation. Vectors and matrices are denoted by bold, lowercase and uppercase letters, respectively, with lowercase letters reserved for scalars. Symbols $\mathbf{1}_d$ and $\mathbf{0}_d$ are the all-ones and all-zeros vectors of dimension d , and $\mathbf{1}_{d \times m}$ and $\mathbf{0}_{d \times m}$ are the analogous matrices. Subscripts of the form $\mathbf{A}_{i:j, p:q}$ identify the submatrix formed by rows from i to j and columns from p to q . Unless stated otherwise, $\|\cdot\|$ denotes the Euclidean norm of vectors and the spectral norm of matrices. We denote by $\text{Tr}(\cdot)$ the trace of a matrix, $\langle \cdot, \cdot \rangle$ the inner product, by $\|\cdot\|_F$ its Frobenius norm, and by $\rho(\cdot)$ its spectral radius. We use \mathbf{e}_i for the i^{th} canonical basis vector and \mathbf{I} for the identity matrix. The notation $\mathbb{S}_{+(+) }^d$ defines the cone of symmetric positive-semidefinite (-definite) matrices in $\mathbb{R}^{d \times d}$, respectively. We use \mathbb{S}^{d-1} to denote the unit sphere in \mathbb{R}^d . We use \odot to denote the Hadamard product. Finally, we use $[n]$ when referencing the set of integers $\{1, 2, \dots, n\}$. We write w.p. to mean “with probability”.

The big picture: filtering, system identification, and linear regression. The KF (Kalman, 1960) computes the optimal estimates $\hat{\mathbf{x}}_i$ of \mathbf{x}_i through the system of recursions

$$\left\{ \begin{array}{l} \text{Predict: } \hat{\mathbf{x}}_{t+1|t} = \mathbf{A}\hat{\mathbf{x}}_t \\ \quad \mathbf{P}_{t+1|t} = \mathbf{A}\mathbf{P}_t\mathbf{A}^\top + \Sigma_w \\ \text{Gain: } \mathbf{k}_{t+1} = \mathbf{P}_{t+1|t}\mathbf{c}(\mathbf{c}^\top\mathbf{P}_{t+1|t}\mathbf{c} + \sigma_v)^{-1} \\ \text{Update: } \hat{\mathbf{x}}_{t+1} = \hat{\mathbf{x}}_{t+1|t} + \mathbf{k}_{t+1}(\mathbf{y}_{t+1} - \mathbf{c}^\top\hat{\mathbf{x}}_{t+1|t}) \\ \quad \mathbf{P}_{t+1} = (\mathbf{I}_d - \mathbf{k}_{t+1}\mathbf{c}^\top)\mathbf{P}_{t+1|t}, \end{array} \right. \quad (2)$$

where $\hat{\mathbf{x}}_0$ and the error covariance estimate \mathbf{P}_0 are inputs. Under the Gaussian errors assumption, the state prediction satisfies $\hat{\mathbf{x}}_t = \mathbb{E}[\mathbf{x}_t | y_t, \dots, y_1]$ and, consequently, the forward observation prediction follows $\hat{y}_{t+1} := \mathbf{c}^\top\mathbf{A}\hat{\mathbf{x}}_t = \mathbb{E}[y_{t+1} | y_t, \dots, y_1]$. The fast, constant-time KF predictions, however, require knowing the LDS parameters — a condition generally not satisfied in practice.

Consequently, “proper learning” approaches seek to reconstruct the underlying model by first estimating \mathbf{A} , \mathbf{c} , Σ_w , σ_v through costly parameter identification techniques, and then producing forward observation predictions using the KF (Hamilton, 1995). In contrast, “improper learning” methods eschew structural constraints and solely seek to achieve low error relative to the underlying data distribution and the learning objective (Kozdoba et al., 2019, and references therein). For LDSs, this boils down to expressing the next observation as a linear function of the recent past. Not only does the latter approach avoid parameter estimation, but it also benefits from convex formulations, thus being amenable to classical optimization techniques. Most importantly, for certain LDS classes, improper learning methods can closely track $\mathbb{E}[y_{t+1} | y_t, \dots, y_1]$, as we describe next.

The data-generating process (1) can be rewritten via the KF quantities in expression (2). In particular, Tsiamis & Pappas (2019) express future observations as a function of their past s predecessors,

$$\begin{aligned} [y_{s+1}, \dots, y_{T-1}] &= \mathbf{c}^\top [(\mathbf{A} - \mathbf{k}\mathbf{c}^\top)^{s-1}\mathbf{k}, \dots, \mathbf{k}] [\bar{\mathbf{y}}_1, \dots, \bar{\mathbf{y}}_{T-s-1}] \\ &\quad + \mathbf{c}^\top (\mathbf{A} - \mathbf{k}\mathbf{c}^\top)^s [\hat{\mathbf{x}}_1, \dots, \hat{\mathbf{x}}_{T-s+1}] + [\varepsilon_{s+1}, \dots, \varepsilon_{T-1}], \end{aligned} \quad (3)$$

where $\bar{\mathbf{y}}_t := [y_t, y_{t+1}, \dots, y_{t+s-1}]^\top$, \mathbf{k} is the steady-state gain, and $e_i \in \mathbb{R}$ are i.i.d, zero-mean Gaussian errors. Under KF convergence conditions, quantity $\rho(\mathbf{A} - \mathbf{k}\mathbf{c}^\top) < 1$ makes the second term vanish exponentially in s and thus renders it negligible. We are now in the familiar setting of noisy linear regression, albeit with non-i.i.d. data. The resulting order- s autoregressive process (AR(s)) is associated with the optimization objective

$$\min_{\mathbf{w} \in \mathbb{R}^s} \mathcal{L}_{AR(s)}(\mathbf{w}) := \frac{1}{2(T-s-1)} \sum_{t=1}^{T-s-1} (y_{t+s} - \mathbf{w}^\top \bar{\mathbf{y}}_t)^2. \quad (4)$$

This simplification is the crux of improper learning approaches to system identification (Kozdoba et al., 2019) and becomes relevant in conjunction with the view of transformers as optimizers of a context-induced least-squares objective. Should this view withstand scrutiny in the non-i.i.d. setting, it would suggest that transformers can learn LDS-based time series in context to arbitrary accuracy as a function of s . This is our incentive for characterizing the few-shot pretraining loss minimizers.

To ensure the validity of the autoregressive process approximation, we make the following assumption.

Assumption 3.1 (System assumptions). *LDS (1) has strictly positive definite noise covariances Σ_w and $\sigma_v > 0$. The transition matrix $\mathbf{A} \in \mathbb{R}^{d \times d}$ is marginally stable, with $\rho(\mathbf{A}) \leq 1$, and the pair (\mathbf{A}, \mathbf{c}) is observable, meaning that*

$$\mathbf{O} = [\mathbf{c}, \mathbf{A}^\top\mathbf{c}, \dots, (\mathbf{A}^{d-1})^\top\mathbf{c}]^\top \quad (5)$$

has a column rank of d .

Assumption 3.1 is standard in the literature, and ensures KF convergence (Harrison, 1997) along with the exponential vanishing of the bias term in (3). Furthermore, it ensures the closeness of forward observation predictions given by the KF with those produced by a linear autoregressive predictor determined by expression (4) (Kozdoba et al., 2019).

Transformer architecture. Transformers (Vaswani et al., 2017) are neural architectures performing sequence-to-sequence mapping. For input tokens $\mathbf{S}_N = [\mathbf{s}_1, \dots, \mathbf{s}_N]^\top \in \mathbb{R}^{N \times p}$, the transformer produces a corresponding $\hat{\mathbf{S}}_N = [\hat{\mathbf{s}}_1, \dots, \hat{\mathbf{s}}_N]^\top \in \mathbb{R}^{N \times p}$ by dynamically mixing tokens via its attention mechanism. An L -layer transformer $\mathcal{T}_\theta : \mathbb{R}^{N \times p} \rightarrow \mathbb{R}^{N \times p}$ parametrized by $\theta = [\theta_i]_{i=1}^L$ is a composition of blocks $\mathcal{T}_L = \mathcal{T}_{\theta_1} \circ \dots \circ \mathcal{T}_{\theta_L}$. Each \mathcal{T}_{θ_i} is a sequence-to-sequence function given by

$$\mathcal{T}_{\theta_i}(\mathbf{S}) := (\text{MLP}_{\theta_i^{\text{MLP}}} \circ \mathcal{A}_{\theta_i^{\text{att}}})(\mathbf{S}),$$

where $\text{MLP}_{\theta_i^{\text{MLP}}}$ is a multilayer perceptron and $\mathcal{A}_{\theta_i^{\text{att}}}$ is the attention mapping. This paper studies the simplified block $\mathcal{T}_\theta(\mathbf{S}) := \mathcal{A}_\theta(\mathbf{S})$, with $L = 1$ and $\text{MLP}_{\theta_1^{\text{MLP}}} = \text{Id}$.

The causal h -headed attention block with residual connections is given by

$$\mathcal{A}_\theta(\mathbf{S}) := \mathbf{S} + \sum_{h=1}^H \sigma \left(\mathbf{M} \odot \frac{1}{\tau} \mathbf{S} \mathbf{W}_Q^h (\mathbf{W}_K^h)^\top \mathbf{S}^\top \right) \mathbf{S} \mathbf{W}_V^h \mathbf{W}_O^h, \quad (6)$$

where the parameters $\theta = [\mathbf{W}_Q^h, \mathbf{W}_K^h, \mathbf{W}_V^h, \mathbf{W}_O^h]_{h=1}^H$ are the query, key, value, and projection matrices, respectively; $\tau > 0$ is a scaling constant; σ is the softmax normalizing function applied row-wise; and $\mathbf{M} \in \mathbb{R}^{N \times N}$, with $M_{i,j} = 1$ if $i \geq j$ and $-\infty$ otherwise, is a causal mask.

Similar to prior works (Von Oswald et al., 2023a; Ahn et al., 2023; Mahankali et al., 2023), we restrict our study to the analytically tractable setting of single-headed linear attention (Katharopoulos et al., 2020). Without loss of expressivity, we drop the projection matrix \mathbf{W}_O and consider the $\mathbf{W}_Q \mathbf{W}_K^\top$ as a single matrix $\mathbf{W}_{QK} \in \mathbb{R}^{p \times p}$. Since we’re working in the few-shot scenario, we’re concerned solely with predicting the final position as

$$\hat{\mathbf{s}}_N := \mathcal{T}_\theta(\mathbf{S})_N = \mathbf{s}_N + \frac{1}{N-1} \mathbf{W}_V^\top \sum_{i=1}^{N-1} \mathbf{s}_i \mathbf{s}_i^\top \mathbf{W}_{QK}^\top \mathbf{s}_N, \quad (7)$$

where we set $\tau = N - 1$ and omit the last sum element due to a token asymmetry discussed next.

Token construction. We use the token construction approach of Von Oswald et al. (2023a); Ahn et al. (2023); Mahankali et al. (2023). The input matrix \mathbf{Y}_0 built for AR(s) data (4) is

$$\mathbf{Y}_0 := \begin{bmatrix} \bar{\mathbf{y}}_1 & \bar{\mathbf{y}}_2 & \cdots & \bar{\mathbf{y}}_{T-s-1} & \bar{\mathbf{y}}_{T-s} \\ y_{s+1} & y_{s+2} & \cdots & y_{T-1} & 0 \end{bmatrix} \quad (8)$$

where $s \geq 1$ is the window size of the AR process and $\mathbf{Y}_0 \in \mathbb{R}^{(s+1) \times (T-s)}$. The corresponding scaling constant τ in (6) becomes $T - s - 1$. The last column represents the “test” token, whose last entry is filled in the transformer’s forward pass by the estimate \hat{y}_T . This asymmetry motivates the last term’s removal in (7).

Lemma 3.1 below ensures the existence (by construction) of a linear attention layer producing \mathbf{Y}_0 from the raw sequence $\{y_t\}_{t=1}^T$. Its proof is deferred to Appendix C.1 due to space constraints.

Lemma 3.1. *For a given $s \geq 1$, there exists an $s + 1$ -headed linear attention layer with positional encoding which transforms input sequences $[y_1, y_2, \dots, y_T]^\top$ into*

$$\begin{bmatrix} \bar{\mathbf{y}}_1 & \cdots & \bar{\mathbf{y}}_{T-s} & \mathbf{0}_{(s-1) \times (T-s-1)} \\ y_{s+1} & \cdots & 0 & \mathbf{0}_{T-s-1}^\top \end{bmatrix}^\top.$$

The latter quantity is essentially equivalent to \mathbf{Y}_0 as defined in equation (8).

Data distribution, loss function, and training paradigm. We consider trajectories $\{y_i\}_{i=1}^T$ sampled from system (1), where each trajectory corresponds to different parameters \mathbf{A} , \mathbf{c} and \mathbf{x}_0 . We assume $\mathbf{x}_0 \sim \mathcal{N}(\mathbf{0}_d, \Sigma_{\mathbf{x}_0})$, and impose the following condition on the distributions of \mathbf{A} and \mathbf{c} .

Assumption 3.2 (LDS family). *The system matrix $\mathbf{A} \in \mathbb{R}^{d \times d}$ is sampled from a centrally symmetric distribution supported on $\{\mathbf{M} \in \mathbb{R}^{d \times d} \mid \rho(\mathbf{M}) \leq 1\}$, for which it holds that*

$$\mathbb{P}(\{\mathbf{A} \mid \exists i, j \in [d], \text{ s.t. } \lambda_i(\mathbf{A}) = \lambda_j(\mathbf{A})\}) = 0. \quad (9)$$

In other words, \mathbf{A} has a simple spectrum almost surely. The observation vector $\mathbf{c} \in \mathbb{R}^d$ is sampled independently, from a distribution that is absolutely continuous w.r.t. the Lebesgue measure over \mathbb{R}^d .

with $\mathbf{R} := \mathbf{1}_{\lfloor \frac{s+1}{2} \rfloor \times \lceil \frac{s+1}{2} \rceil} \otimes \begin{bmatrix} \star & 0 \\ 0 & \star \end{bmatrix}$ and $\mathbf{r} := \mathbf{1}_{\lceil \frac{s+1}{2} \rceil} \otimes \begin{bmatrix} 0 \\ \star \end{bmatrix}$, ensure that $LHS(12)_{r,\ell} = 0$ whenever $r + l \in 2\mathbb{N}$ and $s + j \in 2\mathbb{Z}$, or $r + l \in 2\mathbb{N} + 1$ and $s + j \in 2\mathbb{Z} + 1$.

Note that these parameters recover the zero-structure of $\mathbf{B}_0(s, j)$ and $\mathbf{B}_1(s, j)$. Lemma 4.1 can be understood as a structure-based narrowing of the parameter class likely to hold minimizers of loss (10). The experiments in Section 5 give empirical support for this claim.

Our second step is to use structure (14) to identify a global minimizer of loss (10) for AR(1)-type tokens. The result is stated in Theorem 4.1, whose proof we defer to Appendix D.4.

Theorem 4.1. *Let \mathbf{Y}_0 encode the input tokens for $s = 1$. Then, the optimal parameters $\theta^* = (\mathbf{W}_{QK}^*, \mathbf{W}_V^*)$ of a single linear self-attention layer with respect to loss $\mathcal{L}(\theta)$ are*

$$\mathbf{W}_{QK}^* = \begin{bmatrix} \frac{(T-2)\mathbb{E}[y_{T-1}y_T \sum_{i=1}^{T-2} y_i y_{i+1}]}{\mathbb{E}[y_{T-1}^2 (\sum_{i=1}^{T-2} y_i y_{i+1})^2]} & 0 \\ 0 & 0 \end{bmatrix}, \mathbf{W}_V^* = \begin{bmatrix} 0 & 0 \\ 0 & 1 \end{bmatrix}, \quad (15)$$

up to rescaling with a nonzero constant.

Broadly, the proof of Theorem 4.1 encounters two difficulties compared to the i.i.d. case: the number of terms that need to be matched in satisfying the first-order optimality condition, and the full-history dependence of the data. We address the first obstacle by using the structural simplification of Lemma 4.1, and the second by invoking Isserlis’ theorem (Isserlis, 1918), which provides a tractable decomposition of the higher-order moments of the data. Details are given in Appendix D.3.

Notably, a forward pass using the optimal parameters (15) amounts to the prediction given after one GD step on $\mathcal{L}_{AR(1)}(w)$ starting from $w_0 = 0$. We thus recover the ICL-as-optimization view upheld by works in the i.i.d. setting (Ahn et al., 2023; Mahankali et al., 2023) but for LDS-produced data.

4.1 DISCUSSION

To our knowledge, the only other architecture proposed for handling noisy observations y_t of type (1) is given by Cole et al. (2025). Theirs is part of a proof of existence by construction and, as such, is not accompanied by confirming experimental evidence. Different from us, they propose an attention-only transformer that unrolls a *modified Richardson iteration* meant to estimate $(\frac{1}{T} \sum_{t=1}^T \mathbf{x}_{t+1} \mathbf{x}_t^\top)^{-1} (\frac{1}{T} \sum_{t=1}^T \mathbf{x}_t \mathbf{x}_t^\top)^{-1}$ for a simpler LDS with direct state access. Their construction extends to the setting of objective (4) via the work of Tsiamis & Pappas (2019), who give a high probability result for the existence of $(\sum_{t=1}^{T-s-1} \bar{\mathbf{y}}_t \bar{\mathbf{y}}_t^\top)^{-1}$ under our assumptions. However, their transformer has a minimum of two layers, of which the first is fixed, therefore providing no guarantee that training will recover it. Our results take a first step towards filling this gap.

