Trained Mamba Emulates Online Gradient Descent in In-Context Linear Regression

Jiarui Jiang*¹, Wei Huang*², Miao Zhang[†] 1, Taiji Suzuki³ 2, Liqiang Nie¹

Harbin Institute of Technology, Shenzhen

²RIKEN AIP

³University of Tokyo

Abstract

State-space models (SSMs), particularly Mamba, emerge as an efficient Transformer alternative with linear complexity for long-sequence modeling. Recent empirical works demonstrate Mamba's in-context learning (ICL) capabilities competitive with Transformers, a critical capacity for large foundation models. However, theoretical understanding of Mamba's ICL remains limited, restricting deeper insights into its underlying mechanisms. Even fundamental tasks such as linear regression ICL, widely studied as a standard theoretical benchmark for Transformers, have not been thoroughly analyzed in the context of Mamba. To address this gap, we study the training dynamics of Mamba on the linear regression ICL task. By developing novel techniques tackling non-convex optimization with gradient descent related to Mamba's structure, we establish an exponential convergence rate to ICL solution, and derive a loss bound that is comparable to Transformer's. Importantly, our results reveal that Mamba can perform a variant of *online gradient* descent to learn the latent function in context. This mechanism is different from that of Transformer, which is typically understood to achieve ICL through gradient descent emulation. The theoretical results are verified by experimental simulation.

1 Introduction

State-space models (SSMs), notably Mamba (Gu and Dao, 2024), have recently emerged as compelling alternatives to Transformer-based architectures (Vaswani et al., 2017). Mamba integrates gating, convolutions, and state-space modeling with selection mechanisms, enabling linear-time complexity. This effectively addresses the quadratic computational costs typically associated with self-attention mechanisms in Transformers. Consequently, Mamba demonstrates superior efficiency in processing long sequences while maintaining or even surpassing Transformer performance across diverse benchmarks (Gu and Dao, 2024; Dao and Gu, 2024; Patro and Agneeswaran, 2024; Liu et al., 2024; Ahamed and Cheng, 2024; Li et al., 2024a,b).

In-context learning (ICL) (Brown et al., 2020) is a powerful paradigm that enables models to generalize to unseen tasks by dynamically leveraging contextual examples (such as input-output pairs) without task-specific fine-tuning. This capability has become a defining characteristic of large foundation models, significantly enhancing their flexibility and adaptability. While extensive research has provided substantial insights into Transformer-based ICL mechanisms (Garg et al., 2022; Gatmiry et al., 2024; Sander et al., 2024; Zheng et al., 2024; Zhang et al., 2025), the principles underlying

^{*}Equal contribution

[†]Corresponding author

Mamba's ability to perform in-context learning remain largely unexplored, highlighting a compelling research gap.

Recent empirical studies have examined Mamba's (ICL) capabilities, showing it matches Transformers on many standard ICL tasks, while surpassing them in specialized scenarios like sparse parity (Park et al., 2024; Grazzi et al., 2024). Bondaschi et al. (2025) theoretically analyzed its representational capacity for in-context learning of Markov chains, and Li et al. (2025a) investigated binary classification tasks with outliers. (Yang et al., 2024, 2025; Behrouz et al., 2025b,a) leverage the connection between SSMs and online learning to design new architectures. However, even the linear regression model, a canonical setting widely used for theoretical analysis of Transformer-based ICL mechanisms, remains theoretically underexplored in the context of Mamba. To fill this gap, we analyze Mamba's training dynamics on in-context linear regression tasks. More precisely, following the previous ICL analysis in Transformers (Garg et al., 2022; Zhang et al., 2024; Ahn et al., 2023), this paper focuses on a data generative model with N input-output pairs $(\{x_i, y_i\}_{i=1}^N)$ and a query input (x_a) satisfying $y = f(x) = w^{\top}x$, where x denotes the input and y denotes the output, and w is randomly sampled from Gaussian distribution, termed the *context*. In this work, we develop a rigorous theoretical framework to analyze how randomly initialized Mamba models, when trained through gradient descent, evolve to implement in-context learning. We demonstrate that the trained Mamba architecture dynamically leverages the input context to perform implicit estimation of the vector w. This estimation is achieved through hidden state updates that mimic online gradient descent steps, finally implementing prediction for $y_q = f(x_q) = w^{\top} x_q$. We also provide a loss bound that is comparable to Transformers'. Our contributions are summarized as follows:

- We construct a Mamba architecture (S6: S4 with selection) capable of ICL, establishing its exponential convergence rate to ICL solution, and further derive the loss bound after convergence. The loss matches that of Transformers.
- Technically, we develop novel techniques to address optimization challenges induced by random initialization and gradient descent, rigorously characterizing Mamba's training dynamics when trained from scratch.
- We reveal how trained Mamba achieves in-context linear regression by progressively aligning
 its hidden states with the *context* through sequential token processing. This finding provides
 a new perspective for understanding Mamba's ICL mechanism, distinct from Transformerbased approaches. All the above results are verified by experiments.

2 Related Work

In-Context Learning The seminal work of Brown et al. (2020) demonstrated the in-context learning capability in Transformers, showing their ability to infer functional mappings from input-output exemplars without weight updates. Garg et al. (2022) initiated the investigation of ICL from the perspective of learning particular function classes. Following these, a line of research analyze this phenomenon through the lens of algorithm imitation: Transformers can be trained to implement various learning algorithms that can mimic the latent functions in context, including: a single step of gradient descent (Von Oswald et al., 2023; Akyürek et al., 2023), statistical algorithms (Bai et al., 2023), reinforcement learning algorithm (Lin et al., 2024), multi-step gradient descent (Gatmiry et al., 2024), mesa-optimization (Zheng et al., 2024), Newton's method (Giannou et al., 2025), weighted preconditioned gradient descent (Li et al., 2025b), in context classification (Bu et al., 2024; Shen et al., 2024; Bu et al., 2025) among others.

Recent work extends ICL analysis beyond Transformers: (Lee et al., 2024; Park et al., 2024) empirically compared popular architectures (e.g., RNNs, CNNs, SSMs, Transformers) on synthetic ICL tasks, identifying capability variations across model types and task demands. Tong and Pehlevan (2024) demonstrate that MLPs can learn in-context a series of classical tasks such as regression and classification with less computation than Transformers. Sushma et al. (2024) show that state space models augmented with *local self-attention* can learn linear regression in-context. Unlike existing research on ICL, this work focuses on the ICL mechanism of Mamba (specifically S4 with selection) and its training dynamics.

Theoretical Understanding of SSMs As Gu et al. (2022) introduce structured state spaces models in modeling long sequence and further be extended to Mamba (Gu and Dao, 2024), which gained

significant attention as alternatives to Transformers, extensive research has sought to theoretically understand the mechanisms and capabilities of state-spaces models (SSMs). Dao and Gu (2024) propose the framework of state space duality, which establishes a connection between SSMs and attention variants through the lens of structured matrices. Vankadara et al. (2024) provide a scaling analysis of signal propagation in SSMs through the lens of feature learning. Cirone et al. (2024) draw the link of SSMs to linear CDEs (controlled differential equations) and use tools from rough path theory to study their expressivity. Chen et al. (2025) establish the computational limits of SSMs and Mamba via circuit complexity analysis, questioning the prevailing belief that Mamba possesses superior computational expressivity compared to Transformers. Nishikawa and Suzuki (2025) demonstrate that state space models integrated with nonlinear layers achieve dynamic token selection capabilities comparable to Transformers. Different from the above, we provide theoretical understanding of Mamba from the perspective of ICL.

3 Problem Setup

In this section, we outline the ICL data model, the Mamba model, the prediction strategy, and the gradient descent training algorithm.

Data Model. We consider an in-context linear regression task where each prompt corresponds to a new function $f(\boldsymbol{x}) = \boldsymbol{w}^{\top} \boldsymbol{x}$ with weights $\boldsymbol{w} \sim \mathcal{N}(0, \boldsymbol{I}_d)$ and d > 1. For each task, we generate N i.i.d. input-output pairs $\{(\boldsymbol{x}_i, y_i)\}_{i=1}^N$ and a query \boldsymbol{x}_q , where all inputs $\boldsymbol{x}_i, \boldsymbol{x}_q \sim \mathcal{N}(0, \boldsymbol{I}_d)$ are independent Gaussian vectors, and the outputs satisfy $y_i = f(\boldsymbol{x}_i)$. The goal is to predict $y_q = f(\boldsymbol{x}_q)$ for the query.

To enable sequential processing of prompts in the Mamba model, we implement an embedding strategy where:

- 1. The *i*-th context token is encoded as $e_i = (x_i^\top, y_i)^\top$, formed by concatenating input x_i with its corresponding label y_i .
- 2. The query token is represented as $e_q = (x_q^\top, 0)^\top$, masking the unknown target value with a zero placeholder.

In many theoretical analyses of Transformer-based in-context learning, token embeddings are conventionally concatenated into a single matrix to enable parallel computation of global attention (Zhang et al., 2024; Ahn et al., 2023; Huang et al., 2023; Mahankali et al., 2024; Wu et al., 2024). In contrast, since Mamba operates as a sequential model, we feed the embeddings of context tokens one by one, and finally the query token $(e_1 \rightarrow e_2 \rightarrow \cdots \rightarrow e_N \rightarrow e_q)$.

Mamba Model. We consider a S6 layer of Mamba $o_{1:L} = \text{Mamba}(\theta; u_{1:L})$ with selection, discretization, and linear recurrence components, where $u_l, o_l \in \mathbb{R}^{d_e}$. It can be described as follows:

$$\boldsymbol{h}_l^{(i)} = \overline{\boldsymbol{A}}_l \boldsymbol{h}_{l-1}^{(i)} + \overline{\boldsymbol{B}}_l \boldsymbol{u}_l^{(i)} \in \mathbb{R}^{d_h \times 1}, \quad \text{(1a)} \quad \overline{\boldsymbol{A}}_l = \exp(\Delta_l \boldsymbol{A}) \in \mathbb{R}^{d_h \times d_h}, \qquad \text{(2a)}$$

$$o_l^{(i)} = \boldsymbol{C}_l^{\top} \boldsymbol{h}_l^{(i)}, \quad \boldsymbol{C}_l \in \mathbb{R}^{d_h \times 1}, \quad \text{(1b)} \quad \overline{\boldsymbol{B}}_l = (\Delta_l \boldsymbol{A})^{-1} (\exp(\Delta_l \boldsymbol{A}) - \boldsymbol{I}) \Delta_l \boldsymbol{B}_l \in \mathbb{R}^{d_h \times 1} \quad \text{(2b)}$$
 for $i \in [d_e]$. Here, the superscript (i) denotes the i -th independent processing channel, where each channel operates on a unique feature dimension of the input \boldsymbol{u}_l and output \boldsymbol{o}_l vectors (i.e., $\boldsymbol{u}_l^{(i)}$ and $o_l^{(i)}$ correspond to the i -th elements of \boldsymbol{u}_l and \boldsymbol{o}_l , respectively). The hidden state $\boldsymbol{h}_l^{(i)}$ is initialized as $\boldsymbol{h}_0^{(i)} = \boldsymbol{0}$ and evolves according to $\overline{\boldsymbol{A}}_l \in \mathbb{R}^{d_h \times d_h}, \overline{\boldsymbol{B}}_l \in \mathbb{R}^{d_h \times 1}$ and the input $\boldsymbol{u}_l^{(i)}$. $\boldsymbol{C}_l \in \mathbb{R}^{d_h \times 1}$ maps the hidden state $\boldsymbol{h}_l^{(i)}$ to the output $o_l^{(i)}$. As shown in (2), $\overline{\boldsymbol{A}}_l$ and $\overline{\boldsymbol{B}}_l$ are computed using the zero-order hold (ZOH) discretization method applied to $\boldsymbol{A} \in \mathbb{R}^{d_h \times d_h}, \boldsymbol{B}_l \in \mathbb{R}^{d_h \times 1}$ and the timestep $\Delta_l \in \mathbb{R}$. Next, we describe the selection mechanism.

 $m{B}_l = m{W}_B m{u}_l + m{b}_B, \quad (3) \quad m{C}_l = m{W}_C m{u}_l + m{b}_C, \quad (4) \quad \Delta_l = \mathrm{softplus}(m{w}_\Delta^\top m{u}_l + b_\Delta), \quad (5)$ Here, $\mathrm{softplus}(\mathbf{x}) = \log(1 + \exp(x)). \quad m{W}_B, m{W}_C \in \mathbb{R}^{d_h \times d_e}, \, m{b}_B, m{b}_C \in \mathbb{R}^{d_h \times 1}, \, m{w}_\Delta \in \mathbb{R}^{d_e \times 1}, \, b_\Delta \in \mathbb{R}, \, \text{along with } m{A} \in \mathbb{R}^{d_h \times d_h} \, \text{are the parameters of the Mamba model. We use } m{\theta} \, \text{to denote the collection of all the parameters.}$ Unlike previous work (Sushma et al., 2024) that introduce *local self-attention* component to augment SSMs, which may inherit the Transformer's ICL ability, our model adheres to Mamba's original selective state-space framework (Gu and Dao, 2024). This alignment ensures us to mechanistically analyze how Mamba's architecture enables in-context learning (ICL).

Linear Regression Prediction. In this work, we set $d_e = d+1$, enabling the Mamba model to process the embeddings $e_{1:N}, e_q$. Given the prompt (e_1, \ldots, e_N, e_q) , the Mamba model will output a sequence $o_{1:N+1} = \mathbf{Mamba}(\boldsymbol{\theta}; e_1, \ldots, e_N, e_q)$. The prediction for the linear regression target $y_q = \boldsymbol{w}^{\top} \boldsymbol{x}_q$ is extracted from the terminal position of the output matrix (corresponding to the zero placeholder in the query token $e_q = (\boldsymbol{x}_q^{\top}, 0)^{\top}$). Concretely, $\hat{y}_q = \boldsymbol{o}_{N+1}^{(d+1)}$.

Training Algorithm. To train a Mamba model over the in-context linear regression task, we consider minimizing the following population loss:

$$\mathcal{L}(\boldsymbol{\theta}) = \mathbb{E}_{\boldsymbol{x}_{1:N}, \boldsymbol{x}_q, \boldsymbol{w}} \left[\frac{1}{2} (\hat{y}_q - y_q)^2 \right]. \tag{6}$$

Given a Mamba model, we use gradient descent to minimize population loss $\mathcal{L}(\theta)$, and the update of trainable parameters $\theta' = \{W_B, W_C, b_B, b_C\}$ can be written as follows:

$$\theta'(t+1) = \theta'(t) - \eta \nabla_{\theta'} \mathcal{L}(\theta(t)). \tag{7}$$

4 Main Results

This section presents our main theoretical results that characterize the convergence state of Mamba and its final loss. We also compare the results with other models.

Assumption 4.1 (1) Matrix $A = -I_{d_h}$. (2) The vector \mathbf{w}_{Δ} is fixed as zero $\mathbf{0}$, and b_{Δ} is fixed as $\ln(\exp((\ln 2)/N) - 1)$. (3) Matrices \mathbf{W}_B , \mathbf{W}_C are initialized with entries drawn i.i.d. from the standard Gaussian distribution $\mathcal{N}(0,1)$. (4) The hidden state dimension satisfies: $d_h = \widetilde{\Omega}(d^2)$. (5) The learning rate satisfies: $\eta = O(d^{-2}d_h^{-1})$. (6) Bias vectors \mathbf{b}_B , \mathbf{b}_C are initialized as zero $\mathbf{0}$. (7) Token length $N = \Omega(d)$.

(1) The negative-definite matrix $A = -I_{d_h}$ guarantees the stable convergence of hidden states $h_l^{(i)}$. (2) Given the zero-mean and symmetric distribution of embeddings, w_Δ can naturally converge to 0 during gradient descent, and we fix it as 0 for simplicity. We further fix b_Δ to an appropriate constant to maintain a suitable timestep Δ_l , enabling us to concentrate our theoretical analysis on W_B , W_C , b_B , and b_C . In prior works on Transformer-based in-context learning, merging key-query weights (e.g., $W := W_Q W_K$) and specific initializations (e.g., $W_Q = W_K = I$) are often adopted to simplify optimization analysis (Zhang et al., 2024; Ahn et al., 2023; Huang et al., 2023; Mahankali et al., 2024; Wu et al., 2024). (3, 4, 5) In contrast, our Gaussian initialization of W_B and W_C demonstrates more practicality, which requires a sufficiently large hidden state dimension d_h and a sufficiently small learning rate η to ensure favorable loss landscape properties. Assumption (6) is intended to simplify the analysis. (7) Token length should be larger enough than the dimension of w to capture sufficient contextual information.

Theorem 4.1 Under Assumption 4.1, if the Mamba is trained with gradient descent, and given a new prompt (e_1, \ldots, e_N, e_q) , then with probability at least $1 - \delta$ for some $\delta \in (0, 1)$, the trainable parameters $\theta'(t) = \{W_B(t), W_C(t), b_B(t), b_C(t)\}$ converge as $t \to \infty$ to parameters that satisfies:

(a) Projected hidden state:
$$(\mathbf{W}_{C}^{\top})_{[1:d,:]}(t)\mathbf{h}_{l}^{(d+1)} = \alpha(\mathbf{W}_{C}^{\top}(t))_{[1:d,:]}\mathbf{h}_{l-1}^{(d+1)} + (1-\alpha)\beta y_{l}\mathbf{x}_{l}$$
,

(b) Prediction for target:
$$\hat{y}_q = \boldsymbol{x}_q^{\top} \sum_{i=0}^{N-1} (1-\alpha) \alpha^{i+1} \beta y_{N-i} \boldsymbol{x}_{N-i}$$
,

(c) Population loss:
$$\mathcal{L}(\boldsymbol{\theta}(t)) \leq \frac{3d(d+1)}{2N}$$
,

where
$$\alpha = \exp((-\ln 2)/N)$$
, $\beta = \frac{2(1+\alpha)}{\alpha(3(1-\alpha)d+4-2\alpha)}$.

Theorem 4.1 characterizes the in-context learning (ICL) mechanism of Mamba and establishes an upper bound on its population loss. Specifically, (Thm 4.1 (a)) shows how the hidden state is updated according the given prompt $e_l = (\boldsymbol{x}_l^\top, y_l)^\top$. (Thm 4.1 (b)) presents the final prediction given prompt (e_1, \ldots, e_N, e_q) . (Thm 4.1 (c)) provides the upper bound for the population loss, which is comparable to that of the Transformer (Zhang et al., 2024). Next, we'll discuss it in more detail.

Update of Hidden State. If we define $\tilde{h}_l := (W_C^\top)_{[1:d,:]} h_l^{(d+1)}$, then (Thm 4.1 (a)) can be rewritten as follows:

$$\tilde{\boldsymbol{h}}_{l} = \alpha \tilde{\boldsymbol{h}}_{l-1} + (1 - \alpha)\beta y_{l} \boldsymbol{x}_{l} = \tilde{\boldsymbol{h}}_{l-1} + (1 - \alpha)(\beta y_{l} \boldsymbol{x}_{l} - \tilde{\boldsymbol{h}}_{l-1}). \tag{8}$$

We observe its intrinsic connection to **online gradient descent**, which updates the model parameters (\tilde{h}_l) with only one currently arriving sample $(e_l = (\boldsymbol{x}_l^\top, y_l)^\top)$ at each step. Specifically, the system gradually updates \tilde{h}_l along the pseudo-gradient direction $\beta y_l \boldsymbol{x}_l$, with a fixed step size $(1 - \alpha)$.

For a newly defined task $f(x) = w^{\top}x$, given that $\mathbb{E}[y_lx_l] = w$, the direction of \tilde{h}_l converges toward w as mamba processes multiple prompts. This demonstrates mamba's ability to internalize f(x) through prompt processing, which ultimately ensures that predictions for query token $e_q = (x_q^{\top}, 0)^{\top}$ closely approximate $f(x_q)$.

Previous works have shown that Transformer can mimic a single step of gradient descent to achieve incontext learning ability (Zhang et al., 2024; Mahankali et al., 2024). Concretely, a trained Transformer can be described as follows

Transformer
$$(e_1, \dots, e_N, e_q) \approx \boldsymbol{x}_q^{\top} \left(\frac{1}{N} \sum_{i=1}^N y_i \boldsymbol{x}_i\right) \approx \boldsymbol{x}_q^{\top} \boldsymbol{w}.$$
 (9)

Our theoretical analysis reveals that Mamba and Transformer have different in-context learning mechanisms. This divergence stems from their inherent architectural biases: Transformers process contexts globally through self-attention, while Mamba enforces local sequential dependencies via recurrent state transitions. These findings provide fundamental insights into the contrasting capabilities of Transformer-based and Mamba-based models for in-context learning. As experimental work shows, transformers can learn vector-valued MQAR tasks in the context which Mamba cannot, while Mamba succeeds in sparse-parity in-context learning tasks where Transformers fail(Park et al., 2024).

Prediction Outcome. Comparing equations Thm 4.1 (b) and (9), we found both similarities and distinctions in how Transformer and Mamba implement in-context learning (ICL). Both models leverage a weighted aggregation of $y_i x_i$, aligning with the intuition that learning $f(x) = x^\top w$ from context reduces to estimating the latent parameter w, since $\mathbb{E}[y_i x_i] = w$. Notably, their token weighting strategies diverge: Transformer's global attention mechanism implicitly assigns nearly uniform weights ($\sim \frac{1}{N}$, where N is the token length) to all $y_i x_i$, while Mamba's linear recurrence imposes position-dependent weight variations. This difference arises from Mamba's iterative state update rule, where the influence of prompt tokens e_i on the hidden state h_l depends on their sequential placement, governed by the model's linear recurrence dynamics.

The derived upper bound (Thm 4.1 (c)) establishes an O(1/N) convergence rate for the loss (ignoring dimension factor), demonstrating that Mamba matches the sample complexity scaling of Transformers in linear regression ICL tasks (Zhang et al., 2024).