Tangentially, Akram & Vikalo (2024) construct a transformer emulating the KF, contingent on knowledge of the system parameters and an elaborate token augmentation scheme. While this architecture is capable of computing the forward KF observation \hat{y}_T , it relies on ideal knowledge of LDS (1) which is rarely encountered in practice.

Theorem 4.1 sets forth a plausible hypothesis for prior experiments (Du et al., 2023, Fig. 2) using a GPT-2 architecture trained autoregressively with data (1) for stable $\mathbf{A} \in \mathbb{S}_{++}^d$. Their results highlight the transformer’s competitive performance relative to the KF for predicting the next observation of a previously unseen sequence, in-context. These experiments suggest an implicit form of system identification might be executed in the forward pass, though the mechanism remains unstudied. Through the ICL-as-optimization lens, we can interpret the high accuracy of GPT-2’s in-context predictions as a possible consequence of Theorem 2 of Kozdoba et al. (2019). Importantly, the latter result implies that for an arbitrary, finite family S of LDSs (1) and an $\varepsilon > 0$, there exists a window-length $s(\varepsilon)$ such that the optimal $AR(s(\varepsilon))$ predictor incurs an average error that is at least as good, up to ε , as that of the forward observation prediction \hat{y}_{t+1} of the best KF in S . Our results take the first step in formally exploring this hypothesis.

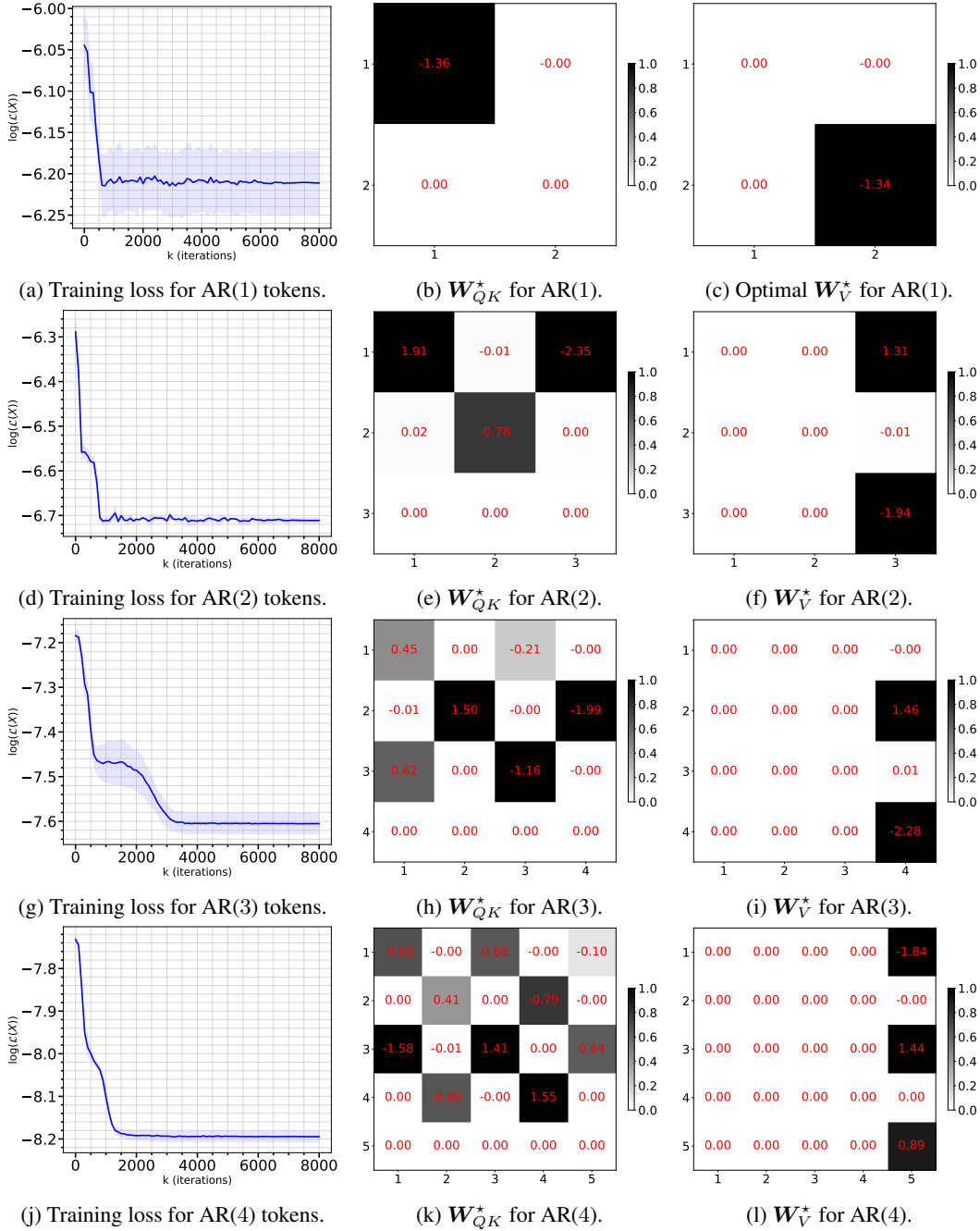


Figure 1: Experimental results for AR(1–4) tokens showing the optimally-trained attention parameters.

5 EXPERIMENTS

We now present numerical evidence supporting our theory. All experiments were implemented in Python 3.12 and run on a ThinkPad T14p with 32 GB RAM and a 22-core Intel Core™ Ultra 9 185H processor. The code is available at <https://github.com/XHZhang01/icl-for-lds-data>.

We train architecture (7) on sequences $\{y_t\}_{t=1}^T$, $T = 30$, each sampled from a different LDS of type (1) with a hidden state dimension $d = 5$. The number of training iterations is 8000 for all cases,

with an increase in the batch size for every increase in window size, starting from 3000 for AR(1). A fresh batch of LDSs is sampled at every iteration. The experiments cover the following four settings.

- For each sequence, we sample \mathbf{A} 's diagonal entries uniformly in the interval $[-1, 1]$ and set $\mathbf{c} = \mathbf{1}_d$. The noise covariances are set to $\Sigma_w = 1e-2\mathbf{I}$ and $\sigma_v^2 = 1e-2$. The results are depicted in Figure 1 for AR(1-4) tokens.
- For each sequence, we sample a vector $\mathbf{v} \sim \text{Unif}([-1, 1]^d)$ and a matrix $\mathbf{Q} \sim \text{Haar}(O(d))$ independently, and set $\mathbf{A} = \mathbf{Q}^\top \text{diag}(\mathbf{v})\mathbf{Q}$. We further sample $\mathbf{c} \sim \text{Unif}([-2, 2]^d)$ independently. The noise covariances are set to $\Sigma_w = 1e-2\mathbf{I}$ and $\sigma_v^2 = 1e-2$. Experiments for AR(1-4)-type tokens are provided in Figure 3 in Appendix B.3.
- For each sequence, we sample a vector $\mathbf{v} \sim \text{Unif}([-1, 1]^d)$ and a matrix $\mathbf{Q} \sim \text{Haar}(O(d))$ independently, and set $\mathbf{A} = \mathbf{Q}^\top \text{diag}(\mathbf{v})\mathbf{Q}$. We further sample $\mathbf{c} \sim \text{Unif}([-0.5, 0.5]^d)$ independently. We fix the covariance $\Sigma_w = \mathbf{Q}_w^\top \text{diag}(1e-2 \cdot [0.9, 0.95, 1.0, 1.05, 1.1])\mathbf{Q}_w$ for all systems, for arbitrary $\mathbf{Q}_w \sim \text{Haar}(O(d))$. We fix $\sigma_v^2 = 1e-2$. Experiments for AR(1-3)-type tokens are provided in Figure 4 in Appendix B.4.
- For each sequence, we independently sample a vector $\mathbf{v} \sim \text{Unif}([-1, 1]^d)$ and a matrix $\mathbf{P} = [p_{i,j}]_{i,j=1}^d$ with $p_{i,j} \stackrel{\text{i.i.d.}}{\sim} \mathcal{U}([-1, 1])$, and set $\mathbf{A} = \mathbf{P}^{-1} \text{diag}(\mathbf{v})\mathbf{P}$. We further sample $\mathbf{c} \sim \text{Unif}([-1, 1]^d)$. The noise covariances are set to $\Sigma_w = 1e-2\mathbf{I}$ and $\sigma_v^2 = 1e-2$. Experiments for AR(1-3)-tokens are provided in Figure 5 in Appendix B.5.

All the settings above have $\mathbf{x}_0 \sim \mathcal{N}(0, \sigma_0^2\mathbf{I})$, $\sigma_0^2 = 1e-2$. Note that we could have used any other centrally symmetric distribution with marginals supported on $[-1, 1]$ for the sampling of the diagonal \mathbf{v} , e.g., $\text{Unif}(\mathbb{S}^{d-1})$ — uniform on the unit sphere; $\text{Unif}(\{x \in \mathbb{R}^d : \|x\|_2 \leq 1\})$ — uniform inside the unit ball, etc. These sampling schemes obey Assumption 3.2, as proven in Appendix E.1. All results are averaged over 3 random seeds. The weights are learned using AdamW (Loshchilov & Hutter, 2017) with gradient clipping and a learning rate schedule consisting of a linear warm-up phase followed by cosine annealing (Loshchilov & Hutter, 2016). A full list of hyperparameters is provided in Appendix B.

Subplots (b,c) in Figures 1, 3, 4 and 5 show optima conforming to Theorem 4.1 for AR(1)-type tokens. Moreover, subplots (e,f,h,i,k,l) in Figures 1 and 3, subplots (e,f,h,i) in Figure 4 and subplots (e,f) in Figure 5 confirm the pattern dictated by Lemma 4.1 for general $s \geq 2$. Finally, we empirically show that in setting (a), the weights converge to the sparsity pattern predicted by Lemma 4.1 in terms of the Jaccard distance of their supports. Setup details and results are given in Appendix B.2.

Interpreting the sparsity pattern for AR(s) $s \geq 2$. A quick calculation of the forward pass reveals that weights trained to optimality with AR(s) tokens for $s \geq 2$ *do not* implement standard GD in the forward pass, but an iteration reminiscent of the Preconditioned Conjugate Gradient method (PCG) (Hestenes et al., 1952; Shewchuk et al., 1994), or more generally, to Krylov subspace methods (Saad, 2003, chapters 6 and 7).

We start by observing that the forward pass factor $\frac{1}{T-s-1}\bar{\mathbf{Y}}$ (with $\bar{\mathbf{Y}}$ defined in (11)) has a meaningful block structure involving the gradient at zero and the Hessian of $\mathcal{L}_{AR(s)}(\mathbf{w})$,

$$\frac{1}{T-s-1}\bar{\mathbf{Y}} = \begin{bmatrix} \nabla^2 \mathcal{L}_{AR(s)} & \nabla \mathcal{L}_{AR(s)}(0) \\ \nabla \mathcal{L}_{AR(s)}(0)^\top & \gamma \end{bmatrix}, \quad (16)$$

where $\gamma := \frac{1}{T-s-1} \sum_{t=1}^{T-s-1} y_{t+s}^2 \in \mathbb{R}$ and $\nabla^2 \mathcal{L}_{AR(s)} \in \mathbb{R}^{s \times s}$ denotes the constant Hessian. Together with the parameter structure of Lemma 4.1 and the experiments, we use expression (16) to rewrite the transformer-induced predictor in a manner that highlights its resemblance to that obtained after two PCG steps on $\mathcal{L}_{AR(s)}(\mathbf{w})$ starting from $\mathbf{w}_0 = \mathbf{0}$. We describe the case for even s , with the odd case following similarly. Let $s = 2k$, $k \in \mathbb{N}$ and $N := \frac{(s+1)^2+1}{2}$. Consider weights

$$\mathbf{W}_{QK} = \left[\begin{array}{ccc|c} c_1 & 0 & c_2 & c_{k+1} \\ 0 & c_{k+2} & 0 & \ddots & 0 \\ \vdots & \vdots & \ddots & \ddots & \vdots \\ 0 & c_{N-2k} & 0 & & 0 \\ \hline & \mathbf{0}_s^\top & & & 0 \end{array} \right] \in \mathbb{R}^{s+1 \times s+1}, \quad \mathbf{W}_V = \left[\begin{array}{c|c} & c_{N-k} \\ \mathbf{0}_{s \times s} & 0 \\ & c_{N-k+1} \\ & \vdots \\ & 0 \\ \hline 0 \dots 0 & c_N \end{array} \right] \in \mathbb{R}^{s+1 \times s+1}$$

where $c_i \in \mathbb{R}, \forall i \in [N]$. Renaming the top left $s \times s$ block of $\mathbf{W}_{\mathbf{QK}}$ as \mathbf{P} , the top right $s \times 1$ block as \mathbf{p} , and the top right $s \times 1$ block of $\mathbf{W}_{\mathbf{V}}$ as \mathbf{q} , the transformer-induced linear predictor $\frac{1}{T-s-1} \mathbf{W}_{\mathbf{QK}} \bar{\mathbf{Y}} \mathbf{W}_{\mathbf{V}}[:,s+1]$ is

$$\mathbf{P} \nabla^2 \mathcal{L}_{AR(s)} \mathbf{q} + \xi_1 \mathbf{P} \nabla \mathcal{L}_{AR(s)}(\mathbf{0}) + \xi_2 \mathbf{p}, \quad (17)$$

where $\xi_1 = c_N \in \mathbb{R}$ and $\xi_2 = \mathbf{q}^\top \nabla \mathcal{L}_{AR(s)}(\mathbf{0}) + c_N \gamma \in \mathbb{R}$.

Comparatively, the PCG (Shewchuk et al., 1994, p. 51) predictor obtained after two steps on loss $\mathcal{L}_{AR(s)}$ with preconditioner \mathbf{P}^{-1} starting from $\mathbf{w}_0 = \mathbf{0}$ is (see Appendix E.2)

$$\tau_1 \mathbf{P} \nabla^2 \mathcal{L}_{AR(s)} \mathbf{P} \nabla \mathcal{L}_{AR(s)}(\mathbf{0}) + \tau_2 \mathbf{P} \nabla \mathcal{L}_{AR(s)}(\mathbf{0}), \quad (18)$$

where $\tau_1, \tau_2 \in \mathbb{R}$ are iteration-dependent constants. Note that we have suspended the requirement for \mathbf{P} 's symmetry here — we will address this drawback shortly.

We make a few observations on the similarities and differences between the two predictors. The second terms of (17) and (18) coincide up to scaling. The first terms of these predictors represent directions mapped by $\mathbf{P} \nabla^2 \mathcal{L}_{AR(s)}$, which is typical of Krylov subspace methods. Moreover, if we initialize the conjugate direction of PCG to the vector \mathbf{q} , the first terms also coincide up to scaling. This prompts the natural question of whether directions $\mathbf{P} \nabla^2 \mathcal{L}_{AR(s)} \mathbf{P} \nabla \mathcal{L}_{AR(s)}(\mathbf{0})$ and $\mathbf{P} \nabla^2 \mathcal{L}_{AR(s)} \mathbf{q}$ are aligned, or at least significantly so. As a preliminary test, we compute this alignment in terms of cosine similarity when replacing $\nabla^2 \mathcal{L}_{AR(s)}$ and $\mathcal{L}_{AR(s)}(\mathbf{0})$ with empirical estimates of their expectations. Using the weights reported for the AR(4) case in settings a) and b) (Figures 1 and 3), batches of 16000 sequences for estimating the expectations and averaging the result over five random seeds, we obtain a cosine similarity of 0.88 ± 0.05 for setting a) and 0.93 ± 0.01 for setting b). The values suggest a strong alignment, which gives weight to this line of interpretation.

In terms of predictor dissimilarities, the set of directions whose linear combination yields (17) contains the additional vector \mathbf{p} . This structure is typical of so-called augmented Krylov methods, whereby “correction directions” are added to the standard Krylov search space to compensate for ill-behaved modes of $\mathbf{P} \nabla^2 \mathcal{L}_{AR(s)}$ (Carpenter et al., 2010, and references therein). Broadly, they can be seen as PCG generalizations, and have the added benefit of accommodating asymmetric \mathbf{P} s, thus resolving the earlier drawback. We test whether \mathbf{p} acts as a correction by evaluating the residual obtained prior to and after incorporating \mathbf{p} . We use the same setup as above, and compute the cosine of the angle between the preconditioned residual evaluated at \mathbf{q} (which coincides with \mathbf{w}_1 in PCG with non-standard initial conjugate direction), and the direction $\mathbf{P} \nabla^2 \mathcal{L}_{AR(s)} \mathbf{p}$. Anti-alignment would suggest that \mathbf{p} moves the predictor in a residual reduction direction. Indeed, we observe a cosine of $-0.99 \pm 7e^{-5}$ for the AR(4) case in setting a), and $-0.99 \pm 3e^{-5}$ for the AR(4) case in setting b). Taken together, these preliminary results warrant further exploration of this interpretation.