Compare with S4. Mamba extends the structured state space model (S4) (Gu et al., 2022) by integrating a selection mechanism, which is critical for enabling ICL. In S4 model, the matrices $A \in \mathbb{R}^{d_h \times d_h}$, and $B, C \in \mathbb{R}^{d_h \times 1}$ are static, leading to a fixed linear combination of inputs:

$$o_l^{(i)} = \sum_{j=1}^l \mathbf{C}^\top \overline{\mathbf{A}}^{l-j} \overline{\mathbf{B}} u_j^{(i)}, \tag{10}$$

where the coefficients $C^{\top}\overline{A}^{l-j}\overline{B}$ are *task-agnostic*. This formulation inherently limits S4's ability to adapt to *task-specific* parameters w in ICL scenarios, as the model cannot adjust its inductive bias to match distinct w across different tasks. Therefore, the S4 model cannot truly learn in-context.

In contrast, Mamba's selection mechanism dynamically adjusts B_l and C_l (and optionally A_l) based on the input tokens (u_1, \ldots, u_N) . This allows the model to implicitly adapt its hidden state to

align with the latent w of each task, effectively transforming the linear combination weights into context-dependent functions $f(x) = x^{\top}w$. Such adaptability is essential for ICL, as it enables Mamba to reconstruct diverse w from input prompts without task-specific fine-tuning.

5 Proof Sketch

This section outlines the main technical ideas to prove Theorem 4.1. The complete proofs are given in the appendix.

Linear Recurrence. To start with, we show how the hidden states update when receiving token $e_l = (x_l^\top, y_l)^\top$. By (Eq. (5)) and Assumption 4.1(2), we have $\Delta_l = (\ln 2)/N$. Combining it with (Eq. (1)(2)) and get:

$$\boldsymbol{h}_{l}^{(d+1)} = \alpha \boldsymbol{h}_{l-1}^{(d+1)} + (1 - \alpha) y_{l} \boldsymbol{B}_{l}, \tag{11}$$

where $\alpha := \exp(-\Delta_l) = \exp((-\ln 2)/N)$, the second equality is by discretization rule (2), the third equality is by Assumption 4.1(2) and $\exp(-\Delta_l \mathbf{I}) = \exp(-\Delta_l) \mathbf{I}$.

Prediction Output. We next derive the expression of \hat{y}_q . By recurring (Eq.(11)), the hidden state after receiving the first l context prompts $e_{1:l}$ is given by $\boldsymbol{h}_l^{(d+1)} = (1-\alpha) \sum_{i=0}^{l-1} \alpha^i y_{l-i} \boldsymbol{B}_{l-i}$. Receiving all the prompt tokens $\boldsymbol{e}_{1:N}$ and the query token $\boldsymbol{e}_q = (\boldsymbol{x}_q^\top, 0)^\top$, we have:

$$\boldsymbol{h}_{N+1}^{(d+1)} = \alpha \boldsymbol{h}_{N}^{(d+1)} + (1 - \alpha) \cdot 0 \cdot \boldsymbol{B}_{N} = (1 - \alpha) \sum_{i=0}^{N-1} \alpha^{i+1} y_{N-i} \boldsymbol{B}_{N-i}.$$
 (12)

Finally, the prediction output is as follows

$$\hat{y}_q = \mathbf{C}_{N+1}^{\top} \mathbf{h}_{N+1}^{(d+1)} = (1 - \alpha) (\mathbf{W}_C \mathbf{e}_q + \mathbf{b}_C)^{\top} \sum_{i=0}^{N-1} \alpha^{i+1} y_{N-i} (\mathbf{W}_B \mathbf{e}_{N-i} + \mathbf{b}_B).$$
(13)

To handle $W_C e_q$ and $W_B e_{N-i}$, we further decompose $W_B = [B b]$ and $W_C = [C c]$, where $B, C \in \mathbb{R}^{d_h \times d}$, $b, c \in \mathbb{R}^{d_h \times 1}$. Then we write another form of (Eq. (13)):

$$\hat{y}_q = (1 - \alpha)(Cx_q + b_C)^{\top} \sum_{i=0}^{N-1} \alpha^{i+1} y_{N-i} (Bx_{N-i} + y_{N-i}b + b_B).$$
 (14)

The loss becomes:

$$\mathcal{L}(\boldsymbol{\theta}) = \frac{1}{2} \mathbb{E} \left[\left((1 - \alpha) (\boldsymbol{C} \boldsymbol{x}_q + \boldsymbol{b}_C)^\top \sum_{i=0}^{N-1} \alpha^{i+1} y_{N-i} (\boldsymbol{B} \boldsymbol{x}_{N-i} + y_{N-i} \boldsymbol{b} + \boldsymbol{b}_B) - \boldsymbol{w}^\top \boldsymbol{x}_q \right)^2 \right]. \tag{15}$$

By computing the gradient of C, b_C , B, b and b_B with respect to $\mathcal{L}(\theta(t))$, we derive the following update rule according to Eq. (7).

Lemma 5.1 (Update Rule) Let η be the learning rate and we use gradient descent to update the weights W_B, W_C, b_B, b_C , for $t \ge 0$ we have

$$\begin{split} \boldsymbol{B}(t+1) &= \boldsymbol{B}(t) + \eta \beta_3 \boldsymbol{C}(t) - \eta \beta_1 \boldsymbol{C}(t) \boldsymbol{C}(t)^\top \boldsymbol{B}(t), \\ \boldsymbol{C}(t+1) &= \boldsymbol{C}(t) + \eta \beta_3 \boldsymbol{B}(t) - \eta \beta_1 \boldsymbol{B}(t) \boldsymbol{B}(t)^\top \boldsymbol{C}(t) - \eta \beta_2 \boldsymbol{b}(t) \boldsymbol{b}(t)^\top \boldsymbol{C}(t), \\ \boldsymbol{b}(t+1) &= \boldsymbol{b}(t) - \eta \beta_2 \boldsymbol{C}(t) \boldsymbol{C}(t)^\top \boldsymbol{b}(t), \quad \boldsymbol{b}_B(t) = \boldsymbol{b}_C(t) = \boldsymbol{0}, \\ where \ \beta_1 &= \mathbb{E} \Big[\sum_{i=0}^{N-1} \sum_{j=0}^{N-1} (1-\alpha)^2 \alpha^{i+j+2} y_{N-i} y_{N-j} \boldsymbol{x}_{N-i} \boldsymbol{x}_{N-j}^\top \Big], \ \beta_2 = \mathbb{E} \Big[\sum_{i=0}^{N-1} \sum_{j=0}^{N-1} (1-\alpha)^2 \alpha^{i+j+2} y_{N-i} y_{N-j} \boldsymbol{x}_{N-i} \boldsymbol{w}^\top \Big]. \end{split}$$

Technical Challenges. Unlike many prior Transformer-based ICL analyses that simplify dynamics via merged weights or special initializations, our Gaussian-initialized W_B , W_C and discrete-time gradient descent introduces more complexity (cf. assumption 4.1). To solve the optimization problem described in Lemma 5.1, we have the following three questions to answer: (1) Convergence Target: Where do the parameters converge? (2) Convergence Proof: How to rigorously establish convergence? (3) Saddle Point Avoidance: How to avoid saddle points? To answer these three questions, we propose two key techniques: *Vector-coupled Dynamic*, *Negative Feedback Convergence*, and apply them with a *Fine-grained Induction*. We next describe them in detail.

5.1 Vector-coupled Dynamics

We can verify by Lemma 5.1 that $C^{\top}B = \text{Diag}(a_1, \dots, a_d)$ with $a_i \in \{0, \frac{\beta_3}{\beta_1}\}$, $C^{\top}b = 0$ are the fixed points for the parameters W_B , W_C .

Combining the loss function Eq. 15 and $b_B(t) = b_C(t) = 0$ in Lemma 5.1, the loss function can be rewritten as

$$\mathcal{L}(\boldsymbol{\theta}) = \frac{1}{2} \mathbb{E} \Big[\Big((1 - \alpha) \sum_{i=0}^{N-1} \alpha^{i+1} (\boldsymbol{x}_q^\top \boldsymbol{C}^\top \boldsymbol{B} y_{N-i} \boldsymbol{x}_{N-i} + y_{N-i}^2 \boldsymbol{x}_q^\top \boldsymbol{C}^\top \boldsymbol{b}) - \boldsymbol{w}^\top \boldsymbol{x}_q \Big)^2 \Big].$$

To minimize loss, the term $(1-\alpha)\sum_{i=0}^{N-1}\alpha^{i+1}(\boldsymbol{x}_q^\top\boldsymbol{C}^\top\boldsymbol{B}y_{N-i}\boldsymbol{x}_{N-i}+y_{N-i}^2\boldsymbol{x}_q^\top\boldsymbol{C}^\top\boldsymbol{b})$ should approximate $\boldsymbol{w}^\top\boldsymbol{x}_q$. Given $\mathbb{E}[y_{N-i}\boldsymbol{x}_{N-i}]=\boldsymbol{w}$ and $\mathbb{E}[y_{N-i}^2]>0$, we derive that $\boldsymbol{C}^\top\boldsymbol{B}$ should converge to $\frac{\beta_3}{\beta_1}\boldsymbol{I}$, while $\boldsymbol{C}^\top\boldsymbol{b}$ converges to $\boldsymbol{0}$ to minimize the loss. However, as mentioned above, $\boldsymbol{C}^\top\boldsymbol{B}=\mathrm{Diag}(a_1,\ldots,a_d)$ with partial $a_i=0$ can also enable convergence, which is an undesirable scenario.

To analyze the convergence behavior of $C^{\top}B$ and $C^{\top}b$, we introduce the *Vector-coupled Dynamics* technique, which studies the inner product dynamics between decomposed column vectors of B and C. Specifically, we decompose B and C into $B = [b_1 \dots b_d]$, $C = [c_1 \dots c_d]$. Then we have another form of Lemma 5.1 for B, C and D as the following lemma.

Lemma 5.2 (Vectors Update Rule) *Let* η *be the learning rate and we use gradient descent to update the weights* W_B, W_C, b_B, b_C , *for* $i \in [d]$, $t \geq 0$ *we have*

$$\boldsymbol{b}_i(t+1) = \boldsymbol{b}_i(t) + \eta \Big(\big(\beta_3 - \beta_1 \boldsymbol{c}_i^\top(t) \boldsymbol{b}_i(t) \big) \boldsymbol{c}_i(t) - \beta_1 \sum_{k \neq i}^d \boldsymbol{c}_k^\top(t) \boldsymbol{b}_i(t) \cdot \boldsymbol{c}_k(t) \Big),$$

$$\boldsymbol{c}_i(t+1) = \boldsymbol{c}_i(t) + \eta \Big(\big(\beta_3 - \beta_1 \boldsymbol{c}_i^\top(t) \boldsymbol{b}_i(t) \big) \boldsymbol{b}_i(t) - \beta_1 \sum_{k \neq i}^d \boldsymbol{c}_i^\top(t) \boldsymbol{b}_k(t) \cdot \boldsymbol{b}_k(t) - \beta_2 \boldsymbol{c}_i^\top(t) \boldsymbol{b}(t) \cdot \boldsymbol{b}(t) \Big),$$

$$\boldsymbol{b}(t+1) = \boldsymbol{b}(t) - \eta \Big(\beta_2 \sum_{k=1}^d \boldsymbol{c}_k^\top(t) \boldsymbol{b}(t) \cdot \boldsymbol{c}_k(t) \Big).$$

With Lemma 5.2, we can further analyze the dynamics of the inner products $c_i^{\top}(t)b_i(t)$, $c_i^{\top}(t)b_j(t)$ and $c_i^{\top}(t)b(t)$, precisely characterizing the behavior of $C^{\top}B$ and $C^{\top}b$. This technique helps answer the question "Where do the parameters converge?"

5.2 Negative Feedback Convergence

As we discuss in Section 5.1, to minimize loss, the following conditions must be satisfied for all $i, j \in [d]$ with $i \neq j$: $\mathbf{c}_i^\top(t)\mathbf{b}_i(t) \to \frac{\beta_3}{\beta_1}$, $\mathbf{c}_i^\top(t)\mathbf{b}_j(t) \to 0$, $\mathbf{c}_i^\top(t)\mathbf{b}(t) \to 0$. To establish the convergence, we introduce the *Negative Feedback Convergence* technique. This technique leverages the negative feedback terms in the dynamical equations of $\mathbf{c}_i^\top(t)\mathbf{b}_i(t)$, $\mathbf{c}_i^\top(t)\mathbf{b}_j(t)$, and $\mathbf{c}_i^\top(t)\mathbf{b}(t)$ to derive an exponential convergence rate. Taking $\mathbf{c}_i^\top(t)\mathbf{b}_i(t)$ as an example, we derive the following

update rule by Lemma 5.2.

$$(\beta_{3} - \beta_{1} \boldsymbol{c}_{i}^{\top}(t+1)\boldsymbol{b}_{i}(t+1)) = \beta_{3} - \beta_{1} \boldsymbol{c}_{i}^{\top}(t)\boldsymbol{b}_{i}(t)$$

$$-\eta \beta_{1} (\beta_{3} - \beta_{1} \boldsymbol{c}_{i}^{\top}(t)\boldsymbol{b}_{i}(t)) \boldsymbol{b}_{i}^{\top}(t)\boldsymbol{b}_{i}(t) - \eta \beta_{1} (\beta_{3} - \beta_{1} \boldsymbol{c}_{i}^{\top}(t)\boldsymbol{b}_{i}(t)) \boldsymbol{c}_{i}^{\top}(t)\boldsymbol{c}_{i}(t)$$

$$= negative\ feedback\ term$$

$$+ \eta \beta_{1}^{2} \sum_{k \neq i}^{d} \boldsymbol{c}_{i}^{\top}(t)\boldsymbol{b}_{k}(t) \cdot \boldsymbol{b}_{k}^{\top}(t)\boldsymbol{b}_{i}(t) + \eta \beta_{1}^{2} \sum_{k \neq i}^{d} \boldsymbol{c}_{k}^{\top}(t)\boldsymbol{b}_{i}(t) \cdot \boldsymbol{c}_{i}^{\top}(t)\boldsymbol{c}_{k}(t)$$

$$+ \eta \beta_{1}\beta_{2}\boldsymbol{c}_{i}^{\top}(t)\boldsymbol{b}(t) \cdot \boldsymbol{b}_{i}^{\top}(t)\boldsymbol{b}(t) - \beta_{1}(\boldsymbol{c}(t+1) - \boldsymbol{c}(t))^{\top}(\boldsymbol{b}(t+1) - \boldsymbol{b}(t)).$$

$$(16)$$

The term $(\beta_3 - \beta_1 \boldsymbol{c}_i^{\top}(t+1)\boldsymbol{b}_i(t+1))$ decomposes into its previous state $(\beta_3 - \beta_1 \boldsymbol{c}_i^{\top}(t)\boldsymbol{b}_i(t))$ (marked with underline) plus the remaining terms (increment terms). The increment terms includes a *negative* feedback term, which induces a tendency to drive $(\beta_3 - \beta_1 \boldsymbol{c}_i^{\top}(t)\boldsymbol{b}_i(t))$ to $(\boldsymbol{c}_i^{\top}(t)\boldsymbol{b}_i(t) \to \frac{\beta_3}{\beta_1})$.

Intuitively, $\boldsymbol{b}_i^\top(t)\boldsymbol{b}_i(t)$ and $\boldsymbol{c}_i^\top(t)\boldsymbol{c}_i(t)$ are much larger than $\boldsymbol{b}_k^\top(t)\boldsymbol{b}_i(t)$, $\boldsymbol{c}_i^\top(t)\boldsymbol{c}_k(t)$ and $\boldsymbol{b}_i^\top(t)\boldsymbol{b}(t)$ at Gaussian initialization with high probability. Also, as $\boldsymbol{c}_i^\top(t)\boldsymbol{b}_j(t)$, $\boldsymbol{b}_i^\top(t)\boldsymbol{b}(t) \to 0$ with $i \neq j$ and η is small enough, the effect of *negative feedback term* is the dominant term in the increment terms. Therefore, denoting $y(t) = \beta_3 - \beta_1 \boldsymbol{c}_i^\top(t)\boldsymbol{b}_i(t)$ and $\xi(t) = y(t+1) - y(t)$ - negative feedback term we can model the update rule of (Eq. 16) as follows:

$$y(t+1) = (1 - \eta \beta_1 (\boldsymbol{b}_i^{\top}(t)\boldsymbol{b}_i(t) + \boldsymbol{c}_i^{\top}(t)\boldsymbol{c}_i(t)))y(t) + \xi(t).$$

Recur this formula from 0 to t, we have:

$$y(t+1) = \prod_{s=0}^{t} \left(1 - \eta \beta_1 \left(\boldsymbol{b}_i^{\top}(s) \boldsymbol{b}_i(s) + \boldsymbol{c}_i^{\top}(s) \boldsymbol{c}_i(s) \right) \right) y(0)$$

$$+ \sum_{s=0}^{t} \prod_{s'=s+1}^{t} \left(1 - \eta \beta_1 \left(\boldsymbol{b}_i^{\top}(s') \boldsymbol{b}_i(s') + \boldsymbol{c}_i^{\top}(s') \boldsymbol{c}_i(s') \right) \right) \xi(s').$$
(17)

Denoting $\gamma = \min\{\boldsymbol{b}_i^\top(s)\boldsymbol{b}_i(s), \boldsymbol{c}_i^\top(s)\boldsymbol{c}_i(s)\}$ for $s \in [0,t]$, the first term on the RHS of (Eq. (17)) can be upper bounded by $(1-2\eta\beta_1\gamma)^{t+1}y(0)$. if $\xi(s')$ has an exponentially decaying upper bound (it can be proved when $\boldsymbol{c}_i^\top(t)\boldsymbol{b}_j(t) \to 0, \boldsymbol{c}_i^\top(t)\boldsymbol{b}(t) \to 0$ with an exponential rate), the second term on the RHS of (Eq. (17)) has an exponentially decaying upper bound. Therefore, we can establish an exponential convergence rate for $\boldsymbol{c}_i^\top(t)\boldsymbol{b}_i(t) \to \frac{\beta_3}{\beta_1}$. The similar method can be used on $\boldsymbol{c}_i^\top(t)\boldsymbol{b}_j(t) \to 0, \boldsymbol{c}_i^\top(t)\boldsymbol{b}(t) \to 0$. This technique helps answer the question "How to rigorously establish convergence?"

5.3 Fine-grained Induction

The exponential convergence of $c_i^{\top}(t)b_i(t) \to \frac{\beta_3}{\beta_1}$ under the *Negative Feedback Convergence* framework requires the following two conditions for all $i, j \in [d]$ with $i \neq j$:

- (1) $\boldsymbol{b}_i^{\top}(t)\boldsymbol{b}_i(t)$ and $\boldsymbol{c}_i^{\top}(t)\boldsymbol{c}_i(t)$ dominate $\boldsymbol{b}_i^{\top}(t)\boldsymbol{b}_j(t)$, $\boldsymbol{c}_i^{\top}(t)\boldsymbol{c}_j(t)$ and $\boldsymbol{b}_i^{\top}(t)\boldsymbol{b}(t)$ in magnitude.
- (2) $c_i^{\top}(t)b_i(t) \to 0, c_i^{\top}(t)b(t) \to 0$ at an exponentially decaying rate.

On the one hand, condition (1) at initialization (t=0) can be established via concentration inequalities, and critically, the preservation of Condition (1) for t>0 relies on the rapid decay of $\boldsymbol{c}_i^\top(t)\boldsymbol{b}_i(t), \boldsymbol{c}_i^\top(t)\boldsymbol{b}_j(t)$, and $\boldsymbol{c}_i^\top(t)\boldsymbol{b}(t)$ (condition (2)). On the other hand, under the framework of Negative Feedback Convergence, $\boldsymbol{c}_i^\top(t)\boldsymbol{b}_j(t)\to 0$ in Condition (2) also relies on Condition (1) and the rapid decay of $\boldsymbol{c}_i^\top(t)\boldsymbol{b}_i(t)\to \frac{\beta_3}{\beta_1}, \boldsymbol{c}_i^\top(t)\boldsymbol{b}(t)\to 0$. This implies mutual dependencies among the bounds of these Vector-coupled inner products.

To handle these dependencies and establish stable bounds, we introduce the technique Fine-grained Induction: Divide the inner products into three groups: (1) Squared norms: $\boldsymbol{b}_i^{\top}(t)\boldsymbol{b}_i(t), \boldsymbol{c}_i^{\top}(t)\boldsymbol{c}_i(t), \boldsymbol{b}^{\top}(t)\boldsymbol{b}(t)$. (2) Target terms: $\boldsymbol{c}_i^{\top}(t)\boldsymbol{b}_i(t), \boldsymbol{c}_i^{\top}(t)\boldsymbol{b}_j(t), \boldsymbol{c}_i^{\top}(t)\boldsymbol{b}(t)$. (3) Crossinteractions: $\boldsymbol{b}_i^{\top}(t)\boldsymbol{b}_j(t), \boldsymbol{c}_i^{\top}(t)\boldsymbol{c}_j(t), \boldsymbol{b}_i^{\top}(t)\boldsymbol{b}(t)$. And then carefully give bounds for them with an induction.

Specifically, denoting $\delta(t) = \max_{s \in [0,t]} \{2\sqrt{d_h \log(4d(2d+1)/\delta)}, |\boldsymbol{b}_i^\top(s)\boldsymbol{b}_j(s)|, |\boldsymbol{c}_i^\top(s)\boldsymbol{c}_j(s)|, |\boldsymbol{b}_i^\top(s)\boldsymbol{b}(s)| \}$ and $\gamma = \frac{1}{2}d_h \leq \min_{t \geq 0} \{\boldsymbol{b}_i^\top(t)\boldsymbol{b}_i(t), \boldsymbol{c}_i^\top(t)\boldsymbol{c}_i(t), \boldsymbol{b}^\top(t)\boldsymbol{b}(t) \}$, we establish the following three properties $\mathcal{A}(t)$, $\mathcal{B}(t)$, and $\mathcal{C}(t)$ simultaneously for $t \geq 0$:

$$\mathcal{A}(t): d_h/2 \leq \boldsymbol{b}_i^{\top}(t)\boldsymbol{b}_i(t), \boldsymbol{c}_i^{\top}(t)\boldsymbol{c}_i(t), \boldsymbol{b}^{\top}(t)\boldsymbol{b}(t) \leq 2d_h.$$

$$\mathcal{B}(t): \qquad |\beta_3 - \beta_1 \boldsymbol{c}_i^{\top}(t) \boldsymbol{b}_i(t)| \le \delta(t) \exp(-\eta \beta_1 \gamma t), \quad |\boldsymbol{c}_i^{\top}(t) \boldsymbol{b}_j(t)| \le 2\delta(t) \exp(-\eta \beta_1 \gamma t),$$

$$|\boldsymbol{c}_i^{\top}(t)\boldsymbol{b}(t)| \le 2\delta(t)\exp(-\eta\beta_2\gamma t) + \frac{\delta(t)}{\beta_2}\exp(-\eta\beta_1\gamma t).$$

$$\mathcal{C}(t): \qquad |\boldsymbol{b}_i^{\top}(t)\boldsymbol{b}_i(t)|, |\boldsymbol{c}_i^{\top}(t)\boldsymbol{c}_i(t)|, |\boldsymbol{b}_i^{\top}(t)\boldsymbol{b}(t)| \leq \delta(t) \leq 3\sqrt{d_h \log(4d(2d+1)/\delta)}.$$

The initial conditions $\mathcal{A}(0)$, $\mathcal{B}(0)$, and $\mathcal{C}(0)$ are established with high probability by concentration inequalities. We also provide the following claims to establish $\mathcal{A}(t)$, $\mathcal{B}(t)$, and $\mathcal{C}(t)$ for $t \geq 0$:

Claim 5.1
$$\mathcal{A}(0), \dots, \mathcal{A}(T), \mathcal{B}(0), \dots, \mathcal{B}(T), \mathcal{C}(0), \dots, \mathcal{C}(T) \Longrightarrow \mathcal{A}(T+1).$$

Claim 5.2
$$\mathcal{A}(0), \dots, \mathcal{A}(T), \mathcal{B}(0), \dots, \mathcal{B}(T), \mathcal{C}(0), \dots, \mathcal{C}(T) \Longrightarrow \mathcal{B}(T+1).$$

Claim 5.3
$$\mathcal{A}(0), \dots, \mathcal{A}(T), \mathcal{B}(0), \dots, \mathcal{B}(T), \mathcal{C}(0), \dots, \mathcal{C}(T) \Longrightarrow \mathcal{C}(T+1)$$
.

This induction answers the question "How to avoid saddle points?" because $\mathcal{B}(t)$ guarantees that $C^{\top}B \to \frac{\beta_3}{\beta_1}I$ and $C^{\top}b \to 0$, preventing stagnation of partial diagonal entries of $C^{\top}B$ at zero. Theorem 4.1 can be proved by substituting $C^{\top}B = \frac{\beta_3}{\beta_1}I$, $C^{\top}b = 0$, $b_B = b_C = 0$ into (Eq. (11), (14), (15))

6 Experimental Results

We present simulation results on synthetic data to verify our theoretical results. More experimental results can be found in Appendix E.

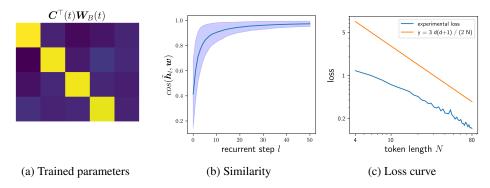


Figure 1: (a) Post-training visualization of matrix product $C^{\top}W_B$; (b) Cosine similarity evolution between w and $\tilde{h}_l = (W_C^{\top})_{[1:d,:]}h_l^{(d+1)}$ across recurrent steps l (after processing prompts $e_{1:l}$); (c) Test loss versus token sequence length N. Blue curve: experimental results; orange curve: theoretical upper bound.

Experiments Setting We follow Section 3 to generate the dateset and initialize the model. Specifically, we set dimension d=4, $d_h=80$, prompt token length N=50, and train the Mamba model on 3000 sequences by gradient descent. After training, we save the model and test it on 1000 new generated sequences, tracking the cosine similarity between $\tilde{\boldsymbol{h}}_l (:= (\boldsymbol{W}_C^\top)_{[1:d,:]} \boldsymbol{h}_l^{(d+1)})$ and \boldsymbol{w} . Moreover, we vary the length of the prompt token N from 4 to 80 and compare the test loss with the theoretical upper bound. For each N, we conduct 10 independent experiments and report the averaged results. All experiments are performed on an NVIDIA A800 GPU.

Experiment Result Recalling that we denote $W_B = [B \ b]$, Figure 1a reveals the convergence of $C^{\top}B$ to a diagonal matrix and $C^{\top}b$ to 0, confirming the theoretical induction presented in Section 5, also consistent with (Thm 4.1 (b)). Figure 1b shows that the projected hidden state \tilde{h}_l gradually aligns with w as more prompt tokens are processed, consistent with (Thm 4.1 (a)). Figure 1c demonstrates that the experimental loss has an upper bound $\frac{3d(d+1)}{2N}$ that decays linearly with N, aligning with (Thm 4.1 (c)).

7 Conclusion

This paper study Mamba's in-context learning mechanism, and rigorously establish its convergence and loss bound. By analysing the *Vector-coupled Dynamics*, we provide an exponential convergence rate with *Negative Feedback Convergence* technique in a *Fine-grained Induction*, and finally establish a O(1/N) loss bound. The loss bound is comparable to that of Transformer. Our theoretical results reveal the different mechanism between Transformer and Mamba on ICL, where Mamba emulates a variant of *online gradient descent* to perform in-context, while Transformers approximate a single step of gradient descent. Furthermore, our comparison with the S4 model demonstrates that the selection components are essential for Mamba to perform ICL.

Limitations and Social Impact Our analysis focuses on one-layer Mamba model, thus the behavior of Mamba with multi-layer or augmented with other components such as MLP is still unclear. We believe that our work will provide insight for those cases and can be used to study more data models such as nonlinear features. This paper is mainly a theoretical investigation, and we do not see an immediate social impact.

Acknowledgements

We thank the anonymous reviewers for their insightful comments to improve the paper. Miao Zhang was partially sponsored by the National Natural Science Foundation of China under Grant 62306084 and U23B2051, Shenzhen College Stability Support Plan under Grant GXWD20231128102243003, and Shenzhen Science and Technology Program under Grant ZDSYS20230626091203008 and KJZD20230923115113026.

References

- Ahamed, M. A. and Cheng, Q. (2024). Timemachine: A time series is worth 4 mambas for long-term forecasting. In ECAI 2024: 27th European Conference on Artificial Intelligence, 19-24 October 2024, Santiago de Compostela, Spain-Including 13th Conference on Prestigious Applications of Intelligent Systems. European Conference on Artificial Intelli, volume 392, page 1688.
- Ahn, K., Cheng, X., Daneshmand, H., and Sra, S. (2023). Transformers learn to implement preconditioned gradient descent for in-context learning. *Advances in Neural Information Processing Systems*, 36:45614–45650.
- Akyürek, E., Schuurmans, D., Andreas, J., Ma, T., and Zhou, D. (2023). What learning algorithm is in-context learning? investigations with linear models. In *The Eleventh International Conference on Learning Representations*.
- Arora, S., Cohen, N., Golowich, N., and Hu, W. (2019). A convergence analysis of gradient descent for deep linear neural networks. In *International Conference on Learning Representations*.
- Bai, Y., Chen, F., Wang, H., Xiong, C., and Mei, S. (2023). Transformers as statisticians: Provable in-context learning with in-context algorithm selection. *Advances in neural information processing systems*, 36:57125–57211.
- Behrouz, A., Li, Z., Kacham, P., Daliri, M., Deng, Y., Zhong, P., Razaviyayn, M., and Mirrokni, V. (2025a). Atlas: Learning to optimally memorize the context at test time. *arXiv preprint arXiv:2505.23735*.

- Behrouz, A., Razaviyayn, M., Zhong, P., and Mirrokni, V. (2025b). It's all connected: A journey through test-time memorization, attentional bias, retention, and online optimization. *arXiv* preprint arXiv:2504.13173.
- Bondaschi, M., Rajaraman, N., Wei, X., Ramchandran, K., Pascanu, R., Gulcehre, C., Gastpar, M., and Makkuva, A. V. (2025). From markov to laplace: How mamba in-context learns markov chains. *arXiv preprint arXiv:2502.10178*.
- Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J. D., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., et al. (2020). Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- Bu, D., Huang, W., Han, A., Nitanda, A., Suzuki, T., Zhang, Q., and Wong, H.-S. (2024). Provably transformers harness multi-concept word semantics for efficient in-context learning. *Advances in Neural Information Processing Systems*, 37:63342–63405.
- Bu, D., Huang, W., Han, A., Nitanda, A., Zhang, Q., Wong, H.-S., and Suzuki, T. (2025). Provable in-context vector arithmetic via retrieving task concepts. In *Forty-second International Conference on Machine Learning*.
- Chen, Y., Li, X., Liang, Y., Shi, Z., and Song, Z. (2025). The computational limits of state-space models and mamba via the lens of circuit complexity. In *The Second Conference on Parsimony and Learning (Proceedings Track)*.
- Cirone, N. M., Orvieto, A., Walker, B., Salvi, C., and Lyons, T. (2024). Theoretical foundations of deep selective state-space models. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*.
- Dao, T. and Gu, A. (2024). Transformers are SSMs: Generalized models and efficient algorithms through structured state space duality. In *International Conference on Machine Learning (ICML)*.
- Du, S. and Hu, W. (2019). Width provably matters in optimization for deep linear neural networks. In *International Conference on Machine Learning*, pages 1655–1664. PMLR.
- Garg, S., Tsipras, D., Liang, P. S., and Valiant, G. (2022). What can transformers learn in-context? a case study of simple function classes. *Advances in Neural Information Processing Systems*, 35:30583–30598.
- Gatmiry, K., Saunshi, N., Reddi, S., Jegelka, S., and Kumar, S. (2024). Can looped transformers learn to implement multi-step gradient descent for in-context learning? In *Proceedings of the 41st International Conference on Machine Learning*, pages 15130–15152.
- Giannou, A., Yang, L., Wang, T., Papailiopoulos, D., and Lee, J. D. (2025). How well can transformers emulate in-context newton's method? In *The 28th International Conference on Artificial Intelligence and Statistics*.
- Grazzi, R., Siems, J. N., Schrodi, S., Brox, T., and Hutter, F. (2024). Is mamba capable of in-context learning? In *Proceedings of the Third International Conference on Automated Machine Learning*, volume 256 of *Proceedings of Machine Learning Research*, pages 1/1–26. PMLR.
- Gu, A. and Dao, T. (2024). Mamba: Linear-time sequence modeling with selective state spaces. In *First Conference on Language Modeling*.
- Gu, A., Goel, K., and Re, C. (2022). Efficiently modeling long sequences with structured state spaces. In *International Conference on Learning Representations*.
- Huang, Y., Cheng, Y., and Liang, Y. (2023). In-context convergence of transformers. *arXiv preprint arXiv:2310.05249*.
- Lee, I., Jiang, N., and Berg-Kirkpatrick, T. (2024). Is attention required for icl? exploring the relationship between model architecture and in-context learning ability. In *The Twelfth International Conference on Learning Representations*.

- Li, H., Lu, S., Cui, X., Chen, P.-Y., and Wang, M. (2025a). Understanding mamba in in-context learning with outliers: A theoretical generalization analysis. In *High-dimensional Learning Dynamics*.
- Li, K., Chen, G., Yang, R., and Hu, X. (2024a). Spmamba: State-space model is all you need in speech separation. *arXiv preprint arXiv:2404.02063*.
- Li, K., Li, X., Wang, Y., He, Y., Wang, Y., Wang, L., and Qiao, Y. (2024b). Videomamba: State space model for efficient video understanding. In *European Conference on Computer Vision*, pages 237–255. Springer.
- Li, Y., Tarzanagh, D. A., Rawat, A. S., Fazel, M., and Oymak, S. (2025b). Gating is weighting: Understanding gated linear attention through in-context learning. *arXiv* preprint arXiv:2504.04308.
- Lin, L., Bai, Y., and Mei, S. (2024). Transformers as decision makers: Provable in-context reinforcement learning via supervised pretraining. In *The Twelfth International Conference on Learning Representations*.
- Liu, Y., Tian, Y., Zhao, Y., Yu, H., Xie, L., Wang, Y., Ye, Q., Jiao, J., and Liu, Y. (2024). Vmamba: Visual state space model. *Advances in neural information processing systems*, 37:103031–103063.
- Mahankali, A. V., Hashimoto, T., and Ma, T. (2024). One step of gradient descent is provably the optimal in-context learner with one layer of linear self-attention. In *The Twelfth International Conference on Learning Representations*.
- Nishikawa, N. and Suzuki, T. (2025). State space models are provably comparable to transformers in dynamic token selection. In *The Thirteenth International Conference on Learning Representations*.
- Park, J., Park, J., Xiong, Z., Lee, N., Cho, J., Oymak, S., Lee, K., and Papailiopoulos, D. (2024). Can mamba learn how to learn? a comparative study on in-context learning tasks. In *Proceedings of* the 41st International Conference on Machine Learning, pages 39793–39812.
- Patro, B. N. and Agneeswaran, V. S. (2024). Mamba-360: Survey of state space models as transformer alternative for long sequence modelling: Methods, applications, and challenges. *arXiv* preprint *arXiv*:2404.16112.
- Sander, M. E., Giryes, R., Suzuki, T., Blondel, M., and Peyré, G. (2024). How do transformers perform in-context autoregressive learning? In *Proceedings of the 41st International Conference on Machine Learning*, pages 43235–43254.
- Shen, W., Zhou, R., Yang, J., and Shen, C. (2024). On the training convergence of transformers for in-context classification of gaussian mixtures. *arXiv preprint arXiv:2410.11778*.
- Sushma, N. M., Tian, Y., Mestha, H., Colombo, N., Kappel, D., and Subramoney, A. (2024). State-space models can learn in-context by gradient descent. *arXiv preprint arXiv:2410.11687*.
- Tong, W. L. and Pehlevan, C. (2024). Mlps learn in-context on regression and classification tasks. *arXiv* preprint arXiv:2405.15618.
- Vankadara, L. C., Xu, J., Haas, M., and Cevher, V. (2024). On feature learning in structured state space models. In The Thirty-eighth Annual Conference on Neural Information Processing Systems.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., and Polosukhin, I. (2017). Attention is all you need. *Advances in neural information processing systems*, 30.
- Von Oswald, J., Niklasson, E., Randazzo, E., Sacramento, J., Mordvintsev, A., Zhmoginov, A., and Vladymyrov, M. (2023). Transformers learn in-context by gradient descent. In *International Conference on Machine Learning*, pages 35151–35174. PMLR.
- Wu, J., Zou, D., Chen, Z., Braverman, V., Gu, Q., and Bartlett, P. (2024). How many pretraining tasks are needed for in-context learning of linear regression? In *The Twelfth International Conference on Learning Representations*.

- Yang, S., Kautz, J., and Hatamizadeh, A. (2025). Gated delta networks: Improving mamba2 with delta rule. In *The Thirteenth International Conference on Learning Representations*.
- Yang, S., Wang, B., Shen, Y., Panda, R., and Kim, Y. (2024). Gated linear attention transformers with hardware-efficient training. In *International Conference on Machine Learning*, pages 56501–56523. PMLR.
- Zhang, R., Frei, S., and Bartlett, P. L. (2024). Trained transformers learn linear models in-context. *Journal of Machine Learning Research*, 25(49):1–55.
- Zhang, Y., Singh, A. K., Latham, P. E., and Saxe, A. (2025). Training dynamics of in-context learning in linear attention. *arXiv preprint arXiv:2501.16265*.
- Zheng, C., Huang, W., Wang, R., Wu, G., Zhu, J., and Li, C. (2024). On mesa-optimization in autoregressively trained transformers: Emergence and capability. *Advances in Neural Information Processing Systems*, 37:49081–49129.

NeurIPS Paper Checklist

1. Claims

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

Answer: [Yes]

Justification: The abstract reflects the paper's scope, and the introduction reflects the paper's contributions.

Guidelines:

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

2. Limitations

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

Answer: [Yes]

Justification: We discuss our limitation in the conclusion Section.

Guidelines:

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

3. Theory assumptions and proofs

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

Answer: [Yes]

Justification: We provide the assumptions in Assumption 4.1. In the main paper, we provide a proof sketch (Section 5). The detailed proofs are in the appendix.

Guidelines:

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

4. Experimental result reproducibility

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

Answer: [Yes]

Justification: The experiments follows our problem setup and assumption. We also provide the experiments setting in this paper, and the codes are in the supplementary material.

Guidelines:

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

5. Open access to data and code

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

Answer: [Yes]

Justification: The codes are in the supplementary material. We also provide a readme file. Guidelines:

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

6. Experimental setting/details

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

Answer: [Yes]

Justification: The experiments follows our problem setup and assumption. We also provide hyperparameters in the experiments setting.

Guidelines:

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

7. Experiment statistical significance

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

Answer: [Yes]

Justification: We plot the 1-sigma error bar. Figure 1b, the error bar for experimental loss can be found in Figure 2c.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).

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

8. Experiments compute resources

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

Answer: [Yes]

Justification: We report the computer resources in experiments setting. All experiments are performed on an NVIDIA A800 GPU within hours.

Guidelines:

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

9. Code of ethics

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

Answer: [Yes]

Justification: This work focuses on theoretical study of Mamba's in-context learning. All the data is synthesized. We see no ethical or potential harms of our work. We will not violate the Code Of Ethics.

Guidelines:

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

10. Broader impacts

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

Answer: [Yes]

Justification: We discuss the societal impacts in Conclusion Section. We do not see an immediate social impact.

Guidelines:

• The answer NA means that there is no societal impact of the work performed.

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

11. Safeguards

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

Answer: [NA]

Justification: This paper is mainly a theoretical work. It poses no such risks.

Guidelines:

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

12. Licenses for existing assets

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

Answer: [NA]

Justification: We perform experiments on synthetic data. It does not use existing assets.

Guidelines:

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

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

13. New assets

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

Answer: [NA]

Justification: We perform experiments on synthetic data. No new assets are introduced in the paper.

Guidelines:

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

14. Crowdsourcing and research with human subjects

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

Answer: [NA]

Justification: This paper does not involve crowdsourcing or research with human subjects. Guidelines:

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

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

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

Answer: [NA]

Justification: This paper does not involve crowdsourcing or research with human subjects. Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.

- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

16. Declaration of LLM usage

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

Answer: [NA]

Justification: This paper is mainly a theoretical work. The core method development in this research does not involve LLMs as any important, original, or non-standard components.

Guidelines:

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

Appendix

Contents

A	Basic Calculations	22
	A.1 Data Statistics	22
	A.2 Output, Loss, Gradient	24
	A.3 Training Dynamics	25
В	Proof of Theorem 4.1	27
C	Complete Proof	31
	C.1 Proof of Lemma A.2	32
	C.2 Proof of Lemma A.3	33
	C.3 Proof of Lemma A.4	36
	C.4 Proof of claim B.1	44
	C.5 Proof of claim B.2	50
	C.6 Proof of claim B.3	57
	C.7 Bounds of η^2 terms	63
D	Discussion	66
E	Additional Experimental Results	67

Table 1: Key notations

Symbols	Definitions
$oldsymbol{x}_i, oldsymbol{x}_q, oldsymbol{w}, y_i, y_q$	$m{x}_i, m{x}_q, m{w}$ are i.i.d. sampled from Gaussian distribution $\mathcal{N}(0, m{I}_d)$. $y_i = m{w}^{ op} m{x}_i, y_q = m{w}^{ op} m{x}_q.$
$ar{m{b}}_i(t),ar{m{c}}_i(t),ar{m{b}}(t)$	$egin{aligned} ar{m{b}}_i(t) &= rac{1}{\eta} ig(m{b}_i(t+1) - m{b}_i(t) ig), \ ar{m{c}}_i(t) &= rac{1}{\eta} ig(m{c}_i(t+1) - m{c}_i(t) ig), \ ar{m{b}}(t) &= rac{1}{\eta} ig(m{b}(t+1) - m{b}(t) ig). \end{aligned}$
B,C,b,c,b_i,c_i	Decompose the matrices $m{W}_B, m{W}_C$ into colums of vectors: $m{W}_B = [m{B} \ b] = [m{b}_1, \dots, m{b}_d \ b], m{W}_C = [m{C} \ c] = [m{c}_1, \dots, m{c}_d \ c]$ where $m{W}_B, m{W}_C \in \mathbb{R}^{d_h \times (d+1)}, \ m{B}, m{C} \in \mathbb{R}^{d_h \times d},$ $m{b}, m{c}, m{b}_i, m{c}_i \in \mathbb{R}^{d_h \times 1}$
$egin{aligned} oldsymbol{b}_i(t)^ op oldsymbol{b}_j(t), oldsymbol{c}_i(t)^ op oldsymbol{c}_j(t), oldsymbol{c}_i(t)^ op oldsymbol{b}_j(t), oldsymbol{c}_i(t)^ op oldsymbol{b}(t), oldsymbol{c}_i(t)^ op oldsymbol{b}(t) \end{aligned}$	inner product of the vectors $\boldsymbol{b}, \boldsymbol{b}_i, \boldsymbol{c}_i$ with $i, j \in [1, d]$. e.g. $\boldsymbol{b}_i(t)^{\top} \boldsymbol{b}_j(t)$ is the inner product of $\boldsymbol{b}_i(t)$ and $\boldsymbol{b}_i(t)$.
α	A factor, $\alpha := \exp(-\Delta_l) = \exp((-\ln 2)/N)$.
eta_1,eta_2,eta_3	The factors appearing in the gradient equation. Specifically, $\beta_1 = \left(\alpha^2 \left(1 - \alpha^N\right)^2 + \frac{(d+1)\alpha^2 (1-\alpha) \left(1 - \alpha^{2N}\right)}{(1+\alpha)}\right),$ $\beta_2 = \left(d^2 \alpha^2 \left(1 - \alpha^N\right)^2 + \frac{(2d^2 + 6d)\alpha^2 (1-\alpha) \left(1 - \alpha^{2N}\right)}{(1+\alpha)}\right),$ $\beta_3 = \alpha \left(1 - \alpha^N\right)$
γ	The lower bound of squared norms $\boldsymbol{b}_i^{\top}(t)\boldsymbol{b}_i(t), \boldsymbol{c}_i^{\top}(t)\boldsymbol{c}_i(t),$ and $\boldsymbol{b}^{\top}(t)\boldsymbol{b}(t).$ Specifically, $\gamma = \frac{1}{2}d_h.$
$\delta(T)$	The upper bound of cross-interactions: $\boldsymbol{b}_i^\top(t)\boldsymbol{b}_j(t), \boldsymbol{c}_i^\top(t)\boldsymbol{c}_j(t),$ and $\boldsymbol{b}_i^\top(t)\boldsymbol{b}(t).$ Specifically, $\delta(t) = \max_{s \in [0,t]} \{2\sqrt{d_h \log(4d(2d+1)/\delta)}, \boldsymbol{b}_i^\top(s)\boldsymbol{b}_j(s) , \boldsymbol{c}_i^\top(s)\boldsymbol{c}_j(s) , \boldsymbol{b}_i^\top(s)\boldsymbol{b}(s) \}.$

A Basic Calculations

This Section provide the data statistics related to gaussian distribution, and compute the expressions for the output, loss, gradient, training dynamics (particularly *Vector-coupled Dynamics*) of the Mamba model. Section B presents the *Fine-grained Induction* with *Negative Feedback Convergence* technique, and finally establish the results for Theorem 4.1. Section C details the complete proofs for Section A and Section B. In Section D, we discuss about orthogonal initialization and compare our framework with other techniques. In Section E, we give more experimental results.