Finally, we remark that our interpretation of the AR(s) case does not contradict the plain GD step observed for AR(1), since PCG variants collapse to GD for one-dimensional covariates.

6 LIMITATIONS & FUTURE DIRECTIONS

We have sketched a path to understanding how attention layers learn LDSs in context by leveraging results from the improper learning literature on system identification. Our study fully characterizes the learning of order-one autoregressive LDS approximations (Theorem 4.1) and narrows the class of weights that hold minimizers for higher-order approximations (Lemma 4.1), with all findings experimentally confirmed. These contributions provide the building blocks for closing further gaps, the most pressing of which is describing the optima of loss (10) for $s \geq 2$. A meaningful intermediate result here is quantifying the approximation power of optima obtained in the stationary LDS regime, where the problem structure further simplifies due to $\bar{\mathbf{Y}}$ being Toeplitz. Such an assumption is reasonable, since systems with $\rho(\mathbf{A}) < 1$ reach stationarity exponentially fast. Further valuable directions concern bridging the theory-practice gap, such as extending the analysis to typical causal pretraining objectives, and empirically probing the parallel between Krylov subspace methods and the predictor modeled by a multi-layer transformer.

ACKNOWLEDGMENTS

The authors are grateful to Xiang Cheng and Anastasios Tsiamis for helpful discussions throughout the development of this work. The authors also sincerely thank the anonymous reviewers for their constructive feedback which helped improve and clarify this manuscript. All the authors of this work were generously supported by the Alexander von Humboldt foundation.

REFERENCES

- Kwangjun Ahn, Xiang Cheng, Hadi Daneshmand, and Suvrit Sra. Transformers learn to implement preconditioned gradient descent for in-context learning. *Advances in Neural Information Processing Systems*, 36:45614–45650, 2023. 1, 2, 3, 5, 6, 7, 28
- Usman Akram and Haris Vikalo. Can transformers in-context learn behavior of a linear dynamical system?, 2024. URL <https://arxiv.org/abs/2410.16546>. 3, 7
- Ekin Akyürek, Dale Schuurmans, Jacob Andreas, Tengyu Ma, and Denny Zhou. What learning algorithm is in-context learning? investigations with linear models. *arXiv preprint arXiv:2211.15661*, 2022. 1, 2
- Peter Beim Graben, Bryan Jurish, Douglas Saddy, and Stefan Frisch. Language processing by dynamical systems. *International Journal of Bifurcation and Chaos*, 14(02):599–621, 2004. 2
- David Belanger and Sham Kakade. A linear dynamical system model for text. In *International Conference on Machine Learning*, pp. 833–842. PMLR, 2015. 2
- 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. *Advances in neural information processing systems*, 33:1877–1901, 2020. 1
- Mark H Carpenter, C Vuik, Peter Lucas, Martin vanGijzen, and Hester Bijl. A general algorithm for reusing krylov subspace information. i. unsteady navier-stokes. Technical report, 2010. 10
- Frank Cole, Yulong Lu, Tianhao Zhang, and Yuxuan Zhao. In-context learning of linear dynamical systems with transformers: Error bounds and depth-separation. *arXiv preprint arXiv:2502.08136*, 2025. 3, 7
- Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Jingyuan Ma, Rui Li, Heming Xia, Jingjing Xu, Zhiyong Wu, Tianyu Liu, et al. A survey on in-context learning. *arXiv preprint arXiv:2301.00234*, 2022. 1
- Zhe Du, Haldun Balim, Samet Oymak, and Necmiye Ozay. Can transformers learn optimal filtering for unknown systems? *IEEE Control Systems Letters*, 7:3525–3530, 2023. 2, 3, 7
- Jeffrey L Elman. Language as a dynamical system. *Mind as motion: Explorations in the dynamics of cognition*, pp. 195–223, 1995. 2
- Shivam Garg, Dimitris Tsipras, Percy S Liang, and Gregory Valiant. What can transformers learn in-context? a case study of simple function classes. *Advances in Neural Information Processing Systems*, 35:30583–30598, 2022. 1, 2
- Angeliki Giannou, Shashank Rajput, Jy-yong Sohn, Kangwook Lee, Jason D Lee, and Dimitris Papailiopoulos. Looped transformers as programmable computers. In *International Conference on Machine Learning*, pp. 11398–11442. PMLR, 2023. 1
- Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedforward neural networks. In *Proceedings of the thirteenth international conference on artificial intelligence and statistics*, pp. 249–256. JMLR Workshop and Conference Proceedings, 2010. 15, 16, 17
- Gautam Goel and Peter Bartlett. Can a transformer represent a kalman filter? In *6th Annual Learning for Dynamics & Control Conference*, pp. 1502–1512. PMLR, 2024. 3
- James D Hamilton. Time series analysis, 1995. 4

- P Jeff Harrison. Convergence and the constant dynamic linear model. *Journal of Forecasting*, 16(5): 287–292, 1997. 4
- Magnus R Hestenes, Eduard Stiefel, et al. Methods of conjugate gradients for solving linear systems. *Journal of research of the National Bureau of Standards*, 49(6):409–436, 1952. 9
- Leon Isserlis. On a formula for the product-moment coefficient of any order of a normal frequency distribution in any number of variables. *Biometrika*, 12(1/2):134–139, 1918. 7, 24
- Paul Jaccard. Étude comparative de la distribution florale dans une portion des alpes et des jura. *Bull Soc Vaudoise Sci Nat*, 37:547–579, 1901. 17
- Hui Jiang. A latent space theory for emergent abilities in large language models. *arXiv preprint arXiv:2304.09960*, 2023. 1
- Rudolph Emil Kalman. A new approach to linear filtering and prediction problems. *Journal of basic Engineering*, 82(1):35–45, 1960. 2, 3, 4
- Angelos Katharopoulos, Apoorv Vyas, Nikolaos Pappas, and François Fleuret. Transformers are rnns: Fast autoregressive transformers with linear attention. In *International conference on machine learning*, pp. 5156–5165. PMLR, 2020. 5
- Mark Kozdoba, Jakub Marecek, Tigran Tchrakian, and Shie Mannor. On-line learning of linear dynamical systems: Exponential forgetting in kalman filters. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pp. 4098–4105, 2019. 4, 7
- Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing. *ACM computing surveys*, 55(9):1–35, 2023. 1
- Ilya Loshchilov and Frank Hutter. Sgdr: Stochastic gradient descent with warm restarts. *arXiv preprint arXiv:1608.03983*, 2016. 9, 15, 16, 17
- Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017. 9, 15, 16, 17
- Arvind Mahankali, Tatsunori B Hashimoto, and Tengyu Ma. One step of gradient descent is provably the optimal in-context learner with one layer of linear self-attention. *arXiv preprint arXiv:2307.03576*, 2023. 1, 2, 3, 5, 7
- Thomas P Minka. From hidden markov models to linear dynamical systems. Technical report, Citeseer, 1999. 2
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019. 3
- H. L. Royden and P. M. Fitzpatrick. *Real Analysis*. Prentice Hall, 4th edition, 2010. ISBN 978-0131437470. 23
- Yousef Saad. *Iterative methods for sparse linear systems*. SIAM, 2003. 9
- Michael E Sander and Gabriel Peyré. Towards understanding the universality of transformers for next-token prediction. *arXiv preprint arXiv:2410.03011*, 2024. 1, 3
- Michael E Sander, Raja Giryes, Taiji Suzuki, Mathieu Blondel, and Gabriel Peyré. How do transformers perform in-context autoregressive learning? *arXiv preprint arXiv:2402.05787*, 2024. 1, 2, 3
- Lingfeng Shen, Aayush Mishra, and Daniel Khashabi. Do pretrained transformers learn in-context by gradient descent? *arXiv preprint arXiv:2310.08540*, 2023. 1
- Jonathan Richard Shewchuk et al. An introduction to the conjugate gradient method without the agonizing pain. Technical report, Carnegie-Mellon University. Department of Computer Science Pittsburgh, 1994. 9, 10, 34

- Whitney Tabor, Christopher Juliano, and Michael Tenenhaus. A dynamical system for language processing. In *Proceedings of the Annual Meeting of the Cognitive Science Society*, volume 18, 1996. URL <https://escholarship.org/uc/item/78r6h0cg>. 2
- Anastasios Tsiamis and George J Pappas. Finite sample analysis of stochastic system identification. In *2019 IEEE 58th Conference on Decision and Control (CDC)*, pp. 3648–3654. IEEE, 2019. 4, 7
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017. 5
- Johannes Von Oswald, Eyvind Niklasson, Ettore Randazzo, João Sacramento, Alexander Mordvintsev, Andrey Zhmoginov, and Max Vladymyrov. Transformers learn in-context by gradient descent. In *International Conference on Machine Learning*, pp. 35151–35174. PMLR, 2023a. 1, 2, 5, 22
- Johannes Von Oswald, Maximilian Schlegel, Alexander Meulemans, Seijin Kobayashi, Eyvind Niklasson, Nicolas Zucchet, Nino Scherrer, Nolan Miller, Mark Sandler, Max Vladymyrov, et al. Uncovering mesa-optimization algorithms in transformers. *arXiv preprint arXiv:2309.05858*, 2023b. 1, 2, 3
- Xinyi Wang, Wanrong Zhu, Michael Saxon, Mark Steyvers, and William Yang Wang. Large language models are latent variable models: Explaining and finding good demonstrations for in-context learning. *Advances in Neural Information Processing Systems*, 36:15614–15638, 2023. 1
- Noam Wies, Yoav Levine, and Amnon Shashua. The learnability of in-context learning. *Advances in Neural Information Processing Systems*, 36:36637–36651, 2023. 1
- Sang Michael Xie, Aditi Raghunathan, Percy Liang, and Tengyu Ma. An explanation of in-context learning as implicit bayesian inference. *arXiv preprint arXiv:2111.02080*, 2021. 1, 2
- Ruiqi Zhang, Spencer Frei, and Peter L Bartlett. Trained transformers learn linear models in-context. *Journal of Machine Learning Research*, 25(49):1–55, 2024. 3
- Chenyu Zheng, Wei Huang, Rongzhen Wang, Guoqiang Wu, Jun Zhu, and Chongxuan Li. On mesa-optimization in autoregressively trained transformers: Emergence and capability. *Advances in Neural Information Processing Systems*, 37:49081–49129, 2024. 3
- Ingvar Ziemann, Nikolai Matni, and George J Pappas. State space models, emergence, and ergodicity: How many parameters are needed for stable predictions? *arXiv preprint arXiv:2409.13421*, 2024. 3

A LLM USAGE DISCLOSURE

LLMs were used in elaborating this paper as follows:

- Finding related work.
- Computing the result of polynomial multiplications.
- Generating LaTeX tables and tikz figures.
- Transferring proofs from pen-and-paper format into LaTeX automatically using the online tool Manus <https://manus.im/>.

B EXPERIMENTS — FURTHER DETAILS

B.1 HYPERPARAMETERS

Table 1: Training hyperparameters of setting (a) in Section 5

Hyperparameter	Value
Weight initialization	Xavier normal distribution (Glorot & Bengio, 2010) with gain = $1e-6$
Optimizer	AdamW (Loshchilov & Hutter, 2017) with $\beta_1 = 0.9$ and $\beta_2 = 0.98$ for all AR(1–4), $\epsilon = 1e-9$
Weight decay	$5e-3$ for AR(2 & 4) and $1e-2$ for AR(1 & 3)
Learning rate (i.e., max. val.)	$2e-2$ for AR(1), $3e-2$ for AR(2), $3e-3$ for AR(3), $4e-3$ for AR(4)
Min. learning rate	$1e-4$ for AR(1 & 2) and $1e-5$ for AR(3 & 4)
Linear warmup	800 iter.
Decay schedule	Cosine annealing (Loshchilov & Hutter, 2016)
Max. decay steps	7200 iter.
Max. grad norm (clipping)	300
Random seeds	{666013, 1, 0}
Batch size / iter.	3000 for AR(1), 4000 for AR(2), 10000 for AR(3), 18000 for AR(4)
Total iter.	8001

Table 2: Training hyperparameters of setting (b) in Section 5

Hyperparameter	Value
Weight initialization	Xavier normal distribution (Glorot & Bengio, 2010) with gain = $1e-6$
Optimizer	AdamW (Loshchilov & Hutter, 2017) with $\beta_1 = 0.98$ for AR(1), 0.96 for AR(2), 0.92 for AR(3), 0.88 for AR(4), $\beta_2 = 0.99$ for AR(1), 0.98 for AR(2), 0.96 for AR(3), 0.92 for AR(4), $\epsilon = 1e-9$
Weight decay	$1e-2$ for AR(1), $1e-3$ for AR(2), $1e-4$ for AR(3 & 4)
Learning rate (i.e., max. val.)	$1e-3$ for AR(1), $4e-3$ for AR(2), $6e-3$ for AR(3), $8e-3$ for AR(4)
Min. learning rate	$1e-5$
Linear warmup	800 iter.
Decay schedule	Cosine annealing (Loshchilov & Hutter, 2016)
Max. decay steps	7200 iter.
Max. grad norm (clipping)	300
Random seeds	{666013, 1, 0}
Batch size / iter.	4000 for AR(1), 8000 for AR(2), 14000 for AR(3), 18000 for AR(4)
Total iter.	8001

Table 3: Training hyperparameters of setting (c) in Section 5

Hyperparameter	Value
Weight initialization	Xavier normal distribution (Glorot & Bengio, 2010) with gain = $1e-6$
Optimizer	AdamW (Loshchilov & Hutter, 2017) with $\beta_1 = 0.98$ for AR(1), 0.96 for AR(2), 0.92 for AR(3), $\beta_2 = 0.99$ for AR(1), 0.98 for AR(2), 0.96 for AR(3), $\epsilon = 1e-9$
Weight decay	$1e-2$ for AR(1), $5e-3$ for AR(2), $1e-3$ for AR(3)
Learning rate (i.e., max. val.)	$2e-3$ for AR(1), $3e-3$ for AR(2), $4e-3$ for AR(3)
Min. learning rate	$1e-5$
Linear warmup	800 iter.
Decay schedule	Cosine annealing (Loshchilov & Hutter, 2016)
Max. decay steps	7200 iter.
Max. grad norm (clipping)	300
Random seeds	{666013, 1, 0}
Batch size / iter.	4000 for AR(1), 8000 for AR(2), 16000 for AR(3)
Total iter.	8001

Table 4: Training hyperparameters of setting (d) in Section 5

Hyperparameter	Value
Weight initialization	Xavier normal distribution (Glorot & Bengio, 2010) with gain = $1e-6$
Optimizer	AdamW (Loshchilov & Hutter, 2017) with $\beta_1 = 0.92$ and $\beta_2 = 0.96$ for all AR(1 - 2), $\beta_1 = 0.88$ and $\beta_2 = 0.94$ for all AR(3), $\epsilon = 1e-9$
Weight decay	$5e-3$ for AR(1), $1e-3$ for AR(2), $1e-3$ for AR(3)
Learning rate (i.e., max. val.)	$5e-3$ for AR(1), $5e-4$ for AR(2), $1e-4$ for AR(3)
Min. learning rate	$1e-5$
Linear warmup	800 iter.
Decay schedule	Cosine annealing (Loshchilov & Hutter, 2016)
Max. decay steps	7200 iter.
Max. grad norm (clipping)	300
Random seeds	{666013, 1, 0}
Batch size / iter.	4000 for AR(1), 10000 for AR(2), 18000 for AR(3)
Total iter.	8001

B.2 EXPERIMENTS SHOWING CONVERGENCE TO THE CHECKERBOARD PATTERN DURING TRAINING

This set of experiments serves to illustrate that parameters \mathbf{W}_{QK} and \mathbf{W}_V converge to the checkerboard pattern across iterations. Since the non-zero values of these parameters are of different magnitudes and we do not have their theoretical expressions for window-sizes greater than 1, we shall only consider their non-zero support, as follows.

Definition B.1. For a matrix $\mathbf{M} \in \mathbb{R}^{d \times m}$, its support is defined as the collection of positions corresponding to non-zero values

$$\text{supp}(\mathbf{M}) := \{(i, j) \in [d] \times [m] \mid a_{i,j} \neq 0\}. \quad (19)$$

Additionally, the support-induced mask is a binary matrix with unit entries on the support

$$\text{mask}(\mathbf{M}) := \left[\mathbf{1}_{(i,j) \in \text{supp}(\mathbf{M})} \right]_{i,j=1}^{i=d, j=m} \quad (20)$$

where $\mathbf{1}_C = 1$ if condition C is true and 0 otherwise, is the indicator function centered at z .