A.1 Data Statistics

Lemma A.1 (Concentration Inequalities) Let $b_i(0)$ be the i-th colum of B(0), $c_i(0)$ be the i-th colum of C(0), and suppose that $\delta > 0$ and $d_h = \Omega(\log(4(2d+1)/\delta))$, with probability at least $1 - \delta$, we have:

$$\frac{3d_h}{4} \le \boldsymbol{b}_i(0)^{\top} \boldsymbol{b}_i(0), \boldsymbol{c}_i(0)^{\top} \boldsymbol{c}_i(0), \boldsymbol{b}(0)^{\top} \boldsymbol{b}(0) \le \frac{5d_h}{4},$$

$$\left| \boldsymbol{c}_i(0)^{\top} \boldsymbol{b}_i(0) \right|, \left| \boldsymbol{c}_i(0)^{\top} \boldsymbol{b}_j(0) \right|, \left| \boldsymbol{c}_i(0)^{\top} \boldsymbol{b}(0) \right| \le 2\sqrt{d_h \log(4d(2d+1)/\delta)},$$

$$\left| \boldsymbol{b}_i(0)^{\top} \boldsymbol{b}_j(0) \right|, \left| \boldsymbol{c}_i(0)^{\top} \boldsymbol{c}_j(0) \right|, \left| \boldsymbol{b}_i(0)^{\top} \boldsymbol{b}(0) \right| \le 2\sqrt{d_h \log(4d(2d+1)/\delta)},$$

for $i, j \in [d], i \neq j$.

Proof of Lemma A.1. By Bernstein's inequality, with probability at least $1 - \delta/2(2d+1)$ we have

$$|\boldsymbol{b}_i(0)^{\top} \boldsymbol{b}_i(0) - d_h| = O(\sqrt{d_h \log(4(2d+1)/\delta)}).$$

Therefore, as long as $d_h = \Omega(\log(4(2d+1)/\delta))$, we have $3d_h/4 \leq \boldsymbol{b}_i(0)^\top \boldsymbol{b}_i(0) \leq 5d_h/4$. Similarly, we have

$$\frac{3d_h}{4} \leq \boldsymbol{c}_i(0)^{\top} \boldsymbol{c}_i(0), \boldsymbol{b}(0)^{\top} \boldsymbol{b}(0) \leq \frac{5d_h}{4}.$$

For $i, j \in [d], i \neq j$, By Bernstein's inequality, with probability at least $1 - \delta/2d(2d + 1)$, we have

$$\begin{aligned} & \left| \boldsymbol{c}_i(0)^{\top} \boldsymbol{b}_i(0) \right|, \left| \boldsymbol{c}_i(0)^{\top} \boldsymbol{b}_j(0) \right|, \left| \boldsymbol{c}_i(0)^{\top} \boldsymbol{b}(0) \right| \leq 2\sqrt{d_h \log(4d(2d+1)/\delta)}, \\ & \left| \boldsymbol{b}_i(0)^{\top} \boldsymbol{b}_j(0) \right|, \left| \boldsymbol{c}_i(0)^{\top} \boldsymbol{c}_j(0) \right|, \left| \boldsymbol{b}_i(0)^{\top} \boldsymbol{b}(0) \right| \leq 2\sqrt{d_h \log(4d(2d+1)/\delta)}. \end{aligned}$$

We can apply a union bound to complete the proof.

Lemma A.2 If vectors \mathbf{x} and \mathbf{w} are iid generated from $\mathcal{N}(0, \mathbf{I}_d)$, $y = \mathbf{x}^{\top} \mathbf{w}$ we have the following expectations:

$$\mathbb{E}\left[\boldsymbol{x}\boldsymbol{x}^{\top}\boldsymbol{w}\boldsymbol{w}^{\top}\boldsymbol{x}\boldsymbol{x}^{\top}\right] = (d+2)\boldsymbol{I},$$
$$\mathbb{E}\left[y^{2}\right] = d,$$
$$\mathbb{E}\left[y^{4}\right] = 3d(d+2).$$

The proof of lemma A.2 is in Section C.1.

Lemma A.3 If vectors \mathbf{x}_i and \mathbf{w} are iid generated from $\mathcal{N}(0, \mathbf{I}_d)$, $y = \mathbf{x}_i^{\top} \mathbf{w}$ we have the following expectations:

$$\begin{split} \mathbb{E}\Big[\sum_{i=0}^{N-1}\sum_{j=0}^{N-1}\alpha^{i+j+2}y_{N-i}y_{N-j}x_{N-i}x_{N-j}^{\top}\Big] &= \Big(\frac{\alpha^2\big(1-\alpha^N\big)^2}{(1-\alpha)^2} + \frac{(d+1)\alpha^2\big(1-\alpha^{2N}\big)}{(1-\alpha)(1+\alpha)}\Big) \cdot \boldsymbol{I}, \\ \mathbb{E}\Big[\sum_{i=0}^{N-1}\sum_{j=0}^{N-1}\alpha^{i+j+2}y_{N-i}y_{N-j}^2\Big] &= 0, \\ \mathbb{E}\Big[\sum_{i=0}^{N-1}\sum_{j=0}^{N-1}\alpha^{i+j+2}y_{N-i}^2y_{N-j}^2\Big] &= \frac{d^2\alpha^2\Big(1-\alpha^N\big)^2}{(1-\alpha)^2} + \frac{(2d^2+6d)\alpha^2\Big(1-\alpha^{2N}\big)}{(1-\alpha)(1+\alpha)}, \\ \mathbb{E}\Big[\sum_{i=0}^{N-1}\alpha^{i+j+2}y_{N-i}w^{\top}\Big] &= \alpha\Big(\frac{1-\alpha^N}{1-\alpha}\Big) \cdot \boldsymbol{I}, \\ \mathbb{E}\Big[\sum_{i=0}^{N-1}\alpha^{i+1}y_{N-i}x_{N-i}w^{\top}\Big] &= 0, \\ \mathbb{E}\Big[\sum_{i=0}^{N-1}\sum_{j=0}^{N-1}\alpha^{i+j+2}x_{N-i}\underbrace{x_{N-i}^{\top}w}_{y_{N-i}}\underbrace{x_{N-j}^{\top}w}_{y_{N-j}}\Big] &= 0, \\ \mathbb{E}\Big[\sum_{i=0}^{N-1}\sum_{j=0}^{N-1}\alpha^{i+j+2}y_{N-i}y_{N-j}\Big] &= 0, \\ \mathbb{E}\Big[\sum_{i=0}^{N-1}\sum_{j=0}^{N-1}\alpha^{i+j+2}y_{N-i}y_{N-j}\Big] &= \frac{d\alpha^2\Big(1-\alpha^{2N}\Big)}{(1-\alpha)(1+\alpha)}, \\ \mathbb{E}\Big[\sum_{i=0}^{N-1}\alpha^{i+j+2}y_{N-i}y_{N-j}\Big] &= 0. \end{split}$$

The proof of lemma A.3 is in Section C.2.

A.2 Output, Loss, Gradient

This section we derive the output of Mamba given sequence $\{e_{1:N}, e_q\}$, and establishe the loss function formulation with its gradient expression.

Linear Recurrence. To start with, we show how the hidden states update when receiving token $e_l = (\boldsymbol{x}_l^\top, y_l)^\top$. By (Eq. (5)) and Assumption 4.1(2), we have $\Delta_l = \operatorname{softplus}(\ln(\exp((\ln 2)/N) - 1)) = (\ln 2)/N$. Combining it with (Eq. (1)(2)) and get:

$$\mathbf{h}_{l}^{(d+1)} = \overline{\mathbf{A}}_{l} \mathbf{h}_{l-1}^{(d+1)} + \overline{\mathbf{B}}_{l} y_{l}
= \exp(\Delta_{l} \mathbf{A}) \mathbf{h}_{l-1}^{(d+1)} + y_{l} (\Delta_{l} \mathbf{A})^{-1} (\exp(\Delta_{l} \mathbf{A}) - \mathbf{I}) \Delta_{l} \mathbf{B}_{l}
= \exp(-\Delta_{l}) \mathbf{I} \mathbf{h}_{l-1}^{(d+1)} - y_{l} \Delta_{l}^{-1} (\exp(-\Delta_{l}) \mathbf{I} - \mathbf{I}) \Delta_{l} \mathbf{B}_{l}
= \exp(-\Delta_{l}) \mathbf{h}_{l-1}^{(d+1)} + (1 - \exp(-\Delta_{l})) y_{l} \mathbf{B}_{l}
= \alpha \mathbf{h}_{l-1}^{(d+1)} + (1 - \alpha) y_{l} \mathbf{B}_{l}$$
(18)

where $\alpha := \exp(-\Delta_l) = \exp((-\ln 2)/N)$, the second equality is by discretization rule (2), the third equality is by Assumption 4.1(2) and $exp(-\Delta_l \mathbf{I}) = exp(-\Delta_l)\mathbf{I}$. (Eq. (18)) is similar to theorem 1 in Gu and Dao (2024)

Prediction Output. We next derive the expression of \hat{y}_q . Based on (Eq.(11)), the hidden state after receiving the first l context prompts $e_{1:l}$ is given by:

$$\mathbf{h}_{l}^{(d+1)} = \alpha \mathbf{h}_{l-1}^{(d+1)} + (1 - \alpha) y_{l} \mathbf{B}_{l}$$

$$= \alpha^{2} \mathbf{h}_{l-2}^{(d+1)} + (1 - \alpha) y_{l} \mathbf{B}_{l} + (1 - \alpha) \alpha y_{l-1} \mathbf{B}_{l-1}$$

$$= \dots$$

$$= (1 - \alpha) \sum_{i=0}^{l-1} \alpha^{i} y_{l-i} \mathbf{B}_{l-i}$$
(19)

Receiving the query token $\boldsymbol{e}_q = (\boldsymbol{x}_q^\top, 0)^\top$, we have:

$$\boldsymbol{h}_{N+1}^{(d+1)} = \alpha \boldsymbol{h}_{N}^{(d+1)} + (1 - \alpha) \cdot 0 \cdot \boldsymbol{B}_{N} = (1 - \alpha) \sum_{i=0}^{N-1} \alpha^{i+1} y_{N-i} \boldsymbol{B}_{N-i}$$
 (20)

Finally, the prediction output is as follows

$$\hat{y}_{q} = \boldsymbol{C}_{N+1}^{\top} \boldsymbol{h}_{N+1}^{(d+1)} = (1 - \alpha) \boldsymbol{C}_{N+1}^{\top} \sum_{i=0}^{N-1} \alpha^{i+1} y_{N-i} \boldsymbol{B}_{N-i}$$

$$= (1 - \alpha) (\boldsymbol{W}_{C} \boldsymbol{e}_{q} + \boldsymbol{b}_{C})^{\top} \sum_{i=0}^{N-1} \alpha^{i+1} y_{N-i} (\boldsymbol{W}_{B} \boldsymbol{e}_{N-i} + \boldsymbol{b}_{B})$$
(21)

To handle $W_C e_q$ and $W_B e_{N-i}$, we further denote $W_B = [B b]$ and $W_C = [C c]$, where $B, C \in \mathbb{R}^{d_h \times d}$, $b, c \in \mathbb{R}^{d_h \times 1}$. Then we write another form of (Eq. (21)):

$$\hat{y}_{q} = (1 - \alpha)(Cx_{q} + b_{C})^{\top} \sum_{i=0}^{N-1} \alpha^{i+1} y_{N-i}(Bx_{N-i} + y_{N-i}b + b_{B})$$
(22)

The loss becomes:

$$\mathcal{L}(\boldsymbol{\theta}) = \frac{1}{2} \mathbb{E} \left[\left((1 - \alpha) (\boldsymbol{C} \boldsymbol{x}_q + \boldsymbol{b}_C)^\top \sum_{i=0}^{N-1} \alpha^{i+1} y_{N-i} (\boldsymbol{B} \boldsymbol{x}_{N-i} + y_{N-i} \boldsymbol{b} + \boldsymbol{b}_B) - \boldsymbol{w}^\top \boldsymbol{x}_q \right)^2 \right]$$
(23)

The following lemma provides the gradient of B, C, b, b_B, b_C with respect to loss (Eq. (23)).

Lemma A.4 (Gradient) The gradient of trainable parameters $\theta' = \{B, C, b, b_B, b_C\}$ with respect to loss (Eq. (23)) are as follows:

$$\nabla_{\boldsymbol{b}_{B}}\mathcal{L}(\boldsymbol{\theta}) = \mathbf{0},$$

$$\nabla_{\boldsymbol{b}_{C}}\mathcal{L}(\boldsymbol{\theta}) = \mathbf{0},$$

$$\nabla_{\boldsymbol{b}_{C}}\mathcal{L}(\boldsymbol{\theta}) = \mathbf{0},$$

$$\nabla_{\boldsymbol{B}}\mathcal{L}(\boldsymbol{\theta}) = \underbrace{\left(\alpha^{2}(1-\alpha^{N})^{2} + \frac{(d+1)\alpha^{2}(1-\alpha)(1-\alpha^{2N})}{(1+\alpha)}\right)}_{:=\beta_{1}}\boldsymbol{C}\boldsymbol{C}^{\top}\boldsymbol{B} - \underbrace{\alpha(1-\alpha^{N})}_{:=\beta_{3}}\boldsymbol{C},$$

$$\nabla_{\boldsymbol{b}}\mathcal{L}(\boldsymbol{\theta}) = \underbrace{\left(d^{2}\alpha^{2}(1-\alpha^{N})^{2} + \frac{(2d^{2}+6d)\alpha^{2}(1-\alpha)(1-\alpha^{2N})}{(1+\alpha)}\right)}_{:=\beta_{2}}\boldsymbol{C}\boldsymbol{C}^{\top}\boldsymbol{b},$$

$$\nabla_{\boldsymbol{C}}\mathcal{L}(\boldsymbol{\theta}) = \underbrace{\left(\alpha^{2}(1-\alpha^{N})^{2} + \frac{(d+1)\alpha^{2}(1-\alpha)(1-\alpha^{2N})}{(1+\alpha)}\right)}_{:=\beta_{1}}\boldsymbol{B}\boldsymbol{B}^{\top}\boldsymbol{C}$$

$$+ \underbrace{\left(d^{2}\alpha^{2}(1-\alpha^{N})^{2} + \frac{(2d^{2}+6d)\alpha^{2}(1-\alpha)(1-\alpha^{2N})}{(1+\alpha)}\right)}_{:=\beta_{2}}\boldsymbol{b}\boldsymbol{b}^{\top}\boldsymbol{C}$$

$$= \underbrace{\alpha(1-\alpha^{N})}_{:=\beta_{2}}\boldsymbol{B}.$$

Here, we denote $\beta_1 = \left(\alpha^2 \left(1 - \alpha^N\right)^2 + \frac{(d+1)\alpha^2 (1-\alpha)\left(1-\alpha^{2N}\right)}{(1+\alpha)}\right)$, $\beta_2 = \left(d^2\alpha^2 \left(1 - \alpha^N\right)^2 + \frac{(2d^2+6d)\alpha^2 (1-\alpha)\left(1-\alpha^{2N}\right)}{(1+\alpha)}\right)$, $\beta_3 = \alpha \left(1 - \alpha^N\right)$ for simplicity. The proof of lemma A.4 is in Section C.3.

A.3 Training Dynamics

With the gradient in lemma A.4, we further provide the update rule for Mamba's parameters and the *Vector-coupled Dynamics*.

Using gradient descent algorithm $\theta'(t+1) = \theta'(t) - \eta \nabla_{\theta'} \mathcal{L}(\theta(t))$ with training rate η , we have the following update rule base on lemma A.4.

Lemma A.5 (Update Rule, restatement of lemma 5.1) Let η be the learning rate and we use gradient descent to update the weights W_B , W_C , b_B , b_C , for $t \ge 0$ we have

$$\boldsymbol{B}(t+1) = \boldsymbol{B}(t) + \eta \beta_3 \boldsymbol{C}(t) - \eta \beta_1 \boldsymbol{C}(t) \boldsymbol{C}(t)^{\top} \boldsymbol{B}(t),$$

$$\boldsymbol{C}(t+1) = \boldsymbol{C}(t) + \eta \beta_3 \boldsymbol{B}(t) - \eta \beta_1 \boldsymbol{B}(t) \boldsymbol{B}(t)^{\top} \boldsymbol{C}(t) - \eta \beta_2 \boldsymbol{b}(t) \boldsymbol{b}(t)^{\top} \boldsymbol{C}(t),$$

$$\boldsymbol{b}(t+1) = \boldsymbol{b}(t) - \eta \beta_2 \boldsymbol{C}(t) \boldsymbol{C}(t)^{\top} \boldsymbol{b}(t),$$

$$\boldsymbol{b}_B(t) = \boldsymbol{b}_C(t) = \boldsymbol{0}.$$

We decompose B and C as $B = [b_1 \dots b_d]$, $C = [c_1 \dots c_d]$, and provide the update rule for b_i , c_i and b with $i \in [1:d]$ as the following lemma.

Lemma A.6 (Vectors Update Rule, restatement of lemma 5.2) *Let* η *be the learning rate and we use gradient descent to update the weights* W_B , W_C , b_B , b_C , for $i \in [d]$, $t \ge 0$ we have

$$\boldsymbol{b}_i(t+1) = \boldsymbol{b}_i(t) + \eta \Big(\big(\beta_3 - \beta_1 \boldsymbol{c}_i^\top(t) \boldsymbol{b}_i(t) \big) \boldsymbol{c}_i(t) - \beta_1 \sum_{k \neq i}^d \boldsymbol{c}_k^\top(t) \boldsymbol{b}_i(t) \cdot \boldsymbol{c}_k(t) \Big)$$

$$=: \boldsymbol{b}_i(t) + \eta \bar{\boldsymbol{b}}_i(t)$$

$$\boldsymbol{c}_{i}(t+1) = \boldsymbol{c}_{i}(t) + \eta \Big((\beta_{3} - \beta_{1}\boldsymbol{c}_{i}^{\top}(t)\boldsymbol{b}_{i}(t))\boldsymbol{b}_{i}(t) - \beta_{1} \sum_{k \neq i}^{d} \boldsymbol{c}_{i}^{\top}(t)\boldsymbol{b}_{k}(t) \cdot \boldsymbol{b}_{k}(t) - \beta_{2}\boldsymbol{c}_{i}^{\top}(t)\boldsymbol{b}(t) \cdot \boldsymbol{b}(t) \Big)$$

$$=: \boldsymbol{c}_{i}(t) + \eta \bar{\boldsymbol{c}}_{i}(t)$$

$$\boldsymbol{b}(t+1) = \boldsymbol{b}(t) - \eta \Big(\beta_2 \sum_{k=1}^d \boldsymbol{c}_k^\top(t) \boldsymbol{b}(t) \cdot \boldsymbol{c}_k(t) \Big) =: \boldsymbol{b}(t) + \eta \bar{\boldsymbol{b}}(t)$$

Here, we denote $\eta \bar{b}_i(t) = b_i(t+1) - b_i(t)$, $\eta \bar{c}_i(t) = c_i(t+1) - c_i(t)$, and $\eta \bar{b}(t) = b(t+1) - b(t)$ for simplicity.

Next, we provide the dynamics for the inner products of these vectors.