We rely on the Jaccard distance (Jaccard, 1901) adapted to binary matrices \mathbf{A}, \mathbf{B}

$$d_{\text{Jac}}(\mathbf{A}, \mathbf{B}) := 1 - \frac{\sum_{i,j} a_{i,j} b_{i,j}}{\sum_{i,j} \max\{a_{i,j}, b_{i,j}\}} \quad (21)$$

to track whether the support-induced masks of our parameters during training converge to the predicted (for AR(1)) or hypothesized (for AR(s) $s \geq 2$) sparsity patterns of Lemma 4.1. Our experiments employ a tolerance level of $1e-1$ when computing the masks of \mathbf{W}_V and \mathbf{W}_{QK} , meaning that any entry below this value is considered zero. The results are depicted in Figure 2 and its subplots for varying window sizes, where $\mathbf{M}_{QK}^{\text{true}}$ and $\mathbf{M}_V^{\text{true}}$ represent the masks posited in Lemma 4.1 for a null tolerance level. The illustrations empirically confirm that our parameters' supports converge to the ones identified in Lemma 4.1.

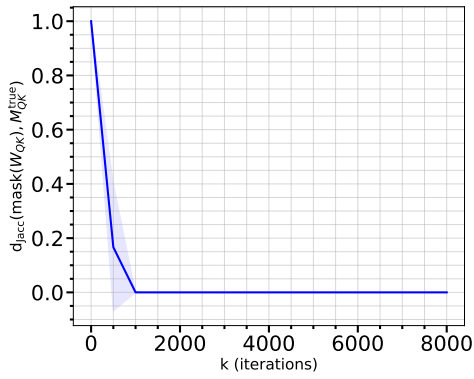
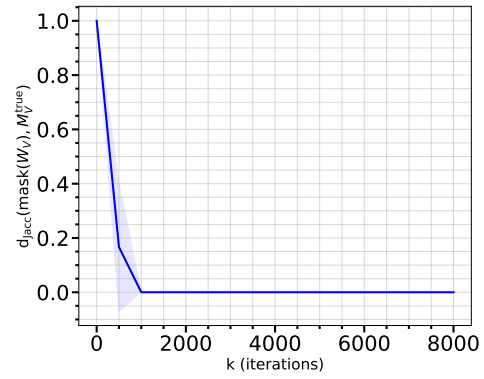
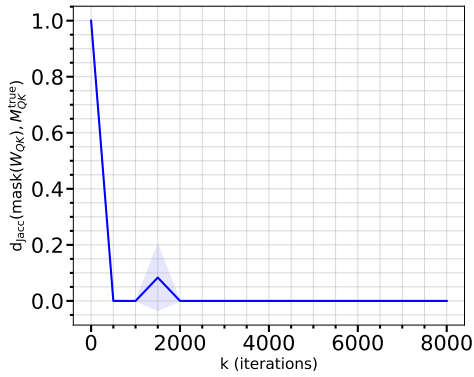
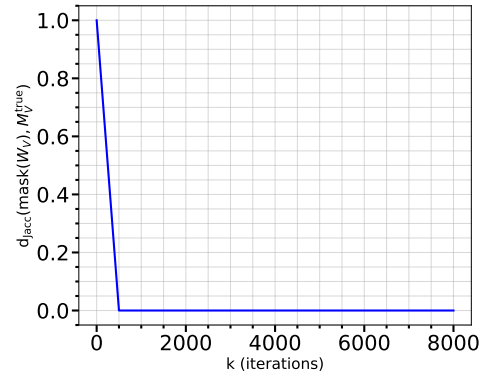
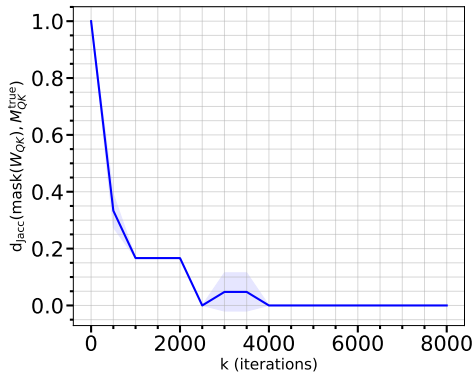
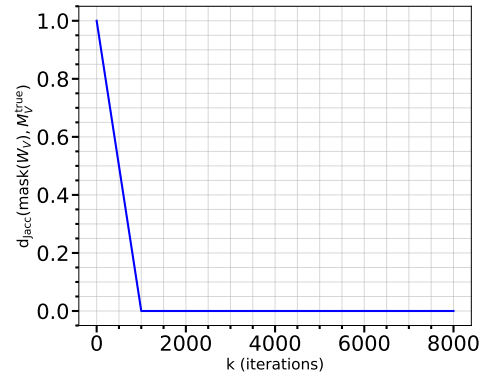
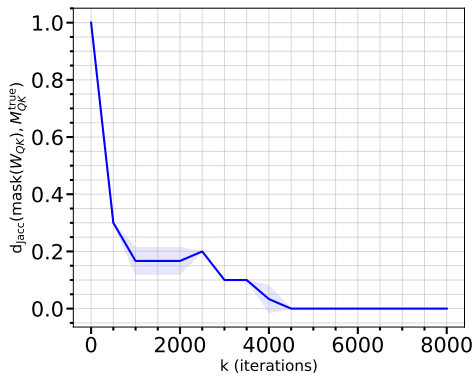
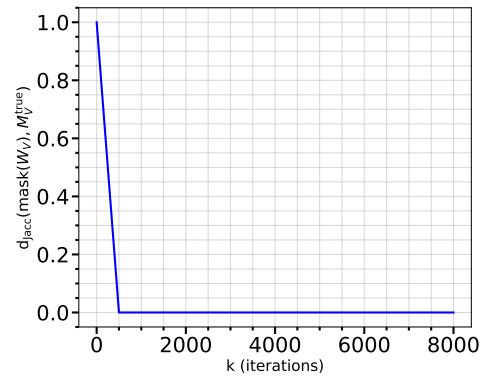
(a) The Jaccard distance of \mathbf{W}_{QK} for AR(1).(b) The Jaccard distance of \mathbf{W}_V for AR(1).(c) The Jaccard distance of \mathbf{W}_{QK} for AR(2).(d) The Jaccard distance of \mathbf{W}_V for AR(2).(e) The Jaccard distance of \mathbf{W}_{QK} for AR(3).(f) The Jaccard distance of \mathbf{W}_V for AR(3).(g) The Jaccard distance of \mathbf{W}_{QK} for AR(4).(h) The Jaccard distance of \mathbf{W}_V for AR(4).

Figure 2: The experiment results of the Jaccard distance between the M_{QK}^{true} and \mathbf{W}_{QK} and the Jaccard distance between the M_V^{true} and \mathbf{W}_V for AR(1–4). Both converge to 0 at the end of the training.

B.3 EXPERIMENTS WITH NON-DIAGONAL, SYMMETRIC \mathbf{A} , RANDOM \mathbf{c} AND ISOTROPIC Σ_w

The LDS which generates the training data is as follows. For each sequence, sample $\mathbf{v} \sim \text{Unif}([-1, 1]^d)$, sample $\mathbf{Q} \sim \text{Haar}(O(d))$ and set $\mathbf{A} = \mathbf{Q}^\top \text{diag}(\mathbf{v})\mathbf{Q}$. Sample $\mathbf{c} \sim \text{Unif}([-2, 2]^d)$. The noise covariances are set to $\Sigma_w = 1e-2\mathbf{I}$ and $\sigma_v^2 = 1e-2$.

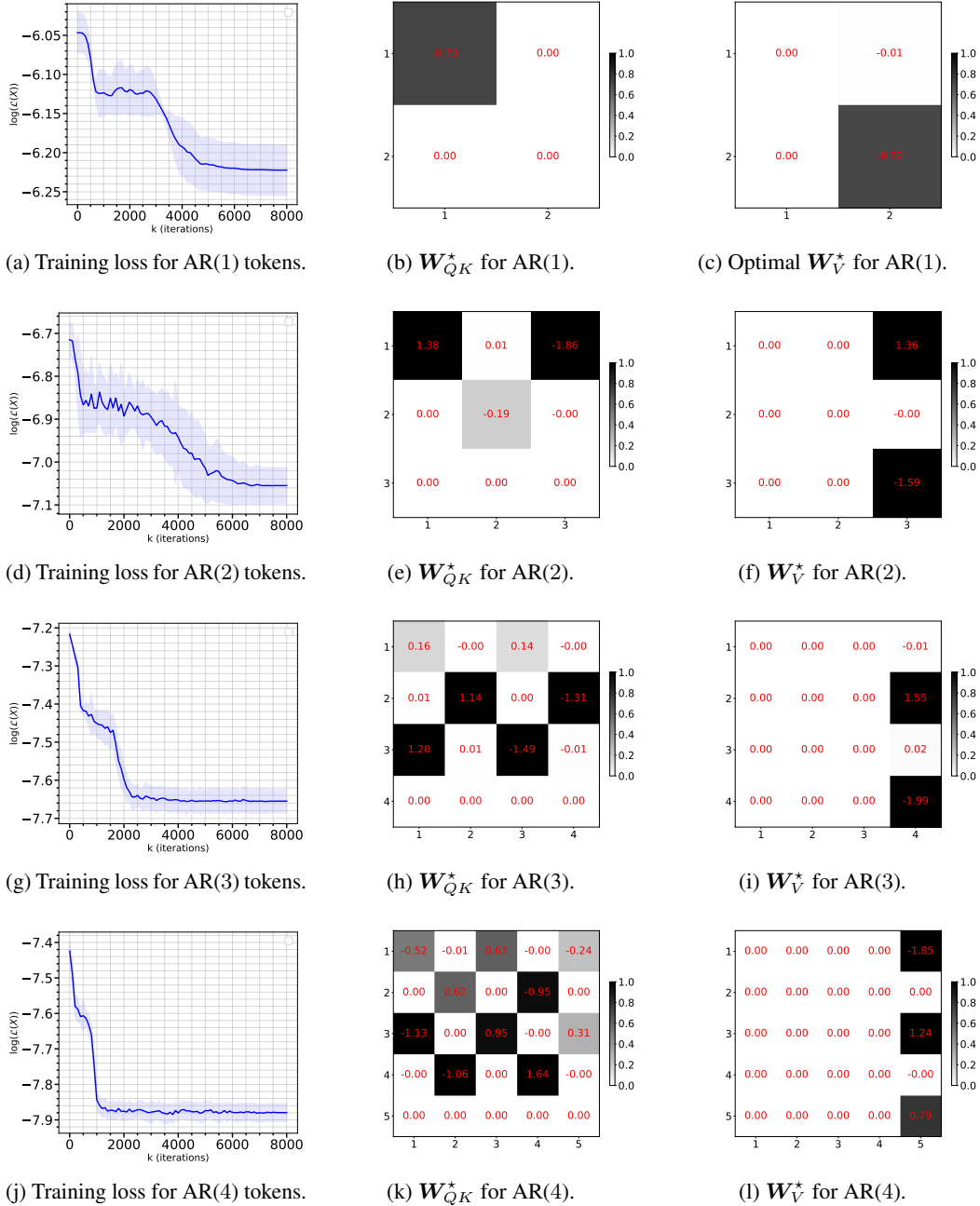


Figure 3: Experimental results for AR(1–4) with non-diagonal, symmetric \mathbf{A} , random \mathbf{c} and isotropic Σ_w , setting (b) in Section 5, which align with the Lemma 4.1.

B.4 EXPERIMENTS WITH NON-DIAGONAL, SYMMETRIC \mathbf{A} , RANDOM \mathbf{c} AND NON-DIAGONAL, ANISOTROPIC Σ_w

The LDS which generates the training data is as follows. For each sequence, sample $\mathbf{v} \sim \text{Unif}([-1, 1]^d)$; sample $\mathbf{Q} \sim \text{Haar}(O(d))$ and set $\mathbf{A} = \mathbf{Q}^\top \text{diag}(\mathbf{v})\mathbf{Q}$; sample $\mathbf{c} \sim \text{Unif}([-0.5, 0.5]^d)$. Set the process noise covariance $\Sigma_w = \mathbf{Q}_w^\top \text{diag}(1e-2 \cdot [0.9, 0.95, 1.0, 1.05, 1.1])\mathbf{Q}_w$, where $\mathbf{Q}_w \sim \text{Haar}(O(d))$. Set $\sigma_v^2 = 1e-2$.

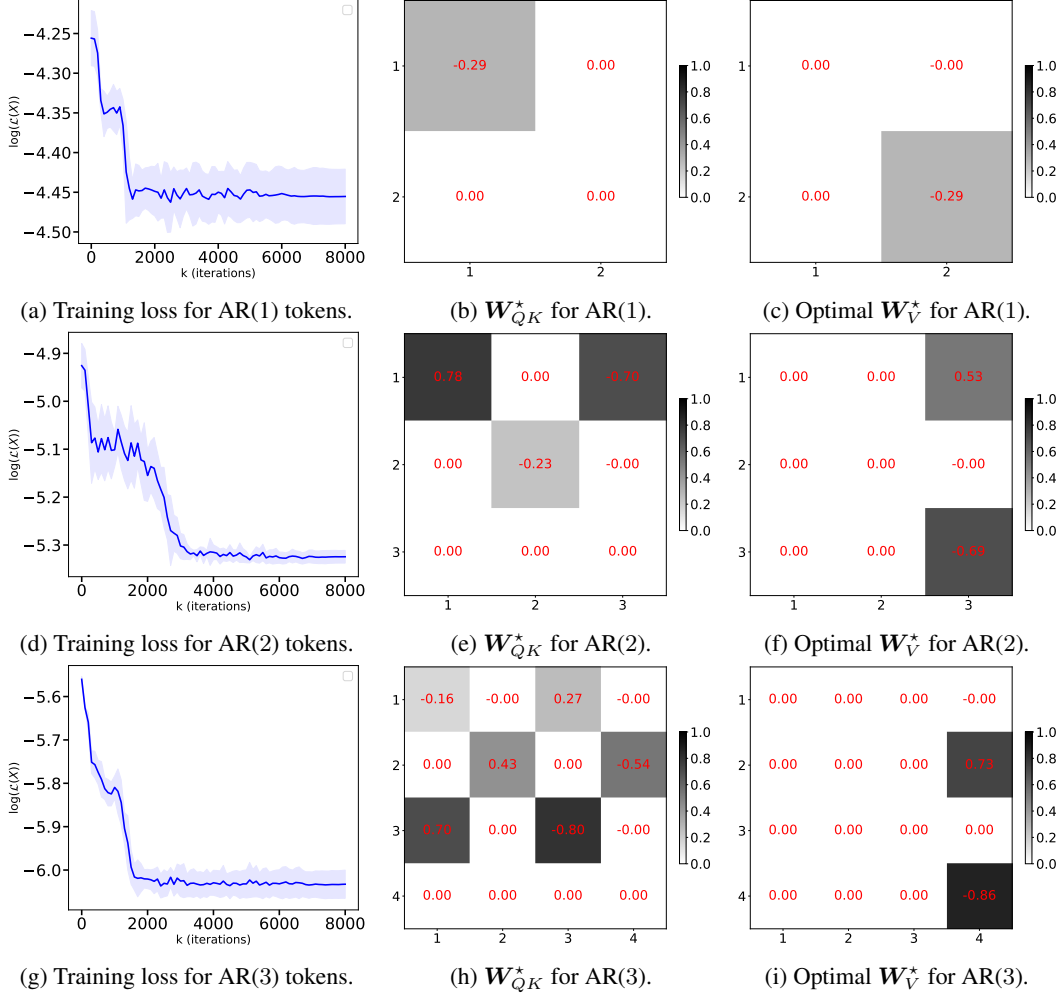


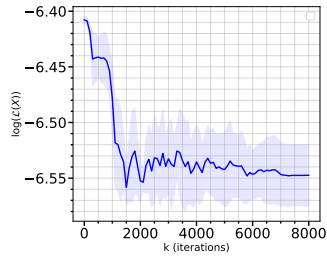
Figure 4: Experimental results for AR(1–3) with non-diagonal, symmetric \mathbf{A} , random \mathbf{c} and non-diagonal, anisotropic Σ_w , setting (c) in Section 5, which align with the Lemma 4.1.

B.5 EXPERIMENTS WITH NON-DIAGONAL, NON-SYMMETRIC \mathbf{A} , RANDOM \mathbf{c} AND ISOTROPIC Σ_w

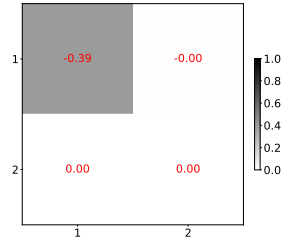
The LDS which generates the training data is as follows.

For each sequence, sample $\mathbf{d} \sim \text{Unif}([-1, 1]^d)$, sample $\mathbf{P} = [p_{i,j}]_{i,j=1}^d$ by sampling $p_{i,j}$ i.i.d. from $\mathcal{U}([-1, 1])$, and set $\mathbf{A} = \mathbf{P}^{-1} \text{diag}(\mathbf{d})\mathbf{P}$. Sample $\mathbf{c} \sim \text{Unif}([-1, 1]^d)$. The noise covariances are set to $\Sigma_w = 1e-2\mathbf{I}$ and $\sigma_v^2 = 1e-2$.

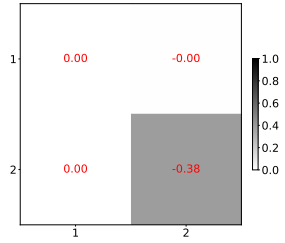
In practice, we need to guarantee \mathbf{P} is well conditioned. After sampling $p_{i,j}$ i.i.d. from $\mathcal{U}([-1, 1])$, we decompose $\mathbf{P} = \mathbf{Q}\mathbf{R}$, where \mathbf{Q} is an orthogonal matrix and \mathbf{R} is an upper-triangle matrix from the QR decomposition of \mathbf{P} . We modify the diagonals of \mathbf{R} manually to make sure $\frac{\max_i R_{ii}}{\min_i R_{ii}} = 2$ and right multiply \mathbf{Q} with the modified \mathbf{R} to have the well conditioned \mathbf{P} .



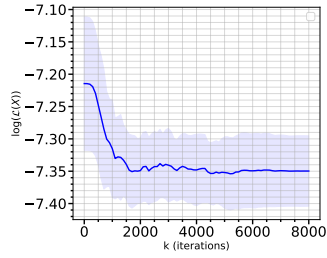
(a) Training loss for AR(1) tokens.



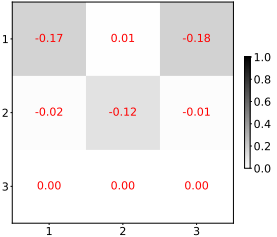
(b) \mathbf{W}_{QK}^* for AR(1).



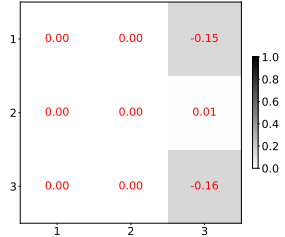
(c) Optimal \mathbf{W}_V^* for AR(1).



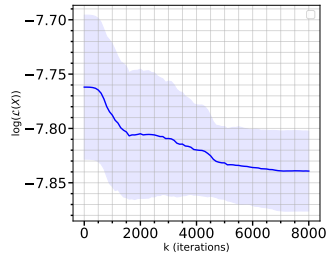
(d) Training loss for AR(2) tokens.



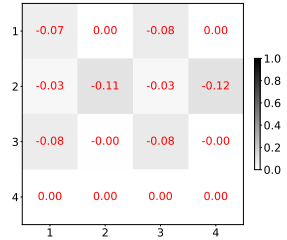
(e) \mathbf{W}_{QK}^* for AR(2).



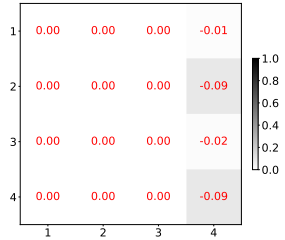
(f) \mathbf{W}_V^* for AR(2).



(g) Training loss for AR(3) tokens.



(h) \mathbf{W}_{QK}^* for AR(3).



(i) \mathbf{W}_V^* for AR(3).

Figure 5: Experimental results for AR(1–3) with non-diagonal, non-symmetric \mathbf{A} , random \mathbf{c} and isotropic Σ_w , setting (d) in Section 5, which align with the Lemma 4.1.

C SECTION 3 PROOFS

C.1 PROOF OF TOKEN CONSTRUCTION LEMMA

Lemma 3.1. *For a given $s \geq 1$, there exists an $s + 1$ -headed linear attention layer with positional encoding which transforms input sequences $[y_1, y_2, \dots, y_T]^\top$ into*

$$\begin{bmatrix} \bar{y}_1 & \dots & \bar{y}_{T-s} & \mathbf{0}_{(s-1) \times (T-s-1)} \\ y_{s+1} & \dots & 0 & \mathbf{0}_{T-s-1}^\top \end{bmatrix}^\top.$$

The latter quantity is essentially equivalent to \mathbf{Y}_0 as defined in equation (8).

Proof. We first define a matrix right-shift operator, which shifts each row one position to the right, padding the first column with zeros. Let $\gg: \mathbb{R}^{m \times n} \rightarrow \mathbb{R}^{m \times n}$ be $\gg(\mathbf{M}) = \mathbf{M}\mathbf{R}$, where

$$\mathbf{R} = \begin{bmatrix} 0 & \mathbf{0}_{n-1}^\top \\ \mathbf{0}_{n-1} & \mathbf{I}_{n-1} \end{bmatrix}. \quad (22)$$

We follow Von Oswald et al. (2023a) in using the one-hot positional encodings, concatenated to the input sequence to obtain tokens $\{[y_t, e_t]\}_{t=1}^T$. We define $s + 1$ attention heads given by

Define $\mathbf{W}_Q \in \mathbb{R}^{T+1 \times T}$, $\mathbf{W}_K \in \mathbb{R}^{T+1 \times T}$ and $\mathbf{W}_V \in \mathbb{R}^{T+1 \times s}$ as follows:

$$\begin{aligned} \mathbf{W}_Q^h &= \begin{bmatrix} \mathbf{0}_T^\top \\ \mathbf{I}_T \end{bmatrix}, \forall h \in [s+1] \\ (\mathbf{W}_K^h)^\top &= \left[\mathbf{0}_T, \underbrace{\gg(\dots \gg(\mathbf{I}_T) \dots)}_{h-1 \text{ times}} \right] \\ \mathbf{W}_V^h &= \begin{bmatrix} 1 & \dots & h & \dots & s+1 \\ \mathbf{0}_{T+1} & \dots & \mathbf{e}_1 & \dots & \mathbf{0}_{T+1} \end{bmatrix}, \forall h \in [s+1] \end{aligned} \quad (23)$$

Each head then computes the following

$$\begin{aligned} & \underbrace{\begin{bmatrix} y_1 & 1 & 0 & \dots & 0 \\ y_2 & 0 & 1 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ y_T & 0 & 0 & \dots & 1 \end{bmatrix}}_{=\mathbf{I}_T} \mathbf{W}_Q^k (\mathbf{W}_K^h)^\top \underbrace{\begin{bmatrix} y_1 & y_2 & y_3 & \dots & y_T \\ 1 & 0 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & 1 \end{bmatrix}}_{=\begin{bmatrix} \mathbf{0}_{T-h+1 \times h-1} & \mathbf{I}_{T-h+1} \\ \mathbf{0}_{h-1 \times h-1} & \mathbf{0}_{h-1 \times T-h+1} \end{bmatrix}} \mathbf{W}_V \underbrace{\begin{bmatrix} y_1 & 1 & 0 & \dots & 0 \\ y_2 & 0 & 1 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ y_T & 0 & 0 & \dots & 1 \end{bmatrix}}_{=\begin{bmatrix} 1 & \dots & h & \dots & s+1 \\ 0 & \dots & y_1 & \dots & 0 \\ 0 & \dots & y_2 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & \dots & y_T & \dots & 0 \end{bmatrix}} \\ &= \begin{bmatrix} 0 & \dots & y_h & \dots & 0 \\ 0 & \dots & y_{h+1} & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & \dots & y_T & \dots & 0 \\ & & \mathbf{0}_{h \times s+1} & & \end{bmatrix} \end{aligned}$$

Summing over the outputs of all heads, we get an equivalent representation to (8). \square

C.2 PROOF OF THE ALMOST SURE OBSERVABILITY OF THE LDS

We seek to show that Assumption 3.2 ensures LDS (1) observability w.p. 1. Note that the central symmetry of the distribution is irrelevant for this statement, and only relevant for the proofs in Section 4. We repeat Assumption 3.2 below for convenience.

Assumption 3.2 (LDS family). *The system matrix $\mathbf{A} \in \mathbb{R}^{d \times d}$ is sampled from a centrally symmetric distribution supported on $\{\mathbf{M} \in \mathbb{R}^{d \times d} \mid \rho(\mathbf{M}) \leq 1\}$, for which it holds that*

$$\mathbb{P}(\{\mathbf{A} \mid \exists i, j \in [d], \text{ s.t. } \lambda_i(\mathbf{A}) = \lambda_j(\mathbf{A})\}) = 0. \quad (9)$$

In other words, \mathbf{A} has a simple spectrum almost surely. The observation vector $\mathbf{c} \in \mathbb{R}^d$ is sampled independently, from a distribution that is absolutely continuous w.r.t. the Lebesgue measure over \mathbb{R}^d .

Lemma C.1. *Assumption 3.2 ensures the pair (\mathbf{A}, \mathbf{c}) is observable w.p. 1.*

Proof. Since \mathbf{A} has distinct eigenvalues w.p. 1 (the simple spectrum condition), it is (block) diagonalizable almost surely, and its eigenvectors $\{\mathbf{v}_1, \dots, \mathbf{v}_d\}$ are linearly independent. Therefore, observability is ensured if $\mathbf{c}^\top \mathbf{v}_i \neq 0$ almost surely for all $i \in [d]$.

Since \mathbf{c} is sampled from a distribution that is absolutely continuous w.r.t. the Lebesgue measure in \mathbb{R}^d , we want to prove that the set

$$\mathcal{U} = \bigcup_{i=1}^d \{\mathbf{c} \in \mathbb{R}^d \mid \mathbf{c}^\top \mathbf{v}_i = 0\}$$

has zero Lebesgue measure in the ambient \mathbb{R}^d . Each collection $\{\mathbf{c} \in \mathbb{R}^d \mid \mathbf{c}^\top \mathbf{v}_i = 0\}$ forms a proper subspace of \mathbb{R}^d with dimension at most $d - 1$ (it can be less, for complex \mathbf{v}_i). Therefore, its Lebesgue measure is null (see, e.g., (Royden & Fitzpatrick, 2010, pg. 435)).

Since \mathcal{U} is a finite union of measure zero sets, it is itself measure zero. Hence, observability holds w.p. 1. \square

D SECTION 4 PROOFS

D.1 PRELIMINARIES

Since we're dealing with data generated from stochastic processes, our proofs will heavily rely on taking expectations conditioned on randomness up to a certain point in the process. In what follows, we formalize the natural filtrations with respect to process (1).

We denote the natural filtration associated with (1). as $\{\mathcal{F}_t\}_{t \geq 0}$, where

$$\mathcal{F}_t := \sigma(\mathbf{A}, \mathbf{c}, \mathbf{x}_0, \mathbf{w}_0, \dots, \mathbf{w}_{t-1}, v_0, \dots, v_{t-1}), \quad t \geq 0. \quad (24)$$

By convention, when $t = 0$ the sets of noise variables are empty, and we define

$$\mathcal{F}_0 = \sigma(\mathbf{A}, \mathbf{c}, \mathbf{x}_0), \quad (25)$$

to illustrate that \mathbf{A} and \mathbf{c} are sampled once at time 0 and then remain fixed.

It follows that

- (a) $\mathcal{F}_t \subseteq \mathcal{F}_{t+1}, \forall t \geq 0$
- (b) \mathbf{x}_t is \mathcal{F}_t -measurable for all $t \geq 0$.
- (c) y_t is \mathcal{F}_{t+1} -measurable (since y_t depends on v_t)
- (d) The noise at time t is independent on the respective filtration: $\mathbf{w}_t \perp\!\!\!\perp \mathcal{F}_t, v_t \perp\!\!\!\perp \mathcal{F}_t$, for all $t \geq 0$.

D.2 AUXILIARY RESULTS AND TECHNICAL LEMMATA

Theorem D.1 (Isserlis (1918)). *Let $\mathbf{y} = [y_1, y_2, \dots, y_n]^\top \sim \mathcal{N}_n(0, \Sigma)$ be an n -dimensional, mean-zero multivariate normal vector. Then, for any even integer n ,*

$$\mathbb{E}[y_1 y_2 \cdots y_n] = \sum_{p \in \text{PP}(n)} \prod_{(\ell, r) \in p} \mathbb{E}[y_\ell y_r],$$

where $\text{PP}(n)$ denotes the set of all pairwise partitions of $[n]$ into disjoint pairs. If n is odd, then $\mathbb{E}[y_1 y_2 \cdots y_n] = 0$.

Lemma D.1. *Given random vectors $\mathbf{z}, \mathbf{w}, \mathbf{q} \in \mathbb{R}^d$ and assuming that \mathbf{w} is independent of \mathbf{z}, \mathbf{q} and the relevant integrability conditions hold, then*

$$\mathbb{E}[\mathbf{z}^\top \mathbf{w} \mathbf{w}^\top \mathbf{q}] = \mathbb{E}[\mathbf{z}^\top \mathbb{E}[\mathbf{w} \mathbf{w}^\top] \mathbf{q}] \quad (26)$$

Proof. We use the towering property of expectations,

$$\begin{aligned} \mathbb{E}[\mathbf{z}^\top \mathbf{w} \mathbf{w}^\top \mathbf{z}] &= \mathbb{E}[\mathbf{z}^\top \mathbb{E}[\mathbf{w} \mathbf{w}^\top | \mathbf{z}, \mathbf{q}] \mathbf{q}] \\ &= \mathbb{E}[\mathbf{z}^\top \mathbb{E}[\mathbf{w} \mathbf{w}^\top] \mathbf{q}], \end{aligned}$$

where the last line follows from the quantities' independence. \square

Lemma D.2. *Let the sequence $\{y_i\}_{i \geq 0}$ be generated by an LDS (1) sampled according to Assumption 3.2. For time indices $0 \leq i \leq j$, it holds that*

$$\mathbb{E}[y_i y_j] = \mathbb{E}[\mathbf{c}^\top \mathbf{A}^i \Sigma_{\mathbf{x}_0} (\mathbf{A}^\top)^j \mathbf{c}] + \sum_{k=0}^{i-1} \mathbb{E}[\mathbf{c}^\top \mathbf{A}^{i-1-k} \Sigma_{\mathbf{w}} (\mathbf{A}^\top)^{j-1-k} \mathbf{c}] + \mathbf{1}_{\{i=j\}} \sigma_v^2, \quad (27)$$

where $\mathbf{1}_{\{i=j\}}$ takes the value 1 if $i = j$ and 0 otherwise.

Proof. For process (1) it holds that

$$\mathbf{x}_j = \mathbf{A}^{j-i} \mathbf{x}_i + \sum_{k=i}^{j-1} \mathbf{A}^{j-1-k} \mathbf{w}_k$$

and therefore

$$y_j = \mathbf{c}^\top \mathbf{A}^{j-i} \mathbf{x}_i + \sum_{k=i}^{j-1} \mathbf{c}^\top \mathbf{A}^{j-1-k} \mathbf{w}_k + v_j.$$

The product of scalars $y_i y_j$ therefore takes the form

$$\begin{aligned} y_i y_j &= y_i y_j^\top \\ &= (\mathbf{c}^\top \mathbf{x}_i + v_i) \left(\mathbf{c}^\top \mathbf{A}^{j-i} \mathbf{x}_i + \sum_{k=i}^{j-1} \mathbf{c}^\top \mathbf{A}^{j-1-k} \mathbf{w}_k + v_j \right)^\top \\ &= \mathbf{c}^\top \mathbf{x}_i \mathbf{x}_i^\top (\mathbf{A}^\top)^{j-i} \mathbf{c} + \sum_{k=i}^{j-1} \mathbf{c}^\top \mathbf{x}_i \mathbf{w}_k^\top (\mathbf{A}^\top)^{j-1-k} \mathbf{c} + \mathbf{c}^\top \mathbf{x}_i v_j \\ &\quad + v_i \mathbf{x}_i^\top (\mathbf{A}^\top)^{j-i} \mathbf{c} + \sum_{k=i}^{j-1} v_i \mathbf{w}_k^\top (\mathbf{A}^\top)^{j-1-k} \mathbf{c} + v_i v_j. \end{aligned} \quad (28)$$

Now, observing that $\mathbb{E}[y_i y_j] = \mathbb{E}[\mathbb{E}[y_i y_j | \mathcal{F}_i]]$ and remembering that $\mathbf{x}_i, \mathbf{A}, \mathbf{c}$ are \mathcal{F}_i -measurable, and that for all i and $p \geq i$, $\mathbf{w}_p \perp\!\!\!\perp \mathcal{F}_i$ and $v_p \perp\!\!\!\perp \mathcal{F}_i$, and $\mathbf{w}_p \perp\!\!\!\perp v_q, \forall p, q \geq 0$, we have

$$\begin{aligned} \mathbb{E}[\mathbf{c}^\top \mathbf{x}_i \mathbf{x}_i^\top (\mathbf{A}^\top)^{j-i} \mathbf{c} | \mathcal{F}_i] &= \mathbf{c}^\top \mathbf{x}_i \mathbf{x}_i^\top (\mathbf{A}^\top)^{j-i} \mathbf{c}, \\ \mathbb{E}\left[\sum_{k=i}^{j-1} \mathbf{c}^\top \mathbf{x}_i \mathbf{w}_k^\top (\mathbf{A}^\top)^{j-1-k} \mathbf{c} \middle| \mathcal{F}_i\right] &= \sum_{k=i}^{j-1} \mathbf{c}^\top \mathbf{x}_i \mathbb{E}[\mathbf{w}_k^\top] (\mathbf{A}^\top)^{j-1-k} \mathbf{c} = 0, \\ \mathbb{E}[\mathbf{c}^\top \mathbf{x}_i v_j | \mathcal{F}_i] &= \mathbf{c}^\top \mathbf{x}_i \mathbb{E}[v_j] = 0, \\ \mathbb{E}[v_i \mathbf{x}_i^\top (\mathbf{A}^\top)^{j-i} \mathbf{c} | \mathcal{F}_i] &= \mathbb{E}[v_i] \mathbf{x}_i^\top (\mathbf{A}^\top)^{j-i} \mathbf{c} = 0, \\ \mathbb{E}\left[\sum_{k=i}^{j-1} v_i \mathbf{w}_k^\top (\mathbf{A}^\top)^{j-1-k} \mathbf{c} \middle| \mathcal{F}_i\right] &= \sum_{k=i}^{j-1} \mathbb{E}[v_i] \mathbb{E}[\mathbf{w}_k^\top] (\mathbf{A}^\top)^{j-1-k} \mathbf{c} = 0, \\ \mathbb{E}[v_i v_j | \mathcal{F}_i] &= \mathbb{E}[v_i v_j] = \mathbf{1}_{\{i=j\}} \sigma_v^2. \end{aligned}$$