Lemma A.7 (Vector-coupled Dynamics) Let η be the learning rate and we use gradient descent to update the weights W_B , W_C , b_B , b_C , we have

$$\begin{aligned} & \boldsymbol{b}_{i}^{\top}(t+1)\boldsymbol{b}_{i}(t+1) \\ &= \boldsymbol{b}_{i}^{\top}(t)\boldsymbol{b}_{i}(t) + 2\eta\Big(\big(\beta_{3} - \beta_{1}\boldsymbol{c}_{i}^{\top}(t)\boldsymbol{b}_{i}(t)\big)\boldsymbol{c}_{i}^{\top}(t)\boldsymbol{b}_{i}(t) - \beta_{1}\sum_{k\neq i}^{d}\big(\boldsymbol{c}_{k}^{\top}(t)\boldsymbol{b}_{i}(t)\big)^{2}\Big) \\ &+ \eta^{2} \left\| \bar{\boldsymbol{b}}_{i}(t) \right\|_{2}^{2} \\ & \boldsymbol{b}_{i}^{\top}(t+1)\boldsymbol{b}_{j}(t+1) \\ &= \boldsymbol{b}_{i}^{\top}(t)\boldsymbol{b}_{j}(t) + \eta\Big(2\big(\beta_{3} - \beta_{1}\boldsymbol{c}_{i}^{\top}(t)\boldsymbol{b}_{i}(t)\big)\boldsymbol{c}_{i}^{\top}(t)\boldsymbol{b}_{j}(t) + 2\big(\beta_{3} - \beta_{1}\boldsymbol{c}_{j}^{\top}(t)\boldsymbol{b}_{j}(t)\big)\boldsymbol{c}_{j}^{\top}(t)\boldsymbol{b}_{i}(t) \\ &- \beta_{3}\big(\boldsymbol{c}_{i}^{\top}(t)\boldsymbol{b}_{j}(t) + \boldsymbol{c}_{j}^{\top}(t)\boldsymbol{b}_{i}(t)\big) - 2\beta_{1}\sum_{k\neq i,k\neq j}^{d}\boldsymbol{c}_{i}^{\top}(t)\boldsymbol{b}_{i}(t) \cdot \boldsymbol{c}_{k}^{\top}(t)\boldsymbol{b}_{j}(t)\Big) + \eta^{2}\bar{\boldsymbol{b}}_{i}^{\top}(t)\bar{\boldsymbol{b}}_{j}(t) \\ &\boldsymbol{c}_{i}^{\top}(t+1)\boldsymbol{c}_{i}(t+1) \\ &= \boldsymbol{c}_{i}^{\top}(t)\boldsymbol{c}_{i}(t) + 2\eta\Big(\big(\beta_{3} - \beta_{1}\boldsymbol{c}_{i}^{\top}(t)\boldsymbol{b}_{i}(t)\big)\boldsymbol{c}_{i}^{\top}(t)\boldsymbol{b}_{i}(t) - \beta_{1}\sum_{k\neq i}^{d}\big(\boldsymbol{c}_{i}^{\top}(t)\boldsymbol{b}_{k}(t)\big)^{2} \\ &- \beta_{2}\big(\boldsymbol{c}_{i}^{\top}(t)\boldsymbol{b}(t)\big)^{2}\big) + \eta^{2}\left\|\bar{\boldsymbol{c}}_{i}(t)\right\|_{2}^{2} \\ &\boldsymbol{c}_{i}^{\top}(t+1)\boldsymbol{c}_{j}(t+1) \\ &= \boldsymbol{c}_{i}^{\top}(t)\boldsymbol{c}_{j}(t) + \eta\Big(2\big(\beta_{3} - \beta_{1}\boldsymbol{c}_{i}^{\top}(t)\boldsymbol{b}_{i}(t)\big)\boldsymbol{c}_{j}^{\top}(t)\boldsymbol{b}_{i}(t) + 2\big(\beta_{3} - \beta_{1}\boldsymbol{c}_{j}^{\top}(t)\boldsymbol{b}_{j}(t)\big)\boldsymbol{c}_{i}^{\top}(t)\boldsymbol{b}_{i}(t) \\ &- \beta_{3}\big(\boldsymbol{c}_{i}^{\top}(t)\boldsymbol{b}_{j}(t) + \boldsymbol{c}_{j}^{\top}(t)\boldsymbol{b}_{i}(t)\big) + \eta^{2}\bar{\boldsymbol{c}}_{i}^{\top}(t)\bar{\boldsymbol{c}}_{j}(t) \\ &\boldsymbol{c}_{i}^{\top}(t+1)\boldsymbol{b}_{i}(t+1) \\ &= \boldsymbol{c}_{i}^{\top}(t)\boldsymbol{b}(t) \cdot \boldsymbol{c}_{j}^{\top}(t)\boldsymbol{b}(t)\Big) + \eta^{2}\bar{\boldsymbol{c}}_{i}^{\top}(t)\bar{\boldsymbol{c}}_{j}(t) \\ &\boldsymbol{c}_{i}^{\top}(t+1)\boldsymbol{b}_{i}(t+1) \\ &= \boldsymbol{c}_{i}^{\top}(t)\boldsymbol{b}_{i}(t) + \eta\Big(\big(\beta_{3} - \beta_{1}\boldsymbol{c}_{i}^{\top}(t)\boldsymbol{b}_{i}(t)\big)\boldsymbol{b}_{i}^{\top}(t)\boldsymbol{b}_{i}(t) - \beta_{1}\sum_{k\neq i}^{d}\boldsymbol{c}_{i}^{\top}(t)\boldsymbol{b}_{k}(t) \cdot \boldsymbol{b}_{k}^{\top}(t)\boldsymbol{b}_{i}(t) \\ &- \beta_{2}\boldsymbol{c}_{i}^{\top}(t)\boldsymbol{b}(t) \cdot \boldsymbol{b}_{i}^{\top}(t)\boldsymbol{b}(t) + \big(\beta_{3} - \beta_{1}\boldsymbol{c}_{i}^{\top}(t)\boldsymbol{b}_{i}(t)\big)\boldsymbol{c}_{i}^{\top}(t)\boldsymbol{b}_{i}(t) \\ &- \beta_{2}\boldsymbol{c}_{i}^{\top}(t)\boldsymbol{b}(t) \cdot \boldsymbol{b}_{i}^{\top}(t)\boldsymbol{b}(t) + \big(\beta_{3} - \beta_{1}\boldsymbol{c}_{i}^{\top}(t)\boldsymbol{b}_{i}(t)\big)\boldsymbol{c}_{i}^{\top}(t)\boldsymbol{c}_{i}(t) \end{aligned}$$

$$\begin{split} &-\beta_1 \sum_{k \neq i}^{d} \boldsymbol{c}_k^\top(t) \boldsymbol{b}_i(t) \cdot \boldsymbol{c}_i^\top(t) \boldsymbol{c}_k(t) \Big) + \eta^2 \bar{\boldsymbol{c}}_i^\top(t) \bar{\boldsymbol{b}}_i(t) \\ & \boldsymbol{c}_i^\top(t+1) \boldsymbol{b}_j(t+1) \\ &= \Big(1 - \eta \beta_1 \big(\boldsymbol{c}_i^\top(t) \boldsymbol{c}_i(t) + \boldsymbol{b}_j^\top(t) \boldsymbol{b}_j(t) \big) \Big) \boldsymbol{c}_i^\top(t) \boldsymbol{b}_j(t) \\ &+ \eta \Big(\big(\beta_3 - \beta_1 \boldsymbol{c}_i^\top(t) \boldsymbol{b}_i(t) \big) \boldsymbol{b}_i^\top(t) \boldsymbol{b}_j(t) - \beta_1 \sum_{k \neq i, k \neq j}^{d} \boldsymbol{c}_i^\top(t) \boldsymbol{b}_k(t) \cdot \boldsymbol{b}_k^\top(t) \boldsymbol{b}_j(t) \\ &- \beta_2 \boldsymbol{c}_i^\top(t) \boldsymbol{b}(t) \cdot \boldsymbol{b}_j^\top(t) \boldsymbol{b}(t) + \big(\beta_3 - \beta_1 \boldsymbol{c}_j^\top(t) \boldsymbol{b}_j(t) \big) \boldsymbol{c}_i^\top(t) \bar{\boldsymbol{b}}_j(t) \\ &- \beta_1 \sum_{k \neq i, k \neq j}^{d} \boldsymbol{c}_k^\top(t) \boldsymbol{b}_j(t) \cdot \boldsymbol{c}_i^\top(t) \boldsymbol{c}_k(t) \Big) + \eta^2 \bar{\boldsymbol{c}}_i^\top(t) \bar{\boldsymbol{b}}_j(t) \\ &\boldsymbol{b}^\top(t+1) \boldsymbol{b}(t+1) = \boldsymbol{b}^\top(t) \boldsymbol{b}(t) - 2\eta \Big(\beta_2 \sum_{k=1}^{d} \big(\boldsymbol{c}_k^\top(t) \boldsymbol{b}(t)\big)^2 \Big) + \eta^2 \Big\| \bar{\boldsymbol{b}}(t) \Big\|_2^2 \\ &\boldsymbol{b}_i^\top(t+1) \boldsymbol{b}(t+1) \\ &= \boldsymbol{b}_i^\top(t) \boldsymbol{b}(t) + \eta \Big(\big(\beta_3 - \beta_1 \boldsymbol{c}_i^\top(t) \boldsymbol{b}_i(t) \big) \boldsymbol{c}_i^\top(t) \boldsymbol{b}(t) - \beta_1 \sum_{k \neq i}^{d} \boldsymbol{c}_k^\top(t) \boldsymbol{b}_i(t) \cdot \boldsymbol{c}_k^\top(t) \boldsymbol{b}(t) \\ &- \beta_2 \sum_{k=1}^{d} \boldsymbol{c}_k^\top(t) \boldsymbol{b}(t) \cdot \boldsymbol{c}_k^\top(t) \boldsymbol{b}_i(t) \Big) + \eta^2 \bar{\boldsymbol{b}}_i^\top(t) \bar{\boldsymbol{b}}(t) \\ &= (1 - \eta \beta_2 \big(\boldsymbol{b}^\top(t) \boldsymbol{b}(t) + \boldsymbol{c}_i^\top(t) \boldsymbol{c}_i(t) \big) \Big) \boldsymbol{c}_i^\top(t) \boldsymbol{b}(t) \\ &+ \eta \Big(\big(\beta_3 - \beta_1 \boldsymbol{c}_i^\top(t) \boldsymbol{b}_i(t) \big) \boldsymbol{b}_i^\top(t) \boldsymbol{b}(t) - \beta_1 \sum_{k \neq i}^{d} \boldsymbol{c}_i^\top(t) \boldsymbol{b}_k(t) \cdot \boldsymbol{b}_k^\top(t) \boldsymbol{b}(t) \\ &- \beta_2 \sum_{k \neq i}^{d} \boldsymbol{c}_k^\top(t) \boldsymbol{b}(t) \cdot \boldsymbol{c}_k^\top(t) \boldsymbol{b}_i(t) \Big) \boldsymbol{b}_i^\top(t) \boldsymbol{b}(t) \\ &- \beta_2 \sum_{k \neq i}^{d} \boldsymbol{c}_k^\top(t) \boldsymbol{b}(t) \cdot \boldsymbol{c}_k^\top(t) \boldsymbol{b}_i(t) \Big) \boldsymbol{b}_i^\top(t) \boldsymbol{b}(t) - \beta_1 \sum_{k \neq i}^{d} \boldsymbol{c}_i^\top(t) \boldsymbol{b}_k(t) \cdot \boldsymbol{b}_k^\top(t) \boldsymbol{b}(t) \\ &- \beta_2 \sum_{k \neq i}^{d} \boldsymbol{c}_k^\top(t) \boldsymbol{b}(t) \cdot \boldsymbol{c}_k^\top(t) \boldsymbol{c}_i(t) \Big) \boldsymbol{b}_i^\top(t) \boldsymbol{b}(t) - \beta_1 \sum_{k \neq i}^{d} \boldsymbol{c}_i^\top(t) \boldsymbol{b}_k(t) \cdot \boldsymbol{b}_k^\top(t) \boldsymbol{b}(t) \\ &- \beta_2 \sum_{k \neq i}^{d} \boldsymbol{c}_k^\top(t) \boldsymbol{b}(t) \cdot \boldsymbol{c}_k^\top(t) \boldsymbol{c}_i(t) \Big) \boldsymbol{b}_i^\top(t) \boldsymbol{b}(t) - \beta_1 \sum_{k \neq i}^{d} \boldsymbol{c}_i^\top(t) \boldsymbol{b}_k(t) \cdot \boldsymbol{b}_k^\top(t) \boldsymbol{b}(t) \\ &- \beta_2 \sum_{k \neq i}^{d} \boldsymbol{c}_k^\top(t) \boldsymbol{b}(t) \cdot \boldsymbol{c}_k^\top(t) \boldsymbol{c}_i(t) \Big) \boldsymbol{b}_i^\top(t) \boldsymbol{b}(t) \\ &- \beta_1 \sum_{k \neq i}^{d} \boldsymbol{c}_k^\top(t) \boldsymbol{b}(t) \cdot \boldsymbol{c}_k^\top(t) \boldsymbol{b}(t) \cdot \boldsymbol{c}_k^\top(t) \boldsymbol{b}(t) \Big) \boldsymbol{c}_i^\top(t) \boldsymbol{b}(t) \\ &- \beta_1 \sum_{k \neq i}^{d} \boldsymbol{c}_k^\top(t) \boldsymbol{b}(t) \cdot \boldsymbol{c}_k^\top(t) \boldsymbol{c}_i(t) \Big) \boldsymbol$$

Lemma A.7 is derive by calculating the inner products of the vectors update rule in lemma A.6. For example, $\mathbf{b}^{\top}(t+1)\mathbf{b}(t+1)$ is derived as follow:

$$\begin{aligned} \boldsymbol{b}^{\top}(t+1)\boldsymbol{b}(t+1) &= \left(\boldsymbol{b}(t) + \eta \bar{\boldsymbol{b}}(t)\right)^{\top} \left(\boldsymbol{b}(t) + \eta \bar{\boldsymbol{b}}(t)\right) \\ &= \boldsymbol{b}^{\top}(t)\boldsymbol{b}(t) - 2\eta \bar{\boldsymbol{b}}(t)^{\top}\boldsymbol{b}(t) + \eta^{2} \left\|\bar{\boldsymbol{b}}(t)\right\|_{2}^{2} \\ &= \boldsymbol{b}^{\top}(t)\boldsymbol{b}(t) - 2\eta \left(\beta_{2} \sum_{k=1}^{d} \left(\boldsymbol{c}_{k}^{\top}(t)\boldsymbol{b}(t)\right)^{2}\right) + \eta^{2} \left\|\bar{\boldsymbol{b}}(t)\right\|_{2}^{2} \end{aligned}$$

The other equations are similar to it.

B Proof of Theorem 4.1

In this section, we present the framework of *Fine-grained Induction*, and establish the results of Theorem 4.1 after convergence.

Fine-gained Induction Specifically, denoting $\delta(t) = \max_{s \in [0,t]} \{|\boldsymbol{b}_i^\top(s)\boldsymbol{b}_j(s)|, |\boldsymbol{c}_i^\top(s)\boldsymbol{c}_j(s)|, |\boldsymbol{b}_i^\top(s)\boldsymbol{b}(s)|\}$ and $\gamma = \min_{t \geq 0} \{\boldsymbol{b}_i^\top(t)\boldsymbol{b}_i(t), \boldsymbol{c}_i^\top(t)\boldsymbol{c}_i(t), \boldsymbol{b}^\top(t)\boldsymbol{b}(t)\}$, we establish the following three properties $\mathcal{A}(t)$, $\mathcal{B}(t)$, and $\mathcal{C}(t)$ simultaneously for $t \geq 0$:

•
$$\mathcal{A}(t)$$
:
$$d_h/2 \leq \boldsymbol{b}_i^\top(t)\boldsymbol{b}_i(t), \boldsymbol{c}_i^\top(t)\boldsymbol{c}_i(t), \boldsymbol{b}^\top(t)\boldsymbol{b}(t) \leq 2d_h$$

•
$$\mathcal{B}(t)$$
:
$$\begin{aligned} |\beta_3 - \beta_1 \boldsymbol{c}_i^\top(t) \boldsymbol{b}_i(t)| &\leq \delta(t) \exp(-\eta \beta_1 \gamma t) \\ |\boldsymbol{c}_i^\top(t) \boldsymbol{b}_j(t)| &\leq 2\delta(t) \exp(-\eta \beta_1 \gamma t) \end{aligned}$$

$$|\boldsymbol{c}_i^\top(t) \boldsymbol{b}(t)| &\leq 2\delta(t) \exp(-\eta \beta_2 \gamma t) + \frac{\delta(t)}{\beta_2} \exp(-\eta \beta_1 \gamma t)$$

• C(t):

$$|\boldsymbol{b}_{i}^{\top}(t)\boldsymbol{b}_{i}(t)|, |\boldsymbol{c}_{i}^{\top}(t)\boldsymbol{c}_{i}(t)|, |\boldsymbol{b}_{i}^{\top}(t)\boldsymbol{b}(t)| \leq \delta(t) \leq 3\sqrt{d_{h}\log(4d(2d+1)/\delta)} =: \delta_{\max}$$

Here, $i, j \in [1, d], i \neq j$. The initial conditions $\mathcal{A}(0)$, $\mathcal{B}(0)$, and $\mathcal{C}(0)$ are established with high probability by concentration inequalities (lemma A.1). We also provide the following claims to establish $\mathcal{A}(t)$, $\mathcal{B}(t)$, and $\mathcal{C}(t)$ for $t \geq 0$:

Claim B.1
$$\mathcal{A}(0), \dots, \mathcal{A}(T), \mathcal{B}(0), \dots, \mathcal{B}(T), \mathcal{C}(0), \dots, \mathcal{C}(T) \Longrightarrow \mathcal{A}(T+1)$$

Claim B.2
$$\mathcal{A}(0), \dots, \mathcal{A}(T), \mathcal{B}(0), \dots, \mathcal{B}(T), \mathcal{C}(0), \dots, \mathcal{C}(T) \Longrightarrow \mathcal{B}(T+1)$$

Claim B.3
$$\mathcal{A}(0), \dots, \mathcal{A}(T), \mathcal{B}(0), \dots, \mathcal{B}(T), \mathcal{C}(0), \dots, \mathcal{C}(T) \Longrightarrow \mathcal{C}(T+1)$$

Remark. Property A(t) establishes the stability of quadratic norms:

$$\min \left\{ \boldsymbol{b}_i(t)^{\top} \boldsymbol{b}_i(t), \, \boldsymbol{c}_i(t)^{\top} \boldsymbol{c}_i(t), \, \boldsymbol{b}(t)^{\top} \boldsymbol{b}(t) \right\} \ge d_h/2.$$

This norm lower bound induces two critical effects:

- 1. Convergence Rate: As we can see in property $\mathcal{B}(t)$, The upper bound of $\boldsymbol{c}_i^{\top}(t)\boldsymbol{b}_i(t)$, $\boldsymbol{c}_i^{\top}(t)\boldsymbol{b}_j(t)$ and $\boldsymbol{c}_i^{\top}(t)\boldsymbol{b}(t)$ is related to γ (lower bound of the squared norms), thus the stability of quadratic norms ensure a stable rapid convergence rate for property $\mathcal{B}(t)$.
- 2. Saddle Point Avoidance: The strict positivity (> 0) of $|\mathbf{b}_i|^2$ and $|\mathbf{c}_i|^2$ prevents the dynamics collapse to undesirable solutions $\mathbf{b}_i = \mathbf{c}_i = \mathbf{0}$, which would permanently make $\mathbf{c}_i^{\top} \mathbf{b}_i = 0$ (saddle points).

Property $\mathcal{B}(t)$ establishes a rapid exponential convergence rate:

$$oldsymbol{C}^ op oldsymbol{B} o rac{eta_3}{eta_1} oldsymbol{I}, \quad oldsymbol{C}^ op oldsymbol{b} o oldsymbol{0}$$

The rapid convergence rate ensures that the variations of *Squared norms* (in property $\mathcal{A}(t)$) and *Cross-interactions* (in property $\mathcal{C}(t)$) remain bounded, thereby establishing their constraints. For example, at initialization, $\boldsymbol{b}_i^{\top}(0)\boldsymbol{b}_i(0)$ is bounded by $3d_h/4 \leq \boldsymbol{b}_i^{\top}(0)\boldsymbol{b}_i(0) \leq 5d_h/4$. Further, thanks to the exponential convergence rate in property $\mathcal{B}(t)$, we can prove that $|\boldsymbol{b}_i^{\top}(t)\boldsymbol{b}_i(t) - \boldsymbol{b}_i^{\top}(0)\boldsymbol{b}_i(0)| \leq d_h/4$, and therefore $d_h/2 \leq \boldsymbol{b}_i^{\top}(t)\boldsymbol{b}_i(t), \boldsymbol{c}_i^{\top}(t)\boldsymbol{c}_i(t), \boldsymbol{b}^{\top}(t)\boldsymbol{b}(t) \leq 3d_h/2 \leq 2d_h$.

Property $\mathcal{C}(t)$ establishes the upper bound for the *Cross-interactions*. As we discuss in section 5.2, if the *Squared norms* ($\boldsymbol{c}_i^{\top}(t)\boldsymbol{b}_i(t)$, $\boldsymbol{c}_i^{\top}(t)\boldsymbol{b}_j(t)$ and $\boldsymbol{c}_i^{\top}(t)\boldsymbol{b}(t)$) are larger enough than the *Cross-interactions* ($\boldsymbol{b}_i^{\top}(t)\boldsymbol{b}_j(t)$, $\boldsymbol{c}_i^{\top}(t)\boldsymbol{c}_j(t)$, and $\boldsymbol{b}_i^{\top}(t)\boldsymbol{b}(t)$), we can make use of the *negative feedback term* to establish an exponential convergence rate. Thus property $\mathcal{C}(t)$ is also important.

The proof of claim B.1, claim B.2, and claim B.3 are in section C.4, section C.5, and section C.6 respectively.

Proof of Theorem 4.1 After convergence $(t \to 0)$, we will have $C^{\top}B = \frac{\beta_3}{\beta_1}I$, $C^{\top}b = 0$ (by property $\mathcal{B}(t)$), and $b_B(t) = b_C(t) = 0$ (by lemma A.5).

We will restate some equality for ease of reference.