Therefore,

$$\mathbb{E}[y_i y_j] = \mathbb{E}[\mathbb{E}[y_i y_j | \mathcal{F}_i]] = \mathbb{E}[\mathbf{c}^\top \mathbf{x}_i \mathbf{x}_i^\top (\mathbf{A}^\top)^{j-i} \mathbf{c}] + \mathbf{1}_{\{i=j\}} \sigma_v^2. \quad (29)$$

Noting that $\mathbf{x}_i = \mathbf{A}^i \mathbf{x}_0 + \sum_{k=0}^{i-1} \mathbf{A}^{i-1-k} \mathbf{w}_k$, we further unroll the first term inside the expectation in (29) and get

$$\begin{aligned} \mathbf{c}^\top \mathbf{x}_i \mathbf{x}_i^\top (\mathbf{A}^\top)^{j-i} \mathbf{c} &= \left[\mathbf{c}^\top \mathbf{A}^i \mathbf{x}_0 + \sum_{k=0}^{i-1} \mathbf{c}^\top \mathbf{A}^{i-1-k} \mathbf{w}_k \right] \left[\mathbf{x}_0^\top (\mathbf{A}^\top)^j \mathbf{c} + \sum_{k=0}^{i-1} \mathbf{w}_k^\top (\mathbf{A}^\top)^{j-1-k} \mathbf{c} \right] \\ &= \mathbf{c}^\top \mathbf{A}^i \mathbf{x}_0 \mathbf{x}_0^\top (\mathbf{A}^\top)^j \mathbf{c} + \sum_{k=0}^{i-1} \mathbf{c}^\top \mathbf{A}^i \mathbf{x}_0 \mathbf{w}_k^\top (\mathbf{A}^\top)^{j-1-k} \mathbf{c} \\ &\quad + \sum_{k=0}^{i-1} \mathbf{c}^\top \mathbf{A}^{i-1-k} \mathbf{w}_k \mathbf{x}_0^\top (\mathbf{A}^\top)^j \mathbf{c} + \sum_{k,l=0}^{i-1} \mathbf{c}^\top \mathbf{A}^{i-1-k} \mathbf{w}_k \mathbf{w}_l^\top (\mathbf{A}^\top)^{j-1-l} \mathbf{c}. \end{aligned} \quad (30)$$

Using $\mathbf{w}_p \perp\!\!\!\perp \mathcal{F}_0 \subset \mathcal{F}_i, \forall p \geq 0$ and $\mathbf{w}_p \perp\!\!\!\perp \mathbf{w}_q, \forall p \neq q$ in conjunction with (30) and Lemma D.1 we get

$$\mathbb{E} [\mathbf{c}^\top \mathbf{x}_i \mathbf{x}_i^\top (\mathbf{A}^\top)^{j-i} \mathbf{c} \mid \mathcal{F}_0] = \mathbf{c}^\top \mathbf{A}^i \mathbf{x}_0 \mathbf{x}_0^\top (\mathbf{A}^\top)^j \mathbf{c} + \sum_{k=0}^{i-1} \mathbf{c}^\top \mathbf{A}^{i-1-k} \Sigma_{\mathbf{w}} (\mathbf{A}^\top)^{j-1-k} \mathbf{c} \quad (31)$$

Furthermore, noting that $\sigma(\mathbf{A}, \mathbf{c}) \subset \mathcal{F}_0$, we have that

$$\mathbb{E} [\mathbf{c}^\top \mathbf{x}_i \mathbf{x}_i^\top (\mathbf{A}^\top)^{j-i} \mathbf{c} \mid \mathbf{A}, \mathbf{c}] = \mathbf{c}^\top \mathbf{A}^i \Sigma_{\mathbf{x}_0} (\mathbf{A}^\top)^j \mathbf{c} + \sum_{k=0}^{i-1} \mathbf{c}^\top \mathbf{A}^{i-1-k} \Sigma_{\mathbf{w}} (\mathbf{A}^\top)^{j-1-k} \mathbf{c}. \quad (32)$$

Taking full expectation in (32), and plugging everything back into (29), we get the stated result

$$\mathbb{E} [y_i y_j] = \mathbb{E} [\mathbf{c}^\top \mathbf{A}^i \Sigma_{\mathbf{x}_0} (\mathbf{A}^\top)^j \mathbf{c}] + \sum_{k=0}^{i-1} \mathbb{E} [\mathbf{c}^\top \mathbf{A}^{i-1-k} \Sigma_{\mathbf{w}} (\mathbf{A}^\top)^{j-1-k} \mathbf{c}] + \mathbf{1}_{\{i=j\}} \sigma_v^2. \quad \square$$

Lemma D.3. Let $\{y_i\}_{i \geq 0}$ be a sequence of observations generated by an LDS (1) sampled according to Assumption 3.2. Then,

- (a) if $i + j = 2p + 1$ for some $p \in \mathbb{N}_+$, $\mathbb{E} [y_i y_j] = 0$;
- (b) if $i + j + k + l = 2p + 1$ for some $p \in \mathbb{N}_+$, $\mathbb{E} [y_i y_j y_k y_l] = 0$;
- (c) if $i + j + k + l + m + n = 2p + 1$ for some $p \in \mathbb{N}_+$, $\mathbb{E} [y_i y_j y_k y_l y_m y_n] = 0$.

Note that there is no condition on the indices being pairwise distinct.

Proof. To prove point (a), we start from the expression derived in Lemma D.2.

$$\mathbb{E} [y_i y_j] = \mathbb{E} [\mathbf{c}^\top \mathbf{A}^i \Sigma_{\mathbf{x}_0} (\mathbf{A}^\top)^j \mathbf{c}] + \sum_{k=0}^{i-1} \mathbb{E} [\mathbf{c}^\top \mathbf{A}^{i-1-k} \Sigma_{\mathbf{w}} (\mathbf{A}^\top)^{j-1-k} \mathbf{c}] + \mathbf{1}_{\{i=j\}} \sigma_v^2$$

Clearly, since $i + j$ is odd, it holds that $i \neq j$ and hence the third term is zero. Furthermore, since \mathbf{A} has a centrally symmetric distribution, we have that

$$\begin{aligned} \mathbb{E} [\mathbf{c}^\top \mathbf{A}^i \Sigma_{\mathbf{x}_0} (\mathbf{A}^\top)^j \mathbf{c}] &= \mathbb{E} [\mathbf{c}^\top (-\mathbf{A})^i \Sigma_{\mathbf{x}_0} (-\mathbf{A}^\top)^j \mathbf{c}] \\ &= (-1)^{i+j} \mathbb{E} [\mathbf{c}^\top \mathbf{A}^i \Sigma_{\mathbf{x}_0} (\mathbf{A}^\top)^j \mathbf{c}], \end{aligned} \quad (33)$$

implying that $\mathbb{E} [\mathbf{c}^\top \mathbf{A}^i \Sigma_{\mathbf{x}_0} (\mathbf{A}^\top)^j \mathbf{c}] = 0$. We apply a similar reasoning for the other term and obtain that

$$\mathbb{E} [y_i y_j] = 0,$$

thus proving the first point.

For both points (b) and (c), we will rely on Isserlis's theorem, which we replicate in Theorem D.1 for convenience. Note that conditioned, on \mathbf{A} and \mathbf{c} , the vectors $[y_i y_j y_k y_l \mid \mathbf{A}, \mathbf{c}]$ and $[y_i y_j y_k y_l y_m y_n \mid \mathbf{A}, \mathbf{c}]$ are jointly Gaussian since they are linear transformations of the jointly Gaussian vectors $\mathbf{r}_1 = [\mathbf{x}_0^\top, \mathbf{w}_0^\top, \dots, \mathbf{w}_{\max\{i,j,k,l\}}^\top, v_0, \dots, v_{\max\{i,j,k,l\}}]^\top$ and $\mathbf{r}_2 = [\mathbf{x}_0^\top, \mathbf{w}_0^\top, \dots, \mathbf{w}_{\max\{i,j,k,l,m,n\}}^\top, v_0, \dots, v_{\max\{i,j,k,l,m,n\}}]^\top$, respectively. We can therefore apply the towering property along with Isserlis's result to get

$$\begin{aligned} \mathbb{E} [y_i y_j y_k y_l] &= \mathbb{E} [\mathbb{E} [y_i y_j y_k y_l \mid \mathbf{A}, \mathbf{c}]] \\ &= \mathbb{E} \left[\mathbb{E} [y_i y_j \mid \mathbf{A}, \mathbf{c}] \mathbb{E} [y_k y_l \mid \mathbf{A}, \mathbf{c}] + \mathbb{E} [y_i y_k \mid \mathbf{A}, \mathbf{c}] \mathbb{E} [y_j y_l \mid \mathbf{A}, \mathbf{c}] \right. \\ &\quad \left. + \mathbb{E} [y_i y_l \mid \mathbf{A}, \mathbf{c}] \mathbb{E} [y_j y_k \mid \mathbf{A}, \mathbf{c}] \right], \end{aligned} \quad (34)$$

since $\text{PP}(\{i, j, k, l\}) = \{\{(i, j), (k, l)\}, \{(i, k), (j, l)\}, \{(i, l), (j, k)\}\}$. Since $i + j + k + l$ is odd, the two pairs inside any $q \in \text{PP}(\{i, j, k, l\})$ must have different parities (i.e., one even, one odd). W.l.o.g, we analyze the first term in (34), assuming $0 \leq i \leq j \leq k \leq l$. From (32), we know that

$$\begin{aligned} \mathbb{E}[y_i y_j | \mathbf{A}, \mathbf{c}] \mathbb{E}[y_k y_l | \mathbf{A}, \mathbf{c}] &= \left[\mathbf{c}^\top \mathbf{A}^i \Sigma_{\mathbf{x}_0} (\mathbf{A}^\top)^j \mathbf{c} + \sum_{t=0}^{i-1} \mathbf{c}^\top \mathbf{A}^{i-1-t} \Sigma_{\mathbf{w}} (\mathbf{A}^\top)^{j-1-t} \mathbf{c} + \mathbf{1}_{\{i=j\}} \sigma_v^2 \right] \\ &\quad \left[\mathbf{c}^\top \mathbf{A}^k \Sigma_{\mathbf{x}_0} (\mathbf{A}^\top)^l \mathbf{c} + \sum_{t=0}^{k-1} \mathbf{c}^\top \mathbf{A}^{k-1-t} \Sigma_{\mathbf{w}} (\mathbf{A}^\top)^{l-1-t} \mathbf{c} + \mathbf{1}_{\{k=l\}} \sigma_v^2 \right] \end{aligned} \quad (35)$$

Assume w.l.o.g that $i + j$ is even, and $k + l$ is odd. This implies that $\mathbf{1}_{\{k=l\}} = 0$. Taking full expectation on both sides and developing the product, we get

$$\begin{aligned} &\mathbb{E}[\mathbb{E}[y_i y_j | \mathbf{A}, \mathbf{c}] \mathbb{E}[y_k y_l | \mathbf{A}, \mathbf{c}]] \\ &= \mathbb{E}[\mathbf{c}^\top \mathbf{A}^i \Sigma_{\mathbf{x}_0} (\mathbf{A}^\top)^j \mathbf{c} \mathbf{c}^\top \mathbf{A}^k \Sigma_{\mathbf{x}_0} (\mathbf{A}^\top)^l \mathbf{c}] \\ &\quad + \sum_{t=0}^{k-1} \mathbb{E}[\mathbf{c}^\top \mathbf{A}^i \Sigma_{\mathbf{x}_0} (\mathbf{A}^\top)^j \mathbf{c} \mathbf{c}^\top \mathbf{A}^{k-1-t} \Sigma_{\mathbf{w}} (\mathbf{A}^\top)^{l-1-t} \mathbf{c}] \\ &\quad + \sum_{t=0}^{i-1} \mathbb{E}[\mathbf{c}^\top \mathbf{A}^{i-1-t} \Sigma_{\mathbf{w}} (\mathbf{A}^\top)^{j-1-t} \mathbf{c} \mathbf{c}^\top \mathbf{A}^k \Sigma_{\mathbf{x}_0} (\mathbf{A}^\top)^l \mathbf{c}] \\ &\quad + \sum_{t=0}^{i-1} \sum_{s=0}^{k-1} \mathbb{E}[\mathbf{c}^\top \mathbf{A}^{i-1-t} \Sigma_{\mathbf{w}} (\mathbf{A}^\top)^{j-1-t} \mathbf{c} \mathbf{c}^\top \mathbf{A}^{k-1-s} \Sigma_{\mathbf{w}} (\mathbf{A}^\top)^{l-1-s} \mathbf{c}] \\ &\quad + \mathbf{1}_{\{i=j\}} \sigma_v^2 \mathbb{E}[\mathbf{c}^\top \mathbf{A}^k \Sigma_{\mathbf{x}_0} (\mathbf{A}^\top)^l \mathbf{c}] \\ &\quad + \mathbf{1}_{\{i=j\}} \sigma_v^2 \sum_{t=0}^{k-1} \mathbb{E}[\mathbf{c}^\top \mathbf{A}^{k-1-t} \Sigma_{\mathbf{w}} (\mathbf{A}^\top)^{l-1-t} \mathbf{c}] \end{aligned} \quad (36)$$

Using the index parity assumptions and the reasoning based on the central symmetry of \mathbf{A}' 's distribution from (33), we get that all the terms on the RHS of (36) are zero. We treat the remaining terms in (34) similarly to get the final result in (b).

Finally, point (c) follows a similar path. We have

$$\begin{aligned} \text{PP}(\{i, j, k, l, m, n\}) &= \{\{(i, j), (k, l), (m, n)\}, \{(i, j), (k, m), (l, n)\}, \{(i, j), (k, n), (l, m)\}, \\ &\quad \{(i, k), (j, l), (m, n)\}, \{(i, k), (j, m), (l, n)\}, \{(i, k), (j, n), (l, m)\}, \\ &\quad \{(i, l), (j, k), (m, n)\}, \{(i, l), (j, m), (k, n)\}, \{(i, l), (j, n), (k, m)\}, \\ &\quad \{(i, m), (j, k), (l, n)\}, \{(i, m), (j, l), (k, n)\}, \{(i, m), (j, n), (k, l)\}, \\ &\quad \{(i, n), (j, k), (l, m)\}, \{(i, n), (j, l), (k, m)\}, \{(i, n), (j, m), (k, l)\}\}. \end{aligned}$$

For the parity hypothesis to be satisfied, not that inside a set $q \in \text{PP}(\{i, j, k, l, m, n\})$, at least one pair must have an odd parity, while the other two must be of the same parity (either even or odd). W.o.l.g let $0 \leq i \leq j \leq k \leq l \leq m \leq n$, pick the first set in $\text{PP}(\{i, j, k, l, m, n\})$ above (the rest follow the same logic) and assume that $m + n$ is odd. By the same logic as before, we have that $\mathbf{1}_{\{m=n\}} = 0$ and

$$\begin{aligned} &\mathbb{E}[\mathbb{E}[y_i y_j | \mathbf{A}, \mathbf{c}] \mathbb{E}[y_k y_l | \mathbf{A}, \mathbf{c}] \mathbb{E}[y_m y_n | \mathbf{A}, \mathbf{c}]] \\ &= \mathbb{E} \left[\left[\mathbf{c}^\top \mathbf{A}^i \Sigma_{\mathbf{x}_0} (\mathbf{A}^\top)^j \mathbf{c} + \sum_{t=0}^{i-1} \mathbf{c}^\top \mathbf{A}^{i-1-t} \Sigma_{\mathbf{w}} (\mathbf{A}^\top)^{j-1-t} \mathbf{c} + \mathbf{1}_{\{i=j\}} \sigma_v^2 \right] \right. \\ &\quad \left[\mathbf{c}^\top \mathbf{A}^k \Sigma_{\mathbf{x}_0} (\mathbf{A}^\top)^l \mathbf{c} + \sum_{t=0}^{k-1} \mathbf{c}^\top \mathbf{A}^{k-1-t} \Sigma_{\mathbf{w}} (\mathbf{A}^\top)^{l-1-t} \mathbf{c} + \mathbf{1}_{\{k=l\}} \sigma_v^2 \right] \\ &\quad \left. \left[\mathbf{c}^\top \mathbf{A}^m \Sigma_{\mathbf{x}_0} (\mathbf{A}^\top)^n \mathbf{c} + \sum_{t=0}^{m-1} \mathbf{c}^\top \mathbf{A}^{m-1-t} \Sigma_{\mathbf{w}} (\mathbf{A}^\top)^{n-1-t} \mathbf{c} \right] \right] \end{aligned}$$

Without computing, one can see that every term in the expanded product will have powers of \mathbf{A} whose sum is odd. Therefore, using the centrally symmetric property of \mathbf{A} 's distribution, all the terms evaluate to zero, and point (c) is proven. \square

D.3 PROOF OF LEMMA 4.1

Lemma 4.1. *For any $s \geq 1$, $j \in [s]$, the parameters \mathbf{W}_{QK} and \mathbf{W}_V having the structure*

$$\mathbf{W}_{QK} = \begin{bmatrix} \mathbf{R}_{1:s, 1:(s+1)} \\ \mathbf{0}_{s+1}^\top \end{bmatrix}, \quad \mathbf{W}_V = [\mathbf{0}_{(s+1) \times s} \quad \mathbf{r}_{2+s \bmod 2 : 2 \lceil \frac{s+1}{2} \rceil}], \quad (14)$$

with $\mathbf{R} := \mathbf{1}_{\lfloor \frac{s+1}{2} \rfloor \times \lceil \frac{s+1}{2} \rceil} \otimes \begin{bmatrix} \star & 0 \\ 0 & \star \end{bmatrix}$ and $\mathbf{r} := \mathbf{1}_{\lceil \frac{s+1}{2} \rceil} \otimes \begin{bmatrix} 0 \\ \star \end{bmatrix}$, ensure that LHS (12) $_{r,\ell} = 0$ whenever $r + l \in 2\mathbb{N}$ and $s + j \in 2\mathbb{Z}$, or $r + l \in 2\mathbb{N} + 1$ and $s + j \in 2\mathbb{Z} + 1$.