Linear Recurrence (restatement of (Eq. (18)))

$$\boldsymbol{h}_{l}^{(d+1)} = \alpha \boldsymbol{h}_{l-1}^{(d+1)} + (1 - \alpha) y_{l} \boldsymbol{B}_{l}$$
 (24)

Prediction Output (restatement of (Eq. (22)))

$$\hat{y}_q = (1 - \alpha)(Cx_q + b_C)^{\top} \sum_{i=0}^{N-1} \alpha^{i+1} y_{N-i} (Bx_{N-i} + y_{N-i}b + b_B)$$
 (25)

Loss (restatement of (Eq. (23)))

$$\mathcal{L}(\boldsymbol{\theta}) = \frac{1}{2} \mathbb{E} \left[\left((1 - \alpha)(\boldsymbol{C}\boldsymbol{x}_q + \boldsymbol{b}_C)^\top \sum_{i=0}^{N-1} \alpha^{i+1} y_{N-i} (\boldsymbol{B}\boldsymbol{x}_{N-i} + y_{N-i}\boldsymbol{b} + \boldsymbol{b}_B) - \boldsymbol{w}^\top \boldsymbol{x}_q \right)^2 \right]$$
(26)

Based on (Eq. (24)), we have

$$(\boldsymbol{W}_{C}^{\top})_{[1:d,:]}(t)\boldsymbol{h}_{l}^{(d+1)} = \alpha(\boldsymbol{W}_{C}^{\top})_{[1:d,:]}(t)\boldsymbol{h}_{l-1}^{(d+1)} + (1-\alpha)y_{l}(\boldsymbol{W}_{C}^{\top})_{[1:d,:]}(t)\boldsymbol{B}_{l}$$

$$= \alpha(\boldsymbol{W}_{C}^{\top})_{[1:d,:]}(t)\boldsymbol{h}_{l-1}^{(d+1)} + (1-\alpha)y_{l}\boldsymbol{C}^{\top}(t)(\boldsymbol{B}(t)\boldsymbol{x}_{l} + y_{l}\boldsymbol{b}(t) + \boldsymbol{b}_{B}(t))$$

$$= \alpha(\boldsymbol{W}_{C}^{\top})_{[1:d,:]}(t)\boldsymbol{h}_{l-1}^{(d+1)} + (1-\alpha)y_{l}\boldsymbol{C}^{\top}(t)\boldsymbol{B}(t)\boldsymbol{x}_{l} + (1-\alpha)y_{l}^{2}\boldsymbol{C}^{\top}(t)\boldsymbol{b}(t)$$

$$= \alpha(\boldsymbol{W}_{C}^{\top})_{[1:d,:]}(t)\boldsymbol{h}_{l-1}^{(d+1)} + (1-\alpha)\frac{\beta_{3}}{\beta_{1}}y_{l}\boldsymbol{x}_{l}$$

$$= \alpha(\boldsymbol{W}_{C}^{\top})_{[1:d,:]}(t)\boldsymbol{h}_{l-1}^{(d+1)} + \frac{2(1+\alpha)(1-\alpha)}{\alpha(3(1-\alpha)d+4-2\alpha)}y_{l}\boldsymbol{x}_{l}$$

$$(27)$$

where the second equality is by selection rule $B_l = W_B e_l + b_B$ (Eq. (3)), and $W_B = [B b]$, $e_l = (x_l^\top, y_l)^\top$. The third equality is by $b_B(t) = 0$. The fourth equality is by $C^\top B = \frac{\beta_3}{\beta_1} I$ and $C^\top b = 0$. (Eq. (27)) establish the first equation (Thm 4.1 (a)) of the Theorem.

Based on (Eq. (25)), we have

$$\hat{y}_{q} = (1 - \alpha)(\mathbf{C}\mathbf{x}_{q} + \mathbf{b}_{C})^{\top} \sum_{i=0}^{N-1} \alpha^{i+1} y_{N-i} (\mathbf{B}\mathbf{x}_{N-i} + y_{N-i}\mathbf{b} + \mathbf{b}_{B})$$

$$= \mathbf{x}_{q}^{\top} \mathbf{C}^{\top} \sum_{i=0}^{N-1} (1 - \alpha) \alpha^{i+1} y_{N-i} (\mathbf{B}\mathbf{x}_{N-i} + y_{N-i}\mathbf{b})$$

$$= \mathbf{x}_{q}^{\top} \sum_{i=0}^{N-1} (1 - \alpha) \alpha^{i+1} y_{N-i} \mathbf{C}^{\top} \mathbf{B}\mathbf{x}_{N-i} + \mathbf{x}_{q}^{\top} \sum_{i=0}^{N-1} (1 - \alpha) \alpha^{i+1} y_{N-i}^{2} \mathbf{C}^{\top} \mathbf{b}$$

$$= \mathbf{x}_{q}^{\top} \sum_{i=0}^{N-1} (1 - \alpha) \alpha^{i+1} \frac{\beta_{3}}{\beta_{1}} y_{N-i} \mathbf{x}_{N-i}$$

$$= \mathbf{x}_{q}^{\top} \sum_{i=0}^{N-1} \frac{2\alpha^{i} (1 + \alpha) (1 - \alpha)}{(3(1 - \alpha)d + 4 - 2\alpha)} y_{N-i} \mathbf{x}_{N-i}$$
(28)

where the second equality is by $b_B(t) = b_C(t) = 0$. The fourth equality is by $C^{\top}B = \frac{\beta_3}{\beta_1}I$ and $C^{\top}b = 0$. (Eq. (28)) establish the second equation (Thm 4.1 (b)) of the Theorem.

Based on (Eq. (25)), we have

$$\mathcal{L}(\boldsymbol{\theta}) = \frac{1}{2} \mathbb{E} \left[\left((1 - \alpha)(\boldsymbol{C} \boldsymbol{x}_{q} + \boldsymbol{b}_{C})^{\top} \sum_{i=0}^{N-1} \alpha^{i+1} y_{N-i} (\boldsymbol{B} \boldsymbol{x}_{N-i} + y_{N-i} \boldsymbol{b} + \boldsymbol{b}_{B}) - \boldsymbol{w}^{\top} \boldsymbol{x}_{q} \right)^{2} \right] \\
= \frac{1}{2} \mathbb{E} \left[\left(\frac{\beta_{3}}{\beta_{1}} \boldsymbol{x}_{q}^{\top} \sum_{i=0}^{N-1} (1 - \alpha) \alpha^{i+1} y_{N-i} \boldsymbol{x}_{N-i} - \boldsymbol{w}^{\top} \boldsymbol{x}_{q} \right)^{2} \right] \\
= \frac{1}{2} \mathbb{E} \left[\left(\frac{\beta_{3}}{\beta_{1}} \boldsymbol{x}_{q}^{\top} \sum_{i=0}^{N-1} (1 - \alpha) \alpha^{i+1} y_{N-i} \boldsymbol{x}_{N-i} \right)^{2} \right] \\
- \mathbb{E} \left[\frac{\beta_{3}}{\beta_{1}} \boldsymbol{x}_{q}^{\top} \sum_{i=0}^{N-1} (1 - \alpha) \alpha^{i+1} y_{N-i} \boldsymbol{x}_{N-i} \cdot \boldsymbol{w}^{\top} \boldsymbol{x}_{q} \right] \\
+ \frac{1}{2} \mathbb{E} \left[\left(\boldsymbol{w}^{\top} \boldsymbol{x}_{q} \right)^{2} \right] \\
= d \text{ (by lemma A.2)}$$

We compute terms \spadesuit and \clubsuit as follows:

$$\begin{split} & \spadesuit = \mathbb{E}\Big[\Big(\frac{\beta_3}{\beta_1} \boldsymbol{x}_q^\top \sum_{i=0}^{N-1} (1-\alpha)\alpha^{i+1} y_{N-i} \boldsymbol{x}_{N-i}\Big)^2\Big] \\ & = \frac{\beta_3^2}{\beta_1^2} \mathbb{E}\Big[\boldsymbol{x}_q^\top \big(\sum_{i=0}^{N-1} (1-\alpha)\alpha^{i+1} y_{N-i} \boldsymbol{x}_{N-i}\big) \big(\sum_{i=0}^{N-1} (1-\alpha)\alpha^{i+1} y_{N-i} \boldsymbol{x}_{N-i}\big)^\top \boldsymbol{x}_q\Big] \\ & = \frac{(1-\alpha)^2 \beta_3^2}{\beta_1^2} \mathbb{E}_{\boldsymbol{x}_q} \Big[\boldsymbol{x}_q^\top \mathbb{E}_{\boldsymbol{x}_{N-i}, \boldsymbol{x}_{N-j}, \boldsymbol{w}} \Big[\sum_{i=0}^{N-1} \sum_{j=0}^{N-1} \alpha^{i+j+2} y_{N-i} y_{N-j} \boldsymbol{x}_{N-i} \boldsymbol{x}_{N-j}^\top \Big] \boldsymbol{x}_q\Big] \\ & = \frac{(1-\alpha)^2 \beta_3^2}{\beta_1^2} \cdot \Big(\frac{\alpha^2 \big(1-\alpha^N\big)^2}{(1-\alpha)^2} + \frac{(d+1)\alpha^2 \big(1-\alpha^{2N}\big)}{(1-\alpha)(1+\alpha)}\Big) \mathbb{E}\Big[\boldsymbol{x}_q^\top \boldsymbol{I} \boldsymbol{x}_q\Big] \\ & = \frac{d\beta_3^2}{\beta_2^2} \end{split}$$

For the fourth equality, $\mathbb{E}_{\boldsymbol{x}_{N-i},\boldsymbol{x}_{N-j},\boldsymbol{w}} \left[\sum_{i=0}^{N-1} \sum_{j=0}^{N-1} \alpha^{i+j+2} y_{N-i} y_{N-j} \boldsymbol{x}_{N-i} \boldsymbol{x}_{N-j}^{\top} \right] = \left(\frac{\alpha^2 \left(1 - \alpha^N \right)^2}{(1 - \alpha)^2} + \frac{(d+1)\alpha^2 \left(1 - \alpha^{2N} \right)}{(1 - \alpha)(1 + \alpha)} \right) \cdot \boldsymbol{I} \text{ by lemma A.3.} \quad \text{The last equality is by } \beta_1 = \left(\alpha^2 \left(1 - \alpha^N \right)^2 + \frac{(d+1)\alpha^2 (1 - \alpha) \left(1 - \alpha^{2N} \right)}{(1 + \alpha)} \right)$

$$\begin{split} & \clubsuit = \mathbb{E}\Big[\frac{\beta_3}{\beta_1} \boldsymbol{x}_q^\top \sum_{i=0}^{N-1} (1-\alpha)\alpha^{i+1} y_{N-i} \boldsymbol{x}_{N-i} \cdot \boldsymbol{w}^\top \boldsymbol{x}_q \Big] \\ & = \frac{(1-\alpha)\beta_3}{\beta_1} \mathbb{E}_{\boldsymbol{x}_q} \Big[\boldsymbol{x}_q^\top \mathbb{E}_{\boldsymbol{x}_{N-i},\boldsymbol{w}} \Big[\sum_{i=0}^{N-1} \alpha^{i+1} y_{N-i} \boldsymbol{x}_{N-i} \boldsymbol{w}^\top \Big] \boldsymbol{x}_q \Big] \\ & = \frac{(1-\alpha)\beta_3}{\beta_1} \cdot \alpha \Big(\frac{1-\alpha^N}{1-\alpha} \Big) \mathbb{E} \Big[\boldsymbol{x}_q^\top \boldsymbol{I} \boldsymbol{x}_q \Big] \\ & = \frac{d\beta_3^2}{\beta_1} \end{split}$$

For the third equality, $\mathbb{E}_{\boldsymbol{x}_{N-i},\boldsymbol{w}}\Big[\sum_{i=0}^{N-1}\alpha^{i+1}y_{N-i}\boldsymbol{x}_{N-i}\boldsymbol{w}^{\top}\Big] = \alpha\Big(\frac{1-\alpha^N}{1-\alpha}\Big)\boldsymbol{I}$ by lemma A.3. The last equality is by $\beta_3 = \alpha\Big(1-\alpha^N\Big)$.

Substituting ♠ and ♣ into (Eq. (29)) and get:

$$\mathcal{L}(\theta) = \frac{d\beta_3^3}{2\beta_1} - \frac{d\beta_3^3}{\beta_1} + \frac{1}{2}d$$

$$= \frac{d}{2} \left(1 - \frac{\beta_3^2}{\beta_1} \right)$$

$$= \frac{d(d+1)\alpha^2 (1-\alpha) \left(1 - \alpha^{2N} \right)}{2(1+\alpha)\beta_1}$$

$$= d(d+1)(1-\alpha) \cdot \frac{\alpha^2}{\beta_1} \cdot \frac{\left(1 - \alpha^{2N} \right)}{2(1+\alpha)}$$

$$\leq \frac{d(d+1)}{N} \cdot 4 \cdot \frac{3}{8}$$

$$= \frac{3d(d+1)}{2N}$$
(30)

Recall $\beta_1 = \left(\alpha^2 \left(1 - \alpha^N\right)^2 + \frac{(d+1)\alpha^2 (1-\alpha) \left(1 - \alpha^{2N}\right)}{(1+\alpha)}\right)$, $\beta_3 = \alpha \left(1 - \alpha^N\right)$ and $\alpha = \exp((-\ln 2)/N)$. For the inequality, $1 - \alpha = 1 - \exp((-\ln 2)/N) \le \frac{\ln 2}{N} \le \frac{1}{N}$, $\beta_1 \ge \alpha^2 \left(1 - \alpha^N\right)^2 = \frac{1}{4}\alpha^2$, $1 - \alpha^{2N} = 1 - \frac{1}{4} = \frac{3}{4}$, thus $d(d+1)(1-\alpha) \le \frac{d(d+1)}{N}$, $\frac{\alpha^2}{\beta_1} \le 4$, $\frac{\left(1 - \alpha^{2N}\right)}{2(1+\alpha)} \le \frac{3}{8(1+\alpha)} \le \frac{3}{8}$. (Eq. (30)) establish the third equation (Thm 4.1 (c)) of the Theorem.

C Complete Proof

This section presents the complete proof for the above results. To begin with, we provide the exact assumptions for N, η and d_h as part of Assumption 4.1.

Assumption

$$N = \Omega(d) \ge \max\{\frac{2\ln 2}{\ln 6 - \ln 5}, \frac{3(d+1)\ln 2}{2}\}$$
$$\eta = O(d^{-2}d_h^{-1}) \le \frac{1}{2d^2d_h} \le \frac{\ln 2}{\beta_2 d_h}$$
$$d_h = \widetilde{\Omega}(d^2) \ge \max\{\lambda_1, \dots, \lambda_{11}\}$$

where

$$\begin{split} \lambda_1 &= \left(1728 \log(4d(2d+1)/\delta) + 576(d-1)\beta_1 \log(4d(2d+1)/\delta)\right)/\beta_1 \\ \lambda_2 &= \left(576 \log(4d(2d+1)/\delta) + 192 \log(4d(2d+1)/\delta)\right)/\beta_1 \\ \lambda_3 &= \left(1728 \log(4d(2d+1)/\delta) + (576d+1872)\beta_1 \log(4d(2d+1)/\delta)\right)/\beta_1 \\ \lambda_4 &= 576 \log(4d(2d+1)/\delta)/\beta_1 + 192(d-1) \log(4d(2d+1)/\delta) \\ &+ 384 \log(4d(2d+1)/\delta)/\beta_1 + 3840 \ln 2 \log(4d(2d+1)/\delta) \\ \lambda_5 &= 2448d \log(4d(2d+1)/\delta) \\ \lambda_6 &= 816d \log(4d(2d+1)/\delta) + 768 \ln 2d^2 \log(4d(2d+1)/\delta) + 48 \log(4d(2d+1)/\delta)/\beta_1 \\ \lambda_7 &= \left(\frac{1}{\sqrt{\log(4d(2d+1)/\delta)}} + 24\sqrt{\log(4d(2d+1)/\delta)} \left(8\beta_1(d-1) + 10 + 6\beta_1 + 12\eta\beta_1/d\right)\right)^2 \\ \lambda_8 &= 36 \log(4d(2d+1)/\delta) \left(\frac{8}{\beta_1} + 8(d-2) + 6 + \frac{12}{d}\right)^2 \\ \lambda_9 &= 36 \log(4d(2d+1)/\delta) \left(4(d-1) + 56d \ln 2\right)^2 \\ \lambda_{10} &= 36 \log(4d(2d+1)/\delta) \left(6 + 4\beta_1(d-1) + 2(d-1)\right)^2 \end{split}$$

$$\lambda_{11} = \frac{36}{(\ln(3/2))^2} \log(4d(2d+1)/\delta) \left(\frac{32d}{\beta_1} + \frac{8\beta_3}{\beta_1} + 32d + \frac{8}{\beta_1} + \frac{4}{\beta_1\beta_2} + 80\ln 2\right)^2$$

Note that under the assumption of $N \geq \max\{\frac{2\ln 2}{\ln 6 - \ln 5}, \frac{3(d+1)\ln 2}{2}\}$ and combining $\alpha = \exp((-\ln 2)/N)$, we have the following:

$$\frac{5}{6} \le \alpha^2$$
, $4\beta_1 \le \beta_2$, $\frac{5}{24} \le \beta_1 \le \frac{3}{4}$, $1 \le \frac{1}{2}d \le \beta_2$, $\beta_3 = \Theta(1)$

These condition will be use to prove some bounds.

C.1 Proof of Lemma A.2

Lemma C.1 (restatement of lemma A.2) If vectors x and w are iid generated from $\mathcal{N}(0, I_d)$, $y = x^\top w$ we have the following expectations:

$$\mathbb{E}\left[\boldsymbol{x}\boldsymbol{x}^{\top}\boldsymbol{w}\boldsymbol{w}^{\top}\boldsymbol{x}\boldsymbol{x}^{\top}\right] = (d+2)\boldsymbol{I},$$
$$\mathbb{E}\left[y^{2}\right] = d,$$
$$\mathbb{E}\left[y^{4}\right] = 3d(d+2).$$

Proof. For (i, j)-th element of $\mathbb{E}\left[\boldsymbol{x}\boldsymbol{x}^{\top}\boldsymbol{w}\boldsymbol{w}^{\top}\boldsymbol{x}\boldsymbol{x}^{\top}\right]$, we have:

$$egin{aligned} \mathbb{E} \Big[oldsymbol{x} oldsymbol{x}^ op oldsymbol{w}^ op oldsymbol{x} oldsymbol{x}^ op \Big]_{[i,j]} &= \mathbb{E} \Big[oldsymbol{x}_{[i]} \sum_{k=1}^d oldsymbol{x}_{[k]} oldsymbol{w}_{[k]} \Big) \sum_{l=1}^d oldsymbol{w}_{[l]} oldsymbol{x}_{[l]} oldsymbol{x}_{[j]} \Big] \ &= \sum_{k=1}^d \sum_{l=1}^d \mathbb{E} \Big[oldsymbol{x}_{[i]} oldsymbol{x}_{[l]} oldsymbol{x}_{[l]} oldsymbol{x}_{[l]} \Big] \mathbb{E} \Big[oldsymbol{w}_{[k]} oldsymbol{w}_{[l]} \Big] \end{aligned}$$

According to the distribution of w, we have $\mathbb{E}\left[w_{[k]}w_{[l]}\right] = \delta_{kl}$, where δ_{kl} is the Kronecker delta defined as:

$$\delta_{kl} = \begin{cases} 1 & \text{if } k = l, \\ 0 & \text{if } k \neq l, \end{cases}$$

By Isserlis Theorem, we have:

$$\begin{split} & \mathbb{E} \Big[\boldsymbol{x}_{[i]} \boldsymbol{x}_{[k]} \boldsymbol{x}_{[l]} \boldsymbol{x}_{[j]} \Big] \\ & = \mathbb{E} \Big[\boldsymbol{x}_{[i]} \boldsymbol{x}_{[k]} \Big] \mathbb{E} \Big[\boldsymbol{x}_{[l]} \boldsymbol{x}_{[j]} \Big] + \mathbb{E} \Big[\boldsymbol{x}_{[i]} \boldsymbol{x}_{[l]} \Big] \mathbb{E} \Big[\boldsymbol{x}_{[k]} \boldsymbol{x}_{[j]} \Big] + \mathbb{E} \Big[\boldsymbol{x}_{[i]} \boldsymbol{x}_{[j]} \Big] \mathbb{E} \Big[\boldsymbol{x}_{[k]} \boldsymbol{x}_{[l]} \Big] \\ & = \delta_{ik} \delta_{lj} + \delta_{il} \delta_{kj} + \delta_{ij} \delta_{kl} \end{split}$$

Then we have:

$$\mathbb{E}\left[\boldsymbol{x}\boldsymbol{x}^{\top}\boldsymbol{w}\boldsymbol{w}^{\top}\boldsymbol{x}\boldsymbol{x}^{\top}\right]_{[i,j]}$$

$$= \sum_{k=1}^{d} \sum_{l=1}^{d} \mathbb{E}\left[\boldsymbol{x}_{[i]}\boldsymbol{x}_{[k]}\boldsymbol{x}_{[l]}\boldsymbol{x}_{[j]}\right] \mathbb{E}\left[\boldsymbol{w}_{[k]}\boldsymbol{w}_{[l]}\right]$$

$$= \sum_{k=1}^{d} \sum_{l=1}^{d} \left(\delta_{ik}\delta_{lj} + \delta_{il}\delta_{kj} + \delta_{ij}\delta_{kl}\right)\delta_{kl}$$

$$= \sum_{k=1}^{d} (2\delta_{ik}\delta_{kj} + \delta_{ij}\delta_{kk})$$

$$= (d+2)\delta_{ij}$$

Then we have:

$$egin{aligned} \mathbb{E}ig[y^2ig] &= \mathbb{E}ig[oldsymbol{x}^ op oldsymbol{w} \cdot oldsymbol{x}^ op oldsymbol{w} \ &= \mathbb{E}ig[\sum_{i=1}^d ig(oldsymbol{x}_{[i]} oldsymbol{w}_{[i]}ig) \sum_{j=1}^d ig(oldsymbol{x}_{[j]} oldsymbol{w}_{[j]}ig) \ &= \sum_{i=1}^d \sum_{j=1}^d eta^2_{ij} \ &= \sum_{i=1}^d \sum_{j=1}^d \delta^2_{ij} \end{aligned}$$