Proof. Recall the in-context loss in (10) with a general AR(s)-constructed input token matrix

$$\mathbf{Y}_0 = \begin{bmatrix} \bar{\mathbf{y}}_1 & \bar{\mathbf{y}}_2 & \cdots & \bar{\mathbf{y}}_{T-s-1} & \bar{\mathbf{y}}_{T-s} \\ y_{s+1} & y_{s+2} & \cdots & y_{T-1} & 0 \end{bmatrix}$$
 is defined as

$$\mathcal{L}(\theta) := \mathbb{E} \left[\left(\mathcal{T}_\theta(\mathbf{Y}_0)_{s+1, T-s} - y_T \right)^2 \right]. \quad (37)$$

For equations (38) to (42) below, we use the same reformulations as Ahn et al. (2023). The last column of the transformer's output above can be written as

$$\begin{bmatrix} \bar{\mathbf{y}}_{T-1} \\ 0 \end{bmatrix} = \begin{bmatrix} \bar{\mathbf{y}}_{T-1} \\ 0 \end{bmatrix} + \frac{1}{T-s-1} \mathbf{W}_V^\top \left(\sum_{i=1}^{T-s-1} \begin{bmatrix} \bar{\mathbf{y}}_i \bar{\mathbf{y}}_i^\top & \bar{\mathbf{y}}_i y_{i+s} \\ \bar{\mathbf{y}}_i^\top y_{i+s} & y_{i+s}^2 \end{bmatrix} \right) \mathbf{W}_{QK}^\top \begin{bmatrix} \bar{\mathbf{y}}_{T-s} \\ 0 \end{bmatrix}, \quad (38)$$

where the summation comes from the causal mask. Therefore, the transformer's prediction of y_T , $\mathcal{T}_\theta(\mathbf{Y}_0)_{s+1, T-s}$ can be written as

$$\frac{1}{T-s-1} \mathbf{b}^\top \left(\underbrace{\sum_{i=1}^{T-s-1} \begin{bmatrix} \bar{\mathbf{y}}_i \bar{\mathbf{y}}_i^\top & \bar{\mathbf{y}}_i y_{i+s} \\ \bar{\mathbf{y}}_i^\top y_{i+s} & y_{i+s}^2 \end{bmatrix}}_{:= \mathbf{Y} \in \mathbb{R}^{(s+1) \times (s+1)}} \right) [\mathbf{z}_1 \mathbf{z}_2 \cdots \mathbf{z}_s] \bar{\mathbf{y}}_{T-s}, \quad (39)$$

where $\mathbf{b}^\top \in \mathbb{R}^{1 \times (s+1)}$ is the last row of \mathbf{W}_V^\top and $\mathbf{z}_j \in \mathbb{R}^{(s+1)}$ is the j^{th} column of \mathbf{W}_{QK}^\top . So the in-context loss $\mathcal{L}(\mathbf{W}_V, \mathbf{W}_{QK})$ can be rewritten as a function of \mathbf{b}^\top and $\mathbf{Z}\mathbf{Z} = [\mathbf{z}_j]_{j=1}^s$

$$\mathcal{L}(\mathbf{b}, \mathbf{Z}\mathbf{Z}) := \mathbb{E} \left[\left(\frac{1}{T-s-1} \mathbf{b}^\top \bar{\mathbf{Y}} \mathbf{Z}\mathbf{Z} \bar{\mathbf{y}}_{T-s} - y_T \right)^2 \right]. \quad (40)$$

Plugging in the expression of $\bar{\mathbf{y}}_{T-s}$, the in-context loss is

$$\mathcal{L}(\mathbf{b}, \mathbf{Z}\mathbf{Z}) = \mathbb{E} \left[\left(\frac{1}{T-s-1} \mathbf{b}^\top \bar{\mathbf{Y}} [\mathbf{z}_1 \mathbf{z}_2 \cdots \mathbf{z}_s] \begin{bmatrix} y_{T-s} \\ y_{T-s+1} \\ \vdots \\ y_{T-1} \end{bmatrix} - y_T \right)^2 \right]$$

$$\begin{aligned}
&= \mathbb{E} \left[\left(\frac{1}{T-s-1} \sum_{k=1}^s \mathbf{b}^\top \bar{\mathbf{Y}} \mathbf{z}_k y_{T-s-1+k} - y_T \right)^2 \right] \\
&= \mathbb{E} \left[\left(\frac{1}{T-s-1} \sum_{k=1}^s \text{Tr}(\bar{\mathbf{Y}} \mathbf{z}_k \mathbf{b}^\top) y_{T-s-1+k} - y_T \right)^2 \right] \\
&= \mathbb{E} \left[\left(\frac{1}{T-s-1} \sum_{k=1}^s \langle \bar{\mathbf{Y}}, \mathbf{b} \mathbf{z}_k^\top \rangle y_{T-s-1+k} - y_T \right)^2 \right]. \tag{41}
\end{aligned}$$

We reparametrize the in-context loss using $\mathbf{X}_k := \mathbf{b} \mathbf{z}_k^\top$

$$\mathcal{L}(\mathbf{X}_{k \in [s]}) = \mathbb{E} \left[\left(\frac{1}{T-s-1} \sum_{k=1}^s \langle \bar{\mathbf{Y}}, \mathbf{X}_k \rangle y_{T-s-1+k} - y_T \right)^2 \right]. \tag{42}$$

Note that the gradient of the in-context loss with respect to each \mathbf{X}_j is

$$\nabla_{\mathbf{X}_j} \mathcal{L}(\mathbf{X}_{k \in [s]}) = 2\mathbb{E} \left[\left(\frac{1}{T-s-1} \sum_{k=1}^s \langle \bar{\mathbf{Y}}, \mathbf{X}_k \rangle y_{T-s-1+k} - y_T \right) y_{T-s-1+j} \bar{\mathbf{Y}} \right]. \tag{43}$$

The gradient $\nabla_{\mathbf{X}_j} \mathcal{L}(\mathbf{X}_{k \in [s]})$ is a sum of two terms, $\nabla_{\mathbf{X}_j} \mathcal{L}(\mathbf{X}_{k=1 \dots s}) = \mathbf{T}_{\mathbf{X}_j}^{(1)} + \mathbf{T}_{\mathbf{X}_j}^{(2)}$, where, replacing $\bar{\mathbf{Y}}$ we have

$$\mathbf{T}_{\mathbf{X}_j}^{(1)} := \frac{2}{T-s-1} \mathbb{E} \left[\sum_{k=1}^s \langle \bar{\mathbf{Y}}, \mathbf{X}_k \rangle y_{T-s-1+k} y_{T-s-1+j} \sum_{i=1}^{T-s-1} \begin{bmatrix} \bar{\mathbf{y}}_i \bar{\mathbf{y}}_i^\top & \bar{\mathbf{y}}_i y_{i+s} \\ \bar{\mathbf{y}}_i^\top y_{i+s} & y_{i+s}^2 \end{bmatrix} \right] \tag{44}$$

$$\mathbf{T}_{\mathbf{X}_j}^{(2)} := -2\mathbb{E} \left[y_T y_{T-s-1+j} \sum_{i=1}^{T-s-1} \begin{bmatrix} \bar{\mathbf{y}}_i \bar{\mathbf{y}}_i^\top & \bar{\mathbf{y}}_i y_{i+s} \\ \bar{\mathbf{y}}_i^\top y_{i+s} & y_{i+s}^2 \end{bmatrix} \right]. \tag{45}$$

Each matrix element of $\mathbf{T}_{\mathbf{X}_j}^{(2)}$ has the form

$$\sum_{i=1}^{T-s-1} 2\mathbb{E} [y_T y_{T-s-1+j} y_{i+m} y_{i+n}] \tag{46}$$

with $j \in [1, s]$, $m \in [0, s]$ and $n \in [0, s]$.

The sum of y 's indices in (46) for each term in the above sum is $2T + 2i + (m + n - s - 1 + j)$. The parity is determined by that of $m + n - s - 1 + j$ and is independent of the sum counter i (i.e., the same for all terms). According to Lemma D.3, (46) is 0 if $(m + n - s - 1 + j)$ is odd, and of arbitrary value if it is even. So a general matrix element of $\mathbf{T}_{\mathbf{X}_j}^{(2)}$ is 0 if $(m + n - s - 1 + j)$ is odd and of arbitrary value if $(m + n - s - 1 + j)$ is even.

For a given AR(s)-type token (s is fixed) and a specific j , whether a matrix element of $\mathbf{T}_{\mathbf{X}_j}^{(2)}$ is 0 only depends on $m + n$ (its position in the matrix). So,

$$\mathbf{T}_{\mathbf{X}_j}^2 = \begin{cases} \begin{bmatrix} \star & 0 & \star & \cdots & \cdots \\ 0 & \star & 0 & \star & \\ \star & 0 & \star & \ddots & \ddots \\ \vdots & \star & \ddots & \ddots & \ddots \\ \vdots & & \ddots & \ddots & \star & 0 \\ \cdots & \cdots & \star & 0 & \star \end{bmatrix}, & \text{if } |j - s - 1| \text{ is even;} \\ \begin{bmatrix} 0 & \star & 0 & \cdots & \cdots \\ \star & 0 & \star & 0 & \\ 0 & \star & 0 & \ddots & \ddots \\ \vdots & 0 & \ddots & \ddots & \ddots \\ \vdots & & \ddots & \ddots & 0 & \star \\ \cdots & \cdots & 0 & \star & 0 \end{bmatrix}, & \text{if } |j - s - 1| \text{ is odd.} \end{cases}$$

We now turn to $\mathbf{T}_{\mathbf{X}_j}^{(1)}$ with the end goal of finding a parameter configuration that matches the sparsity pattern of $\mathbf{T}_{\mathbf{X}_j}^{(2)}$. For this section, assume s is odd (the other case follows similarly). First, let

$\mathbf{X}_k := \left[x_{i,j}^{(k)} \right]_{i,j=1}^{s+1}$ and unfold the expression of the matrix inner product

$$\langle \bar{\mathbf{Y}}, \mathbf{X}_k \rangle = \sum_{r=0}^s \sum_{l=0}^s \sum_{p=0}^{T-s-1} x_{l+1,r+1}^{(k)} y_{p+r} y_{p+l}, \quad (47)$$

where r, l are the indices traversing $\bar{\mathbf{Y}}$.

Furthermore, each matrix element of $\mathbf{T}_{\mathbf{X}_j}^{(1)}$ inside the expectation has the form

$$\frac{2}{T-s-1} \sum_{i=1}^{T-s-1} \sum_{k=1}^s \langle \bar{\mathbf{Y}}, \mathbf{X}_k \rangle y_{T-s-1+k} y_{T-s-1+j} y_{i+n} y_{i+m}, \quad (48)$$

where $n, m \in \{0, 1, \dots, s\}$ are the indices traversing $\bar{\mathbf{Y}}$.

In what follows, we'll use the notion of a matrix's support, as follows.

Definition D.1 (Matrix support). *Consider matrix $\mathbf{B} \in \mathbb{R}^{m \times n}$. The support of \mathbf{B} is the set of index pairs corresponding to its nonzero entries,*

$$\text{supp}(\mathbf{B}) := \{ (r, \ell) \in [m] \times [n] \mid \mathbf{B}_{r\ell} \neq 0 \}. \quad (49)$$

Assume that j is fixed and odd (we discuss the even case afterwards). Note that the sparsity of each position in $\mathbf{T}_{\mathbf{X}_j}^{(2)}$ dictated by the parity of $(m + n - s - 1 + j)$ where, when s, j -odd, the respective element is zero whenever $m + n$ is even. Notice that except for the contribution of the matrix inner product, the sum of indices for the y -factors in (48) is $2(T - s - 1 + i) + k + j + n + m$ so the parity is determined by that of $k + j + n + m$. We distinguish two cases:

- (a) when k is even, $k + j$ is odd, and we wish that the term zeroes out for even $m + n$. This means that \mathbf{X}_k must select in (47) only pairs $y_{p+r} y_{p+l}$ for which $r + \ell$ is even and zero-out the others. In other words, the support of \mathbf{X}_k should satisfy

$$\text{supp}(\mathbf{X}_k) \subseteq \{ (r, \ell) \in [s+1] \times [s+1] \mid r + \ell \in 2\mathbb{Z} \}. \quad (50)$$

Such an \mathbf{X}_k may look like

$$\mathbf{X}_k = \begin{bmatrix} x_{11}^{(k)} & 0 & x_{13}^{(k)} & \cdots & x_{1,s}^{(k)} & 0 \\ 0 & x_{22}^{(k)} & 0 & \cdots & 0 & x_{2,s+1}^{(k)} \\ x_{31}^{(k)} & 0 & x_{33}^{(k)} & \cdots & x_{3,s}^{(k)} & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ x_{s,1}^{(k)} & 0 & x_{s,3}^{(k)} & \cdots & x_{s,s}^{(k)} & 0 \\ 0 & x_{s+1,2}^{(k)} & 0 & \cdots & 0 & x_{s+1,s+1}^{(k)} \end{bmatrix}, \quad (51)$$

with arbitrary (possibly also zero) values for constants $x_{i,j}^{(k)}$. Note that, as a consequence, the entries corresponding to odd $m+n$ are not forced to zero by the distributional symmetry of \mathbf{A} and can generally take arbitrary values.

- (b) when k is odd, $k+j$ is even, and we wish that the term zeroes out for even $m+n$. This means that \mathbf{X}_k must select in (47) only pairs $y_{p+r}y_{p+\ell}$ for which $r+\ell$ is odd and zero-out the others. In other words, the support of \mathbf{X}_k should satisfy

$$\text{supp}(\mathbf{X}_k) \subseteq \{(r, \ell) \in [s+1] \times [s+1] \mid r+\ell \in 2\mathbb{Z}+1\}. \quad (52)$$

Such an \mathbf{X}_k may look like

$$\mathbf{X}_k = \begin{bmatrix} 0 & x_{12}^{(k)} & 0 & \cdots & 0 & x_{1,s+1}^{(k)} \\ x_{21}^{(k)} & 0 & x_{23}^{(k)} & \cdots & x_{2,s}^{(k)} & 0 \\ 0 & x_{32}^{(k)} & 0 & \cdots & 0 & x_{3,s+1}^{(k)} \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & x_{s,2}^{(k)} & 0 & \cdots & 0 & x_{s,s+1}^{(k)} \\ x_{s+1,1}^{(k)} & 0 & x_{s+1,3}^{(k)} & \cdots & x_{s+1,s}^{(k)} & 0 \end{bmatrix}, \quad (53)$$

with arbitrary (possibly also zero) values for constants $x_{i,j}^{(k)}$. Note that, as a consequence, the entries corresponding to odd $m+n$ are not forced to zero by the distributional symmetry of \mathbf{A} and can generally take arbitrary values.

These patterns need to be coherent with the case of j -even. Note that in $\mathbf{T}_{\mathbf{X}_j}^{(2)}$, when s -odd, j -even, the respective element is zero whenever $m+n$ is odd. We again distinguish two cases:

- (a) when k is even, $k+j$ is even, and we wish that the term zeroes out for odd $m+n$. This means that \mathbf{X}_k must select in (47) only pairs $y_{p+r}y_{p+\ell}$ for which $r+\ell$ is even and zero-out the others. Notice that the pattern of \mathbf{X}_k in (51) for even k satisfies this requirement and we have coherence.
- (b) when k is odd, $k+j$ is odd, and we wish that the term zeroes out for odd $m+n$. This means that \mathbf{X}_k must select in (47) only pairs $y_{p+r}y_{p+\ell}$ for which $r+\ell$ is odd and zero-out the others. Notice that the pattern of \mathbf{X}_k in (53) for odd k satisfies this requirement and we have coherence.