 $\mathbb{E} \left[\boldsymbol{x} \boldsymbol{x}^{\top} \boldsymbol{w} \boldsymbol{w}^{\top} \boldsymbol{x} \boldsymbol{x}^{\top} \right] = (d+2) \boldsymbol{I}$

$$\mathbb{E}\left[y^{4}\right] = \mathbb{E}\left[\boldsymbol{x}^{\top}\boldsymbol{w}\cdot\boldsymbol{x}^{\top}\boldsymbol{w}\cdot\boldsymbol{x}^{\top}\boldsymbol{w}\cdot\boldsymbol{x}^{\top}\boldsymbol{w}\cdot\boldsymbol{x}^{\top}\boldsymbol{w}\right]$$

$$= \mathbb{E}\left[\sum_{i=1}^{d}\left(\boldsymbol{x}_{[i]}\boldsymbol{w}_{[i]}\right)\sum_{j=1}^{d}\left(\boldsymbol{x}_{[j]}\boldsymbol{w}_{[j]}\right)\sum_{k=1}^{d}\left(\boldsymbol{x}_{[k]}\boldsymbol{w}_{[k]}\right)\sum_{l=1}^{d}\left(\boldsymbol{x}_{[l]}\boldsymbol{w}_{[l]}\right)\right]$$

$$= \sum_{i=1}^{d}\sum_{j=1}^{d}\sum_{k=1}^{d}\sum_{l=1}^{d}\mathbb{E}\left[\boldsymbol{x}_{[i]}\boldsymbol{x}_{[j]}\boldsymbol{x}_{[k]}\boldsymbol{x}_{[l]}\right]\mathbb{E}\left[\boldsymbol{w}_{[i]}\boldsymbol{w}_{[j]}\boldsymbol{w}_{[k]}\boldsymbol{w}_{[l]}\right]$$

$$= \sum_{i=1}^{d}\sum_{j=1}^{d}\sum_{k=1}^{d}\sum_{l=1}^{d}\left(\delta_{ik}\delta_{lj} + \delta_{il}\delta_{kj} + \delta_{ij}\delta_{kl}\right)^{2}$$

$$= \left(\sum_{i=j=k=l}+\sum_{j=1}+\sum_{i=j\neq k=l}+\sum_{i=k\neq j=l}+\sum_{i=k\neq j=l}\right)\cdot\left(\delta_{ik}\delta_{lj} + \delta_{il}\delta_{kj} + \delta_{ij}\delta_{kl}\right)^{2}$$

$$= d\cdot 3^{2} + 3\cdot d(d-1)\cdot 1^{2}$$

$$= 3d(d+2)$$

C.2 Proof of Lemma A.3

Lemma C.2 (restatement of lemma A.3) If vectors \mathbf{x}_i and \mathbf{w} are iid generated from $\mathcal{N}(0, \mathbf{I}_d)$, $y = \mathbf{x}_i^{\top} \mathbf{w}$ we have the following expectations:

$$\mathbb{E}\Big[\sum_{i=0}^{N-1}\sum_{j=0}^{N-1}\alpha^{i+j+2}y_{N-i}y_{N-j}\boldsymbol{x}_{N-i}\boldsymbol{x}_{N-j}^{\top}\Big] = \Big(\frac{\alpha^{2}(1-\alpha^{N})^{2}}{(1-\alpha)^{2}} + \frac{(d+1)\alpha^{2}(1-\alpha^{2N})}{(1-\alpha)(1+\alpha)}\Big) \cdot \boldsymbol{I},$$

$$\mathbb{E}\Big[\sum_{i=0}^{N-1}\sum_{j=0}^{N-1}\alpha^{i+j+2}y_{N-i}y_{N-j}^{2}\boldsymbol{x}_{N-i}\Big] = \boldsymbol{0},$$

$$\mathbb{E}\Big[\sum_{i=0}^{N-1}\sum_{j=0}^{N-1}\alpha^{i+j+2}y_{N-i}^{2}y_{N-j}^{2}\Big] = \frac{d^{2}\alpha^{2}(1-\alpha^{N})^{2}}{(1-\alpha)^{2}} + \frac{(2d^{2}+6d)\alpha^{2}(1-\alpha^{2N})}{(1-\alpha)(1+\alpha)},$$

$$\mathbb{E}\Big[\sum_{i=0}^{N-1}\alpha^{i+1}y_{N-i}\boldsymbol{x}_{N-i}\boldsymbol{w}^{\top}\Big] = \alpha\Big(\frac{1-\alpha^{N}}{1-\alpha}\Big) \cdot \boldsymbol{I},$$

$$\mathbb{E}\Big[\sum_{i=0}^{N-1} \alpha^{i+1} y_{N-i}^2 \boldsymbol{w}\Big] = \boldsymbol{0},$$

$$\mathbb{E}\Big[\sum_{i=0}^{N-1} \sum_{j=0}^{N-1} \alpha^{i+j+2} \boldsymbol{x}_{N-i} \boldsymbol{x}_{N-i}^{\top} \boldsymbol{w} \boldsymbol{x}_{N-j}^{\top} \boldsymbol{w}\Big] = \boldsymbol{0},$$

$$\mathbb{E}\Big[\sum_{i=0}^{N-1} \sum_{j=0}^{N-1} \alpha^{i+j+2} y_{N-i}^2 y_{N-j}\Big] = \boldsymbol{0},$$

$$\mathbb{E}\Big[\sum_{i=0}^{N-1} \sum_{j=0}^{N-1} \alpha^{i+j+2} y_{N-i} y_{N-j}\Big] = \frac{d\alpha^2 \left(1 - \alpha^{2N}\right)}{(1 - \alpha)(1 + \alpha)},$$

$$\mathbb{E}\Big[\sum_{i=0}^{N-1} \alpha^{i+1} y_{N-i} \boldsymbol{w}\Big] = \boldsymbol{0}.$$

Proof.

$$\mathbb{E}\left[\sum_{i=0}^{N-1}\sum_{j=0}^{N-1}\alpha^{i+j+2}y_{N-i}y_{N-j}\boldsymbol{x}_{N-i}\boldsymbol{x}_{N-j}^{\top}\right]$$

$$=\mathbb{E}\left[\sum_{i=0}^{N-1}\sum_{j=0}^{N-1}\alpha^{i+j+2}\boldsymbol{x}_{N-i}\underbrace{\boldsymbol{x}_{N-i}^{\top}\boldsymbol{w}}_{y_{N-i}}\boldsymbol{w}^{\top}\boldsymbol{x}_{N-j}\boldsymbol{x}_{N-j}^{\top}\right]$$

$$=\sum_{i\neq j}\alpha^{i+j+2}\underbrace{\mathbb{E}\left[\boldsymbol{x}_{N-i}\boldsymbol{x}_{N-i}^{\top}\right]}_{=I}\underbrace{\mathbb{E}\left[\boldsymbol{w}\boldsymbol{w}^{\top}\right]}_{y_{N-j}}\mathbb{E}\left[\boldsymbol{x}_{N-j}\boldsymbol{x}_{N-j}^{\top}\right]$$

$$+\sum_{i=j}\alpha^{i+j+2}\underbrace{\mathbb{E}\left[\boldsymbol{x}_{N-i}\boldsymbol{x}_{N-i}^{\top}\boldsymbol{w}\boldsymbol{w}^{\top}\boldsymbol{x}_{N-j}\boldsymbol{x}_{N-j}^{\top}\right]}_{=(d+2)I,\text{ by lemma }A.2}$$

$$=\left(\left(\sum_{i=0}^{N-1}\alpha^{i+1}\right)^{2}-\left(\sum_{i=0}^{N-1}\alpha^{2i+2}\right)\right)\boldsymbol{I}$$

$$+\left(\sum_{i=0}^{N-1}\alpha^{2i+2}\right)(d+2)\boldsymbol{I}$$

$$=\left(\frac{\alpha^{2}\left(1-\alpha^{N}\right)^{2}}{(1-\alpha)^{2}}+\frac{(d+1)\alpha^{2}\left(1-\alpha^{2N}\right)}{(1-\alpha)(1+\alpha)}\right)\cdot\boldsymbol{I}$$
Here, $\sum_{i\neq j}\alpha^{i+j+2}=\left(\sum_{i=0}^{N-1}\alpha^{i+1}\right)^{2}-\left(\sum_{i=0}^{N-1}\alpha^{2i+2}\right)\text{ for the third equality.}$

$$\mathbb{E}\left[\sum_{i=0}^{N-1}\sum_{j=0}^{N-1}\alpha^{i+j+2}y_{N-i}y_{N-j}^{2}\boldsymbol{x}_{N-i}\right]$$

$$=\mathbb{E}\left[\sum_{i=0}^{N-1}\sum_{j=0}^{N-1}\alpha^{i+j+2}x_{N-i}\underbrace{\boldsymbol{x}_{N-i}^{\top}\boldsymbol{w}\boldsymbol{x}_{N-j}^{\top}\boldsymbol{w}\boldsymbol{x}_{N-j}^{\top}\boldsymbol{w}\boldsymbol{x}_{N-j}^{\top}\boldsymbol{w}}_{y_{N-j}}\right]$$

$$=0$$

Notice that w appears three (odd) times in the second equality, if we define a function $g(w) = x_{N-i}x_{N-i}^{\top}wx_{N-j}^{\top}wx_{N-j}^{\top}w$, we can see that g(-w) = -g(w), and further $\mathbb{E}_{w}[g(w)] = 0$. Therefore, the above expectation equals to 0. We will use the similar property in some of the following equations.

$$\mathbb{E}\left[\sum_{i=0}^{N-1} \sum_{j=0}^{N-1} \alpha^{i+j+2} y_{N-i}^2 y_{N-j}^2\right]$$

$$= \sum_{i \neq j} \alpha^{i+j+2} \mathbb{E}\left[y_{N-i}^2\right] \mathbb{E}\left[y_{N-j}^2\right] + \sum_{i=j} \alpha^{i+j+2} \mathbb{E}\left[y_{N-i}^4\right]$$

$$= \left(\left(\sum_{i=0}^{N-1} \alpha^{i+1}\right)^2 - \left(\sum_{i=0}^{N-1} \alpha^{2i+2}\right)\right) d^2 + \left(\sum_{i=0}^{N-1} \alpha^{2i+2}\right) \cdot 3d(d+2)$$

$$= \frac{d^2 \alpha^2 \left(1 - \alpha^N\right)^2}{(1 - \alpha)^2} + \frac{(2d^2 + 6d)\alpha^2 \left(1 - \alpha^{2N}\right)}{(1 - \alpha)(1 + \alpha)}$$

where the second equality is by $\mathbb{E}\left[y_{N-i}^2\right] = \mathbb{E}\left[y_{N-j}^2\right] = d$ and $\mathbb{E}\left[y_{N-i}^4\right] = 3d(d+2)$ (lemma A.2).

$$\mathbb{E}\left[\sum_{i=0}^{N-1} \alpha^{i+1} y_{N-i} \boldsymbol{x}_{N-i} \boldsymbol{w}^{\top}\right]$$

$$= \sum_{i=0}^{N-1} \alpha^{i+1} \mathbb{E}\left[\boldsymbol{x}_{N-i} \boldsymbol{x}_{N-i}^{\top} \boldsymbol{w} \boldsymbol{w}^{\top}\right]$$

$$= \sum_{i=0}^{N-1} \alpha^{i+1} \mathbb{E}\left[\boldsymbol{x}_{N-i} \boldsymbol{x}_{N-i}^{\top}\right] \mathbb{E}\left[\boldsymbol{w} \boldsymbol{w}^{\top}\right]$$

$$= \sum_{i=0}^{N-1} \alpha^{i+1} \cdot \boldsymbol{I}$$

$$= \alpha \left(\frac{1-\alpha^{N}}{1-\alpha}\right) \cdot \boldsymbol{I}$$

$$\mathbb{E}\left[\sum_{i=0}^{N-1} \alpha^{i+1} y_{N-i}^{2} \boldsymbol{w}\right]$$

$$= \mathbb{E}\left[\sum_{i=0}^{N-1} \alpha^{i+1} \boldsymbol{w} \boldsymbol{x}_{N-i}^{\top} \boldsymbol{w} \boldsymbol{x}_{N-i}^{\top} \boldsymbol{w}\right]$$

Notice that w appears three (odd) times in the second equality, thus this expectation equals to 0.

$$\mathbb{E}\left[\sum_{i=0}^{N-1}\sum_{j=0}^{N-1}\alpha^{i+j+2}\boldsymbol{x}_{N-i}\boldsymbol{x}_{N-i}^{\top}\boldsymbol{w}\boldsymbol{x}_{N-j}^{\top}\boldsymbol{w}\right]$$

$$=\sum_{i=j}\alpha^{i+j+2}\mathbb{E}\left[\boldsymbol{x}_{N-i}\boldsymbol{x}_{N-i}^{\top}\boldsymbol{w}\boldsymbol{x}_{N-i}^{\top}\boldsymbol{w}\right]$$

$$+\sum_{i\neq j}\alpha^{i+j+2}\mathbb{E}\left[\boldsymbol{x}_{N-i}\boldsymbol{x}_{N-i}^{\top}\right]\mathbb{E}\left[\boldsymbol{w}\boldsymbol{x}_{N-j}^{\top}\boldsymbol{w}\right]$$

$$=\mathbf{0}+\mathbf{0}=\mathbf{0}$$

Notice that \boldsymbol{x}_{N-i} appears three (odd) times in $\mathbb{E}\Big[\boldsymbol{x}_{N-i}\boldsymbol{x}_{N-i}^{\top}\boldsymbol{w}\boldsymbol{x}_{N-i}^{\top}\boldsymbol{w}\Big]$, and $\boldsymbol{x}_{N-j}^{\top}$ appears once (odd) in $\mathbb{E}\Big[\boldsymbol{w}\boldsymbol{x}_{N-j}^{\top}\boldsymbol{w}\Big]$, thus this expectation equals to $\boldsymbol{0}$.

$$\mathbb{E}\left[\sum_{i=0}^{N-1}\sum_{j=0}^{N-1}\alpha^{i+j+2}y_{N-i}^{2}y_{N-j}\right]$$

$$=\sum_{i=0}^{N-1}\sum_{j=0}^{N-1}\alpha^{i+j+2}\mathbb{E}\left[\underbrace{\boldsymbol{x}_{N-i}^{\top}\boldsymbol{w}}_{y_{N-i}}\underbrace{\boldsymbol{x}_{N-i}^{\top}\boldsymbol{w}}_{y_{N-i}}\underline{\boldsymbol{x}_{N-j}^{\top}\boldsymbol{w}}\right]$$

$$=\mathbf{0}$$

Notice that w appears three (odd) times in the second equality, thus this expectation equals to 0.

$$\mathbb{E}\left[\sum_{i=0}^{N-1}\sum_{j=0}^{N-1}\alpha^{i+j+2}y_{N-i}y_{N-j}\right]$$

$$=\sum_{i\neq j}\alpha^{i+j+2}\mathbb{E}\left[y_{N-i}y_{N-j}\right] + \sum_{i=j}\alpha^{i+j+2}\mathbb{E}\left[y_{N-i}^{2}\right]$$

$$=\sum_{i\neq j}\alpha^{i+j+2}\mathbb{E}\left[\underbrace{\boldsymbol{x}_{N-i}^{\top}\boldsymbol{w}}_{y_{N-i}}\boldsymbol{x}_{N-j}^{\top}\boldsymbol{w}\right] + \sum_{i=0}^{N-1}\alpha^{2i+2}\mathbb{E}\left[y_{N-i}^{2}\right]$$

$$=\frac{d\alpha^{2}\left(1-\alpha^{2N}\right)}{(1-\alpha)(1+\alpha)}$$

Notice that $\boldsymbol{x}_{N-i}^{\top}$ appears once (odd) in $\mathbb{E}\left[\boldsymbol{x}_{N-i}^{\top}\boldsymbol{w}\boldsymbol{x}_{N-j}^{\top}\boldsymbol{w}\right]$ where $i\neq j$, thus $\sum_{i\neq j}\alpha^{i+j+2}\mathbb{E}\left[y_{N-i}y_{N-j}\right]=0$. Moreover, $\mathbb{E}\left[y_{N-i}^{2}\right]=d$ by lemma A.2.

$$\mathbb{E}\left[\sum_{i=0}^{N-1} \alpha^{i+1} y_{N-i} \boldsymbol{w}\right]$$

$$= \sum_{i=0}^{N-1} \alpha^{i+1} \mathbb{E}\left[\boldsymbol{w} \underbrace{\boldsymbol{w}^{\top} \boldsymbol{x}_{N-i}}_{y_{N-i}}\right]$$

$$= \sum_{i=0}^{N-1} \alpha^{i+1} \mathbb{E}\left[\boldsymbol{w} \boldsymbol{w}^{\top}\right] \mathbb{E}\left[\boldsymbol{x}_{N-i}\right]$$

$$= 0$$

C.3 Proof of Lemma A.4

Lemma C.3 (restatement of lemma A.4) The gradient of trainable parameters $\theta' = \{B, C, b, b_B, b_C\}$ with respect to loss (Eq. (23)) are as follows:

$$\nabla_{\boldsymbol{b}_{B}}\mathcal{L}(\boldsymbol{\theta}) = \boldsymbol{0},$$

$$\nabla_{\boldsymbol{b}_{C}}\mathcal{L}(\boldsymbol{\theta}) = \boldsymbol{0},$$

$$\nabla_{\boldsymbol{b}_{C}}\mathcal{L}(\boldsymbol{\theta}) = \boldsymbol{0},$$

$$\nabla_{\boldsymbol{B}}\mathcal{L}(\boldsymbol{\theta}) = \underbrace{\left(\alpha^{2}(1-\alpha^{N})^{2} + \frac{(d+1)\alpha^{2}(1-\alpha)\left(1-\alpha^{2N}\right)}{(1+\alpha)}\right)}_{:=\beta_{1}}\boldsymbol{C}\boldsymbol{C}^{\top}\boldsymbol{B} - \underbrace{\alpha\left(1-\alpha^{N}\right)}_{:=\beta_{3}}\boldsymbol{C},$$

$$\nabla_{\boldsymbol{b}}\mathcal{L}(\boldsymbol{\theta}) = \underbrace{\left(d^{2}\alpha^{2}(1-\alpha^{N})^{2} + \frac{(2d^{2}+6d)\alpha^{2}(1-\alpha)\left(1-\alpha^{2N}\right)}{(1+\alpha)}\right)}_{\boldsymbol{\theta}}\boldsymbol{C}\boldsymbol{C}^{\top}\boldsymbol{b},$$

$$\nabla_{\boldsymbol{C}} \mathcal{L}(\boldsymbol{\theta}) = \underbrace{\left(\alpha^{2} \left(1 - \alpha^{N}\right)^{2} + \frac{(d+1)\alpha^{2}(1-\alpha)\left(1 - \alpha^{2N}\right)}{(1+\alpha)}\right)}_{:=\beta_{1}} \boldsymbol{B} \boldsymbol{B}^{\top} \boldsymbol{C}$$

$$+ \underbrace{\left(d^{2}\alpha^{2} \left(1 - \alpha^{N}\right)^{2} + \frac{(2d^{2} + 6d)\alpha^{2}(1-\alpha)\left(1 - \alpha^{2N}\right)}{(1+\alpha)}\right)}_{:=\beta_{2}} \boldsymbol{b} \boldsymbol{b}^{\top} \boldsymbol{C}$$

$$- \underbrace{\alpha\left(1 - \alpha^{N}\right)}_{:=\beta_{2}} \boldsymbol{B}.$$

Proof of lemma C.3. Recalling the loss:

$$\mathcal{L}(\boldsymbol{\theta}) = \frac{1}{2} \mathbb{E} \Big[\Big((1 - \alpha) (\boldsymbol{C} \boldsymbol{x}_q + \boldsymbol{b}_C)^\top \sum_{i=0}^{N-1} \alpha^{i+1} y_{N-i} (\boldsymbol{B} \boldsymbol{x}_{N-i} + y_{N-i} \boldsymbol{b} + \boldsymbol{b}_B) - \boldsymbol{w}^\top \boldsymbol{x}_q \Big)^2 \Big]$$

We will compute the gradient of $\{B, C, b, b_B, b_C\}$ with respect to $\mathcal{L}(\theta)$. Some expectation calculation are detailed in Section C.3.1.

$$\nabla_{\boldsymbol{b}_{C}} \mathcal{L}(\boldsymbol{\theta}) = \mathbb{E}\Big[(1 - \alpha) \Big((1 - \alpha) (\boldsymbol{C} \boldsymbol{x}_{q} + \boldsymbol{b}_{C})^{\top} \underbrace{\sum_{i=0}^{N-1} \alpha^{i+1} y_{N-i} (\boldsymbol{B} \boldsymbol{x}_{N-i} + y_{N-i} \boldsymbol{b} + \boldsymbol{b}_{B})}_{:=\boldsymbol{v}} - \boldsymbol{w}^{\top} \boldsymbol{x}_{q} \Big)$$

$$\cdot \underbrace{\sum_{i=0}^{N-1} \alpha^{i+1} y_{N-i} (\boldsymbol{B} \boldsymbol{x}_{N-i} + y_{N-i} \boldsymbol{b} + \boldsymbol{b}_{B})}_{:=\boldsymbol{v}} \Big]$$

$$= (1 - \alpha)^{2} \mathbb{E} \Big[\boldsymbol{v} \boldsymbol{v}^{\top} (\boldsymbol{C} \boldsymbol{x}_{q} + \boldsymbol{b}_{C}) \Big] - (1 - \alpha) \mathbb{E} \Big[\boldsymbol{v} \boldsymbol{w}^{\top} \boldsymbol{x}_{q} \Big]$$

$$= (1 - \alpha)^{2} \mathbb{E} \Big[\boldsymbol{v} \boldsymbol{v}^{\top} \boldsymbol{C} \Big] \mathbb{E} \Big[\boldsymbol{x}_{q} \Big] + (1 - \alpha)^{2} \mathbb{E} \Big[\boldsymbol{v} \boldsymbol{v}^{\top} \Big] \boldsymbol{b}_{C} - (1 - \alpha) \mathbb{E} \Big[\boldsymbol{v} \boldsymbol{w}^{\top} \Big] \mathbb{E} \Big[\boldsymbol{x}_{q} \Big]$$

It is clear that $\mathbb{E}\left[x_q\right]=\mathbf{0}$. Thus, if $b_C=\mathbf{0}$, then $\nabla_{b_C}\mathcal{L}(\boldsymbol{\theta})=\mathbf{0}$. Notice that we assume $b_C(0)=\mathbf{0}$ at initialization, so by induction, $b_C(t)=\mathbf{0}$ and $\nabla_{b_C}\mathcal{L}(\boldsymbol{\theta}(t))=\mathbf{0}$ for $t\geq 0$. We will consider $b_C=\mathbf{0}$ when computing other gradients.