The same approach goes through for even window size s . To sum up, weights satisfying equations (50) and (52) for the appropriate parity of triplet (s, j, k) ensure the symmetry-induced sparsity pattern complies with that of $\mathbf{T}_{\mathbf{X}_j}^{(2)}$.

Finally, recall that $\mathbf{X}_k := \mathbf{b}z_k^\top$. It can be easily shown that vectors \mathbf{b} and z_k^\top yielding an \mathbf{X}_k whose support has maximum cardinality are those whose support consists of indices of opposing parities (e.g., even indices for \mathbf{b} and odd indices for z_k , or vice versa) in the case of (52) or the same parities in the case of (50). For our case of odd window sizes s , the sparsity pattern of \mathbf{b} and z_k^\top yielding a

complying \mathbf{X}_k is

$$\mathbf{b} = \begin{bmatrix} 0 \\ b_2 \\ \vdots \\ 0 \\ b_{s+1} \end{bmatrix} \quad \mathbf{z}_k^\top = \begin{cases} [0, z_2^{(k)}, \dots, 0, z_{s+1}^{(k)}], & \text{if } k \text{ is even} \\ [z_1^{(k)}, 0, \dots, z_s^{(k)}, 0], & \text{if } k \text{ is odd} \end{cases} \quad (54)$$

For even window size s , the patterns are

$$\mathbf{b} = \begin{bmatrix} b_1 \\ 0 \\ b_2 \\ \vdots \\ 0 \\ b_{s+1} \end{bmatrix} \quad \mathbf{z}_k^\top = \begin{cases} [0, z_2^{(k)}, \dots, 0, z_s^{(k)}, 0], & \text{if } k \text{ is even} \\ [z_1^{(k)}, 0, \dots, z_{s-1}^{(k)}, 0, z_{s+1}^{(k)}], & \text{if } k \text{ is odd} \end{cases} \quad (55)$$

Arranging these vectors inside \mathbf{W}_{QK} and \mathbf{W}_V gives the stated result. \square

D.4 PROOF OF THEOREM 4.1

Theorem 4.1. *Let \mathbf{Y}_0 encode the input tokens for $s = 1$. Then, the optimal parameters $\theta^* = (\mathbf{W}_{QK}^*, \mathbf{W}_V^*)$ of a single linear self-attention layer with respect to loss $\mathcal{L}(\theta)$ are*

$$\mathbf{W}_{QK}^* = \begin{bmatrix} \frac{(T-2)\mathbb{E}[y_{T-1}y_T \sum_{i=1}^{T-2} y_i y_{i+1}]}{\mathbb{E}[y_{T-1}^2 (\sum_{i=1}^{T-2} y_i y_{i+1})^2]} & 0 \\ 0 & 0 \end{bmatrix}, \quad \mathbf{W}_V^* = \begin{bmatrix} 0 & 0 \\ 0 & 1 \end{bmatrix}, \quad (15)$$

up to rescaling with a nonzero constant.

Proof. For the transformer parameters in (15), the corresponding $\mathbf{b}^\top = [0 \ 1]$ and the corresponding $\mathbf{F} = [c \ 0]$, where $c := \frac{(T-2)\mathbb{E}_{\mathcal{D}}[y_{T-1}y_T \sum_{i=1}^{T-2} y_i y_{i+1}]}{\mathbb{E}_{\mathcal{D}}[y_{T-1}^2 (\sum_{i=1}^{T-2} y_i y_{i+1})^2]}$.

So $\mathbf{X} = \mathbf{X}_1 = \mathbf{b}\mathbf{f}_1^\top = \mathbf{b}\mathbf{F}^\top = \begin{bmatrix} 0 & 0 \\ c & 0 \end{bmatrix}$. The gradient of the in-context loss $\nabla_{\mathbf{X}} \mathcal{L}(\mathbf{X})$ is

$$\begin{aligned} \mathbf{T}_{\mathbf{X}_j}^{(1)} &= \frac{2}{T-2} \mathbb{E}_{\mathcal{D}} [\langle \bar{\mathbf{Y}}, \mathbf{X} \rangle y_{T-1}^2 \bar{\mathbf{Y}}] \\ &= \frac{2}{T-2} \mathbb{E}_{\mathcal{D}} \left[\left\langle \sum_{r=1}^{T-2} \begin{bmatrix} y_r^2 & y_r y_{r+1} \\ y_{r+1} y_r & y_{r+1}^2 \end{bmatrix}, \begin{bmatrix} 0 & 0 \\ c & 0 \end{bmatrix} \right\rangle y_{T-1}^2 \sum_{i=1}^{T-2} \begin{bmatrix} y_i^2 & y_i y_{i+1} \\ y_{i+1} y_i & y_{i+1}^2 \end{bmatrix} \right] \\ &= \frac{2}{T-2} \mathbb{E}_{\mathcal{D}} \left[c \sum_{r=1}^{T-2} y_r y_{r+1} y_{T-1}^2 \sum_{i=1}^{T-2} \begin{bmatrix} y_i^2 & y_i y_{i+1} \\ y_{i+1} y_i & y_{i+1}^2 \end{bmatrix} \right] \\ &= \frac{2}{T-2} \mathbb{E}_{\mathcal{D}} \left[c \sum_{r=1}^{T-2} y_r y_{r+1} y_{T-1}^2 \sum_{i=1}^{T-2} \begin{bmatrix} 0 & y_i y_{i+1} \\ y_{i+1} y_i & 0 \end{bmatrix} \right]. \end{aligned} \quad (56)$$

According to Lemma D.3, the two diagonal elements in (56) $\mathbb{E}_{\mathcal{D}} \left[c \sum_{r=1}^{T-2} y_r y_{r+1} y_{T-1}^2 \sum_{i=1}^{T-2} y_i^2 \right]$ and $\mathbb{E}_{\mathcal{D}} \left[c \sum_{r=1}^{T-2} y_r y_{r+1} y_{T-1}^2 \sum_{i=1}^{T-2} y_{i+1}^2 \right]$ are 0, since the sums of y indices are both odd.

$$\begin{aligned} \mathbf{T}_{\mathbf{X}_j}^{(2)} &= -2 \mathbb{E}_{\mathcal{D}} \left[y_T y_{T-1} \sum_{i=1}^{T-2} \bar{\mathbf{Y}} \right] \\ &= -2 \mathbb{E}_{\mathcal{D}} \left[y_T y_{T-1} \sum_{i=1}^{T-2} \begin{bmatrix} y_i^2 & y_i y_{i+1} \\ y_{i+1} y_i & y_{i+1}^2 \end{bmatrix} \right] \end{aligned}$$

$$= -2\mathbb{E}_{\tilde{\mathcal{D}}}\left[y_T y_{T-1} \sum_{i=1}^{T-2} \begin{bmatrix} 0 & y_i y_{i+1} \\ y_{i+1} y_i & 0 \end{bmatrix}\right]. \quad (57)$$

According to Lemma D.3, the two diagonal elements in (57) $\mathbb{E}_{\tilde{\mathcal{D}}}\left[y_T y_{T-1} \sum_{i=1}^{T-2} y_i^2\right]$ and $\mathbb{E}_{\tilde{\mathcal{D}}}\left[y_T y_{T-1} \sum_{i=1}^{T-2} y_{i+1}^2\right]$ are 0, since the sums of y indices are both odd.

Plugging in the expression of c , it can be easily found that

$$\nabla_{\mathbf{X}} \mathcal{L}(\mathbf{X}) = \mathbf{T}_{\mathbf{X}_j}^1 + \mathbf{T}_{\mathbf{X}_j}^2 = 0. \quad (58)$$

Since the in-context loss is convex in \mathbf{X} and the \mathbf{X} resulting from the $\mathbf{W}_{\tilde{V}}^*$ and $\mathbf{W}_{\tilde{Q}K}^*$ above makes $\nabla_{\mathbf{X}} \mathcal{L}(\mathbf{X}) = 0$, the $\mathbf{W}_{\tilde{V}}^*$ and $\mathbf{W}_{\tilde{Q}K}^*$ above is a global minimizer for the in-context loss. \square

E PROOFS FOR SECTION 5

E.1 PROOF THAT OUR EXPERIMENTS' SAMPLING SCHEMES OBEY ASSUMPTION 3.2

All our experiments use a sampling schemes whose generalization is the following:

- (a) \mathbf{A} constructed by sampling $\mathbf{v} \sim \mathcal{P}$, where \mathcal{P} is centrally symmetric and absolutely continuous w.r.t. the Lebesgue measure on \mathbb{R}^d with marginals supported on $[-1, 1]$, and independently sampling \mathbf{P} , whose every entry is drawn i.i.d. from any absolutely continuous distribution w.r.t. Lebesgue measure in \mathbb{R} . Matrix \mathbf{A} is then formed as $\mathbf{P}\text{diag}(\mathbf{v})\mathbf{P}^{-1}$.
- (b) \mathbf{c} is sampled from an absolutely continuous distribution w.r.t. Lebesgue measure in \mathbb{R}^d , or otherwise fixed with $\mathbf{c} \neq \mathbf{0}_d$.

We need to show that

- (a) \mathbf{A} 's distribution is centrally symmetric, i.e., that $\mathbf{A} \stackrel{d}{=} -\mathbf{A}$;
- (b) \mathbf{A} 's spectrum is simple w.p. 1;
- (c) observability still holds when \mathbf{c} is fixed according to the above condition.

The first point is achieved since, by the central symmetry of \mathbf{v} 's distribution,

$$-\mathbf{A} = -\mathbf{P}^{-1}\text{diag}(\mathbf{v})\mathbf{P} = -\mathbf{P}^{-1}\text{diag}(-\mathbf{v})\mathbf{P} \stackrel{d}{=} \mathbf{P}^{-1}\text{diag}(\mathbf{v})\mathbf{P} = \mathbf{A}. \quad (59)$$

The second point is ensured by \mathbf{v} 's distribution being absolutely continuous w.r.t. the Lebesgue measure in \mathbb{R}^d , and hence the probability of \mathbf{v} belonging to $(d-1)$ -dimensional subspaces (and lower) such as $\{\mathbf{x} \in \mathbb{R}^d \mid \exists i, j \in [d] \text{ s.t. } x_i = x_j\}$ is null. In conjunction with the above, when we sample \mathbf{c} from a continuous distribution in \mathbb{R}^d , Assumption (3.2) is satisfied.

However, our proofs and experiments go through even if \mathbf{c} is fixed, as follows. First, the theoretical results rest on \mathbf{A} 's distributional symmetry and are invariant to the linear transformation induced by \mathbf{c} . Second, observability is ensured since $\det(\mathbf{O})$ in expression (5) is not zero w.p. 1, as follows.

We use $\det(\mathbf{OP}) \neq 0$ w.p. 1 $\iff \det(\mathbf{O}) \neq 0$ w.p. 1, since $\det(\mathbf{P}) \neq 0$ w.p. 1.

$$\det(\mathbf{OP}) \stackrel{z:=\mathbf{c}^\top \mathbf{P}}{=} \det([\mathbf{z}; \text{diag}(\mathbf{v})\mathbf{z}; \dots \text{diag}(\mathbf{v})^{d-1}\mathbf{z}]) \quad (60)$$

$$= \det(\text{diag}(\mathbf{z})) \det \left(\begin{bmatrix} 1 & v_1 & \dots & v_1^{d-1} \\ 1 & v_2 & \dots & v_2^{d-1} \\ \dots & \dots & \dots & \dots \\ 1 & v_d & \dots & v_d^{d-1} \end{bmatrix} \right). \quad (61)$$

Since \mathbf{P} 's entries are drawn i.i.d. from an absolutely continuous distribution w.r.t. Lebesgue measure in \mathbb{R} , it holds that $z_i \neq 0$ w.p. 1. The remaining matrix is Vandermonde with $v_i \neq v_j, \forall i, j \in [d]$ w.p. 1. Hence, the determinant is nonzero w.p. 1 and observability holds almost surely.

E.2 RELATION OF TRANSFORMER FORWARD PASS WITH PCG

For convenience, we reproduce below the PCG iteration of Shewchuk et al. (1994) for minimizing an objective

$$f(\mathbf{w}) = \frac{1}{2} \mathbf{w}^\top \mathbf{A} \mathbf{w} + \mathbf{b}^\top \mathbf{w} + c$$

We compute the first two steps of the algorithm with respect to the loss (4), which can be rewritten as

$$\mathcal{L}_{AR(s)}(\mathbf{w}) := \frac{1}{2(T-s-1)} \sum_{t=1}^{T-s-1} (y_{t+s} - \mathbf{w}^\top \bar{\mathbf{y}}_t)^2 \quad (62)$$

Algorithm 1 Preconditioned Conjugate Gradient

```

1: Input: preconditioner  $\mathbf{H}$ ,  $\mathbf{w}_0$ ,  $\mathbf{r}_0 = \mathbf{b} - \mathbf{A}\mathbf{w}_0$ ,  $\mathbf{d}_0 = \mathbf{H}^{-1}\mathbf{r}_0$ ,  $\delta_{\text{new}} = \mathbf{r}_0^\top \mathbf{d}_0$ 
2: for  $i = 0, 1, \dots$  do
3:    $\mathbf{z}_i = \mathbf{A}\mathbf{d}_i$ 
4:    $\alpha_i = \frac{\delta_{\text{new}}}{\mathbf{d}_i^\top \mathbf{z}_i}$ 
5:    $\mathbf{w}_{i+1} = \mathbf{w}_i + \alpha_i \mathbf{d}_i$ 
6:    $\mathbf{r}_{i+1} = \mathbf{r}_i - \alpha_i \mathbf{z}_i$ 
7:    $\mathbf{v}_{i+1} = \mathbf{H}^{-1}\mathbf{r}_{i+1}$ 
8:    $\delta_{\text{old}} = \delta_{\text{new}}, \delta_{\text{new}} = \mathbf{r}_{i+1}^\top \mathbf{v}_{i+1}$ 
9:    $\beta_{i+1} = \frac{\delta_{\text{new}}}{\delta_{\text{old}}}$ 
10:   $\mathbf{d}_{i+1} = \mathbf{v}_{i+1} + \beta_{i+1} \mathbf{d}_i$ 
11: end for
12: return  $\theta_T$ 

```

$$= \frac{1}{2(T-s-1)} \sum_{t=1}^{T-s-1} \mathbf{w}^\top \bar{\mathbf{y}}_t \bar{\mathbf{y}}_t^\top \mathbf{w} - 2y_{t+s} \mathbf{w}^\top \bar{\mathbf{y}}_t + y_{t+s}^2 \quad (63)$$

$$= \frac{1}{2} \mathbf{w}^\top \nabla^2 \mathcal{L}_{AR(s)} \mathbf{w} - \mathbf{w}^\top \nabla \mathcal{L}_{AR(s)}(0) + y_{t+s}^2 \quad (64)$$

Two iterations of Algorithm 1 starting from $\mathbf{w}_0 = \mathbf{0}$ and using $\mathbf{H} = \mathbf{P}^{-1}$ yield the following predictor.

$$\begin{aligned}
\mathbf{w}_1 &= \alpha_0 \mathbf{d}_0 = \alpha_0 \mathbf{H}^{-1} \mathbf{r}_0 \\
\mathbf{w}_2 &= \mathbf{w}_1 + \alpha_1 \mathbf{d}_1 \\
&= \alpha_0 \mathbf{d}_0 + \alpha_1 [\mathbf{P}\mathbf{r}_1 + \beta_1 \mathbf{d}_0] \\
&= \alpha_0 \mathbf{P} \nabla \mathcal{L}_{AR(s)}(\mathbf{0}) + \alpha_1 [\mathbf{P}(\mathbf{r}_0 - \alpha_0 \mathbf{z}_0) + \beta_1 \mathbf{P} \nabla \mathcal{L}_{AR(s)}(\mathbf{0})] \\
&= \alpha_0 \mathbf{P} \nabla \mathcal{L}_{AR(s)}(\mathbf{0}) + \alpha_1 [\mathbf{P}(\nabla \mathcal{L}_{AR(s)}(0) - \alpha_0 \nabla^2 \mathcal{L}_{AR(s)} \mathbf{d}_0) + \beta_1 \mathbf{P} \nabla \mathcal{L}_{AR(s)}(\mathbf{0})] \\
&= \alpha_0 \mathbf{P} \nabla \mathcal{L}_{AR(s)}(\mathbf{0}) + \alpha_1 [\mathbf{P}(\nabla \mathcal{L}_{AR(s)}(0) - \alpha_0 \nabla^2 \mathcal{L}_{AR(s)} \mathbf{P} \nabla \mathcal{L}_{AR(s)}(0)) + \beta_1 \mathbf{P} \nabla \mathcal{L}_{AR(s)}(\mathbf{0})] \\
&= (\alpha_0 + \alpha_1 + \alpha_1 \beta_1) \mathbf{P} \nabla \mathcal{L}_{AR(s)}(\mathbf{0}) - \alpha_1 \alpha_0 \mathbf{P} \nabla^2 \mathcal{L}_{AR(s)} \mathbf{P} \nabla \mathcal{L}_{AR(s)}(0)
\end{aligned}$$