$$\nabla_{\boldsymbol{b}_{B}} \mathcal{L}(\boldsymbol{\theta}) = \mathbb{E}\Big[(1 - \alpha) \Big((1 - \alpha) (\boldsymbol{C} \boldsymbol{x}_{q} + \boldsymbol{b}_{C})^{\top} \sum_{i=0}^{N-1} \alpha^{i+1} y_{N-i} (\boldsymbol{B} \boldsymbol{x}_{N-i} + y_{N-i} \boldsymbol{b} + \boldsymbol{b}_{B}) - \boldsymbol{w}^{\top} \boldsymbol{x}_{q} \Big)$$

$$\cdot (\boldsymbol{C} \boldsymbol{x}_{q} + \boldsymbol{b}_{C}) \cdot \sum_{i=0}^{N-1} \alpha^{i+1} y_{N-i} \Big]$$

$$= \mathbb{E}\Big[(1 - \alpha) \Big((1 - \alpha) (\boldsymbol{C} \boldsymbol{x}_{q})^{\top} \sum_{i=0}^{N-1} \alpha^{i+1} y_{N-i} (\boldsymbol{B} \boldsymbol{x}_{N-i} + y_{N-i} \boldsymbol{b} + \boldsymbol{b}_{B}) - \boldsymbol{w}^{\top} \boldsymbol{x}_{q} \Big)$$

$$\cdot \boldsymbol{C} \boldsymbol{x}_{q} \cdot \sum_{i=0}^{N-1} \alpha^{i+1} y_{N-i} \Big]$$

$$= (1 - \alpha)^{2} \mathbb{E}\Big[\boldsymbol{x}_{q}^{\top} \boldsymbol{C}^{\top} \Big(\sum_{i=0}^{N-1} \alpha^{i+1} y_{N-i} (\boldsymbol{B} \boldsymbol{x}_{N-i} + y_{N-i} \boldsymbol{b} + \boldsymbol{b}_{B}) \Big) \cdot \boldsymbol{C} \boldsymbol{x}_{q} \cdot \sum_{i=0}^{N-1} \alpha^{i+1} y_{N-i} \Big]$$

$$- (1 - \alpha) \mathbb{E}\Big[\boldsymbol{w}^{\top} \boldsymbol{x}_{q} \cdot \boldsymbol{C} \boldsymbol{x}_{q} \cdot \sum_{i=0}^{N-1} \alpha^{i+1} y_{N-i} \Big]$$

$$= \frac{d\alpha^{2} (1 - \alpha) \Big(1 - \alpha^{2N} \Big)}{(1 + \alpha)} \boldsymbol{C} \boldsymbol{C}^{\top} \boldsymbol{b}_{B}$$

The last equality follows from lemma C.4 where we have: $\mathbb{E}\left[\boldsymbol{x}_q^{\top}\boldsymbol{C}^{\top}\left(\sum_{i=0}^{N-1}\alpha^{i+1}y_{N-i}(\boldsymbol{B}\boldsymbol{x}_{N-i}+y_{N-i})\boldsymbol{b}+\boldsymbol{b}_B)\right)\cdot\boldsymbol{C}\boldsymbol{x}_q\cdot\sum_{i=0}^{N-1}\alpha^{i+1}y_{N-i}\right]=\frac{d\alpha^2\left(1-\alpha^{2N}\right)}{(1-\alpha)(1+\alpha)}\boldsymbol{C}\boldsymbol{C}^{\top}\boldsymbol{b}_B, \text{ and } \mathbb{E}\left[\boldsymbol{w}^{\top}\boldsymbol{x}_q\cdot\boldsymbol{C}\boldsymbol{x}_q\cdot\sum_{i=0}^{N-1}\alpha^{i+1}y_{N-i}\right]=\boldsymbol{0}.$ Similar to \boldsymbol{b}_C , notice that \boldsymbol{b}_B is initialized as $\boldsymbol{0}$, thus by induction, $\boldsymbol{b}_B(t)=\boldsymbol{0}$ and $\nabla_{\boldsymbol{b}_B}\mathcal{L}(\boldsymbol{\theta}(t))=\boldsymbol{0}$ for $t\geq 0$. We will consider $\boldsymbol{b}_B=\boldsymbol{b}_C=\boldsymbol{0}$ when computing other gradients.

$$\nabla_{\boldsymbol{C}}\mathcal{L}(\boldsymbol{\theta}) = \mathbb{E}\Big[(1-\alpha) \Big((1-\alpha)(\boldsymbol{C}\boldsymbol{x}_{q} + \boldsymbol{b}_{C})^{\top} \sum_{i=0}^{N-1} \alpha^{i+1} y_{N-i} (\boldsymbol{B}\boldsymbol{x}_{N-i} + y_{N-i}\boldsymbol{b} + \boldsymbol{b}_{B}) - \boldsymbol{w}^{\top}\boldsymbol{x}_{q} \Big)$$

$$\cdot \sum_{i=0}^{N-1} \alpha^{i+1} y_{N-i} (\boldsymbol{B}\boldsymbol{x}_{N-i} + y_{N-i}\boldsymbol{b} + \boldsymbol{b}_{B}) \boldsymbol{x}_{q}^{\top} \Big]$$

$$= (1-\alpha)^{2} \mathbb{E}\Big[\boldsymbol{x}_{q}^{\top} \boldsymbol{C}^{\top} \Big(\sum_{i=0}^{N-1} \alpha^{i+1} y_{N-i} (\boldsymbol{B}\boldsymbol{x}_{N-i} + y_{N-i}\boldsymbol{b}) \Big)$$

$$\cdot \Big(\sum_{i=0}^{N-1} \alpha^{i+1} y_{N-i} (\boldsymbol{B}\boldsymbol{x}_{N-i} + y_{N-i}\boldsymbol{b}) \Big) \boldsymbol{x}_{q}^{\top} \Big]$$

$$- (1-\alpha) \mathbb{E}\Big[\boldsymbol{w}^{\top} \boldsymbol{x}_{q} \cdot \sum_{i=0}^{N-1} \alpha^{i+1} y_{N-i} (\boldsymbol{B}\boldsymbol{x}_{N-i} + y_{N-i}\boldsymbol{b}) \boldsymbol{x}_{q}^{\top} \Big]$$

$$= \underbrace{\Big(\alpha^{2} (1-\alpha^{N})^{2} + \frac{(d+1)\alpha^{2}(1-\alpha)(1-\alpha^{2N})}{(1+\alpha)} \Big)}_{:=\beta_{1}} \boldsymbol{B} \boldsymbol{B}^{\top} \boldsymbol{C}$$

$$+ \underbrace{\Big(d^{2}\alpha^{2} (1-\alpha^{N})^{2} + \frac{(2d^{2}+6d)\alpha^{2}(1-\alpha)(1-\alpha^{2N})}{(1+\alpha)} \Big)}_{:=\beta_{2}} \boldsymbol{b} \boldsymbol{b}^{\top} \boldsymbol{C}$$

$$- \underbrace{\alpha(1-\alpha^{N})}_{:=\beta_{2}} \boldsymbol{B}$$

The last equality follows from lemma C.4, where we have:

$$\begin{split} & \mathbb{E}\Big[\boldsymbol{x}_{q}^{\top}\boldsymbol{C}^{\top}\Big(\sum_{i=0}^{N-1}\alpha^{i+1}y_{N-i}(\boldsymbol{B}\boldsymbol{x}_{N-i}+y_{N-i}\boldsymbol{b})\Big) \cdot \Big(\sum_{i=0}^{N-1}\alpha^{i+1}y_{N-i}(\boldsymbol{B}\boldsymbol{x}_{N-i}+y_{N-i}\boldsymbol{b})\Big)\boldsymbol{x}_{q}^{\top}\Big] \\ & = \Big(\frac{\alpha^{2}\big(1-\alpha^{N}\big)^{2}}{(1-\alpha)^{2}} + \frac{(d+1)\alpha^{2}\big(1-\alpha^{2N}\big)}{(1-\alpha)(1+\alpha)}\Big)\boldsymbol{B}\boldsymbol{B}^{\top}\boldsymbol{C} \\ & + \Big(\frac{d^{2}\alpha^{2}\Big(1-\alpha^{N}\Big)^{2}}{(1-\alpha)^{2}} + \frac{(2d^{2}+6d)\alpha^{2}\Big(1-\alpha^{2N}\Big)}{(1-\alpha)(1+\alpha)}\Big)\boldsymbol{b}\boldsymbol{b}^{\top}\boldsymbol{C} \\ & \text{and } \mathbb{E}\Big[\boldsymbol{w}^{\top}\boldsymbol{x}_{q} \cdot \sum_{i=0}^{N-1}\alpha^{i+1}y_{N-i}(\boldsymbol{B}\boldsymbol{x}_{N-i}+y_{N-i}\boldsymbol{b})\boldsymbol{x}_{q}^{\top}\Big] = \alpha\Big(\frac{1-\alpha^{N}}{1-\alpha}\Big)\boldsymbol{B} \\ & \nabla_{\boldsymbol{B}}\mathcal{L}(\boldsymbol{\theta}) = \mathbb{E}\Big[(1-\alpha)\Big((1-\alpha)(\boldsymbol{C}\boldsymbol{x}_{q}+\boldsymbol{b}_{C})^{\top}\sum_{i=0}^{N-1}\alpha^{i+1}y_{N-i}(\boldsymbol{B}\boldsymbol{x}_{N-i}+y_{N-i}\boldsymbol{b}+\boldsymbol{b}_{B}) - \boldsymbol{w}^{\top}\boldsymbol{x}_{q}\Big) \\ & \cdot (\boldsymbol{C}\boldsymbol{x}_{q}+\boldsymbol{b}_{C})\sum_{i=0}^{N-1}\alpha^{i+1}y_{N-i}\boldsymbol{x}_{N-i}^{\top}\Big] \end{split}$$

$$= (1 - \alpha)^{2} \mathbb{E} \left[\boldsymbol{x}_{q}^{\top} \boldsymbol{C}^{\top} \sum_{i=0}^{N-1} \alpha^{i+1} y_{N-i} (\boldsymbol{B} \boldsymbol{x}_{N-i} + y_{N-i} \boldsymbol{b}) \cdot \boldsymbol{C} \boldsymbol{x}_{q} \sum_{i=0}^{N-1} \alpha^{i+1} y_{N-i} \boldsymbol{x}_{N-i}^{\top} \right]$$

$$- (1 - \alpha) \mathbb{E} \left[\boldsymbol{w}^{\top} \boldsymbol{x}_{q} \cdot \boldsymbol{C} \boldsymbol{x}_{q} \sum_{i=0}^{N-1} \alpha^{i+1} y_{N-i} \boldsymbol{x}_{N-i}^{\top} \right]$$

$$= \underbrace{\left(\alpha^{2} (1 - \alpha^{N})^{2} + \frac{(d+1)\alpha^{2} (1 - \alpha) (1 - \alpha^{2N})}{(1+\alpha)} \right)}_{:=\beta_{1}} \boldsymbol{C} \boldsymbol{C}^{\top} \boldsymbol{B}$$

$$- \underbrace{\alpha (1 - \alpha^{N})}_{:=\beta_{2}} \boldsymbol{C}$$

The last equality follows from lemma C.4, where we have:

$$\begin{split} \mathbb{E}\Big[\boldsymbol{x}_q^\top \boldsymbol{C}^\top \sum_{i=0}^{N-1} \alpha^{i+1} y_{N-i} (\boldsymbol{B} \boldsymbol{x}_{N-i} + y_{N-i} \boldsymbol{b}) \cdot \boldsymbol{C} \boldsymbol{x}_q \sum_{i=0}^{N-1} \alpha^{i+1} y_{N-i} \boldsymbol{x}_{N-i}^\top \Big] \\ &= \Big(\frac{\alpha^2 \big(1 - \alpha^N\big)^2}{(1 - \alpha)^2} + \frac{(d+1)\alpha^2 \big(1 - \alpha^{2N}\big)}{(1 - \alpha)(1 + \alpha)} \Big) \boldsymbol{C} \boldsymbol{C}^\top \boldsymbol{B} \\ \text{and } \mathbb{E}\Big[\boldsymbol{w}^\top \boldsymbol{x}_q \cdot \boldsymbol{C} \boldsymbol{x}_q \sum_{i=0}^{N-1} \alpha^{i+1} y_{N-i} \boldsymbol{x}_{N-i}^\top \Big] = \alpha \Big(\frac{1 - \alpha^N}{1 - \alpha}\Big) \boldsymbol{C}. \end{split}$$

$$\nabla_{\boldsymbol{b}} \mathcal{L}(\boldsymbol{\theta}) = \mathbb{E}\Big[(1 - \alpha) \Big((1 - \alpha) (\boldsymbol{C} \boldsymbol{x}_{q} + \boldsymbol{b}_{C})^{\top} \sum_{i=0}^{N-1} \alpha^{i+1} y_{N-i} (\boldsymbol{B} \boldsymbol{x}_{N-i} + y_{N-i} \boldsymbol{b} + \boldsymbol{b}_{B}) - \boldsymbol{w}^{\top} \boldsymbol{x}_{q} \Big)$$

$$\cdot (\boldsymbol{C} \boldsymbol{x}_{q} + \boldsymbol{b}_{C}) \cdot \sum_{i=0}^{N-1} \alpha^{i+1} y_{N-i}^{2} \Big]$$

$$= \mathbb{E}\Big[(1 - \alpha) \Big((1 - \alpha) (\boldsymbol{C} \boldsymbol{x}_{q})^{\top} \sum_{i=0}^{N-1} \alpha^{i+1} y_{N-i} (\boldsymbol{B} \boldsymbol{x}_{N-i} + y_{N-i} \boldsymbol{b}) - \boldsymbol{w}^{\top} \boldsymbol{x}_{q} \Big)$$

$$\cdot \boldsymbol{C} \boldsymbol{x}_{q} \cdot \sum_{i=0}^{N-1} \alpha^{i+1} y_{N-i}^{2} \Big]$$

$$= (1 - \alpha)^{2} \mathbb{E}\Big[\Big(\boldsymbol{x}_{q}^{\top} \boldsymbol{C}^{\top} \sum_{i=0}^{N-1} \alpha^{i+1} y_{N-i} (\boldsymbol{B} \boldsymbol{x}_{N-i} + y_{N-i} \boldsymbol{b}) \Big) \cdot \boldsymbol{C} \boldsymbol{x}_{q} \cdot \sum_{i=0}^{N-1} \alpha^{i+1} y_{N-i}^{2} \Big]$$

$$- (1 - \alpha) \mathbb{E}\Big[\boldsymbol{w}^{\top} \boldsymbol{x}_{q} \cdot \boldsymbol{C} \boldsymbol{x}_{q} \cdot \sum_{i=0}^{N-1} \alpha^{i+1} y_{N-i}^{2} \Big]$$

$$= \underbrace{\Big(d^{2} \alpha^{2} (1 - \alpha^{N})^{2} + \frac{(2d^{2} + 6d) \alpha^{2} (1 - \alpha) (1 - \alpha^{2N})}{(1 + \alpha)} \Big)}_{:=\beta_{2}} \boldsymbol{C} \boldsymbol{C}^{\top} \boldsymbol{b}$$

The last equality follows from lemma C.4, where we have:

$$\begin{split} \mathbb{E}\Big[\Big(\boldsymbol{x}_q^{\top}\boldsymbol{C}^{\top}\sum_{i=0}^{N-1}\alpha^{i+1}y_{N-i}(\boldsymbol{B}\boldsymbol{x}_{N-i}+y_{N-i}\boldsymbol{b})\Big)\cdot\boldsymbol{C}\boldsymbol{x}_q\cdot\sum_{i=0}^{N-1}\alpha^{i+1}y_{N-i}^2\Big]\\ &=\Big(\frac{d^2\alpha^2\Big(1-\alpha^N\Big)^2}{(1-\alpha)^2}+\frac{(2d^2+6d)\alpha^2\Big(1-\alpha^{2N}\Big)}{(1-\alpha)(1+\alpha)}\Big)\boldsymbol{C}\boldsymbol{C}^{\top}\boldsymbol{b}\\ \text{and } \mathbb{E}\Big[\boldsymbol{w}^{\top}\boldsymbol{x}_q\cdot\boldsymbol{C}\boldsymbol{x}_q\cdot\sum_{i=0}^{N-1}\alpha^{i+1}y_{N-i}^2\Big]=\mathbf{0}. \end{split}$$

C.3.1 Auxiliary Lemma for Lemma A.4

Lemma C.4 If vectors \mathbf{x}_i , \mathbf{x}_q and \mathbf{w} are iid generated from $\mathcal{N}(0, \mathbf{I}_d)$, $y = \mathbf{x}_i^{\top} \mathbf{w}$ we have the following expectations:

$$\begin{split} &\mathbb{E}\Big[\mathbf{x}_{q}^{\top}C^{\top}\Big(\sum_{i=0}^{N-1}\alpha^{i+1}y_{N-i}(B\mathbf{x}_{N-i}+y_{N-i}b)\Big) \cdot \Big(\sum_{i=0}^{N-1}\alpha^{i+1}y_{N-i}(B\mathbf{x}_{N-i}+y_{N-i}b)\Big)\mathbf{x}_{q}^{\top}\Big] \\ &= \Big(\frac{\alpha^{2}(1-\alpha^{N})^{2}}{(1-\alpha)^{2}} + \frac{(d+1)\alpha^{2}(1-\alpha^{2N})}{(1-\alpha)(1+\alpha)}\Big)\mathbf{B}\mathbf{B}^{\top}\mathbf{C} \\ &+ \Big(\frac{d^{2}\alpha^{2}\Big(1-\alpha^{N}\Big)^{2}}{(1-\alpha)^{2}} + \frac{(2d^{2}+6d)\alpha^{2}\Big(1-\alpha^{2N}\Big)}{(1-\alpha)(1+\alpha)}\Big)\mathbf{b}\mathbf{b}^{\top}\mathbf{C} \\ &\mathbb{E}\Big[\mathbf{w}^{\top}\mathbf{x}_{q} \cdot \sum_{i=0}^{N-1}\alpha^{i+1}y_{N-i}(B\mathbf{x}_{N-i}+y_{N-i}b)\mathbf{x}_{q}^{\top}\Big] = \alpha\Big(\frac{1-\alpha^{N}}{1-\alpha}\Big)\mathbf{B} \\ &\mathbb{E}\Big[\mathbf{x}_{q}^{\top}\mathbf{C}^{\top}\sum_{i=0}^{N-1}\alpha^{i+1}y_{N-i}(B\mathbf{x}_{N-i}+y_{N-i}b) \cdot \mathbf{C}\mathbf{x}_{q}\sum_{i=0}^{N-1}\alpha^{i+1}y_{N-i}\mathbf{x}_{N-i}^{\top}\Big] \\ &= \Big(\frac{\alpha^{2}(1-\alpha^{N})^{2}}{(1-\alpha)^{2}} + \frac{(d+1)\alpha^{2}(1-\alpha^{2N})}{(1-\alpha)(1+\alpha)}\Big)\mathbf{C}\mathbf{C}^{\top}\mathbf{B} \\ &\mathbb{E}\Big[\mathbf{w}^{\top}\mathbf{x}_{q} \cdot \mathbf{C}\mathbf{x}_{q}\sum_{i=0}^{N-1}\alpha^{i+1}y_{N-i}\mathbf{x}_{N-i}^{\top}\Big] = \alpha\Big(\frac{1-\alpha^{N}}{1-\alpha}\Big)\mathbf{C} \\ &\mathbb{E}\Big[\Big(\mathbf{x}_{q}^{\top}\mathbf{C}^{\top}\sum_{i=0}^{N-1}\alpha^{i+1}y_{N-i}(B\mathbf{x}_{N-i}+y_{N-i}b)\Big) \cdot \mathbf{C}\mathbf{x}_{q} \cdot \sum_{i=0}^{N-1}\alpha^{i+1}y_{N-i}^{2}\Big] \\ &= \Big(\frac{d^{2}\alpha^{2}(1-\alpha^{N})^{2}}{(1-\alpha)^{2}} + \frac{(2d^{2}+6d)\alpha^{2}(1-\alpha^{2N})}{(1-\alpha)(1+\alpha)}\Big)\mathbf{C}\mathbf{C}^{\top}\mathbf{b} \\ &\mathbb{E}\Big[\mathbf{w}^{\top}\mathbf{x}_{q} \cdot \mathbf{C}\mathbf{x}_{q} \cdot \sum_{i=0}^{N-1}\alpha^{i+1}y_{N-i}^{2}\Big] = \mathbf{0} \\ &\mathbb{E}\Big[\mathbf{x}_{q}^{\top}\mathbf{C}^{\top}\Big(\sum_{i=0}^{N-1}\alpha^{i+1}y_{N-i}(B\mathbf{x}_{N-i}+y_{N-i}b+b_{B})\Big) \cdot \mathbf{C}\mathbf{x}_{q} \cdot \sum_{i=0}^{N-1}\alpha^{i+1}y_{N-i}\Big] \\ &= \frac{d\alpha^{2}\Big(1-\alpha^{2N}\Big)}{(1-\alpha)(1+\alpha)}\mathbf{C}\mathbf{C}^{\top}\mathbf{b}_{B} \end{aligned}$$

$$\mathbb{E}\Big[\boldsymbol{w}^{\top}\boldsymbol{x}_{q}\cdot\boldsymbol{C}\boldsymbol{x}_{q}\cdot\sum_{i=0}^{N-1}\alpha^{i+1}y_{N-i}\Big]=\boldsymbol{0}$$

Proof of lemma C.4. We will use the results of lemma A.3 to prove the above equation.

$$\mathbb{E}\left[\underbrace{x_{q}^{\top}C^{\top}\left(\sum_{i=0}^{N-1}\alpha^{i+1}y_{N-i}(Bx_{N-i}+y_{N-i}b)\right)}_{\bullet} \cdot \left(\sum_{i=0}^{N-1}\alpha^{i+1}y_{N-i}(Bx_{N-i}+y_{N-i}b)\right)x_{q}^{\top}\right] \\ = \mathbb{E}\left[\sum_{i=0}^{N-1}\alpha^{i+1}y_{N-i}(Bx_{N-i}+y_{N-i}b)\right] \\ \left(\underbrace{x_{q}^{\top}C^{\top}\sum_{i=0}^{N-1}\alpha^{i+1}y_{N-i}(Bx_{N-i}+y_{N-i}b)}_{\bullet}\right]^{\top}x_{q}^{\top}\right] \\ = \mathbb{E}\left[\sum_{i=0}^{N-1}\sum_{j=0}^{N-1}\alpha^{i+j+2}\left(\left(By_{N-i}x_{N-i}+y_{N-i}^{2}b\right)\right)\right] \\ \left(y_{N-j}x_{N-j}^{\top}B^{\top}+y_{N-j}^{2}b^{\top}\right)\right]C\mathbb{E}\left[x_{q}x_{q}^{\top}\right] \\ = B\mathbb{E}\left[\sum_{i=0}^{N-1}\sum_{j=0}^{N-1}\alpha^{i+j+2}y_{N-i}y_{N-j}x_{N-i}x_{N-j}^{\top}\right]B^{\top}C \\ + B\mathbb{E}\left[\sum_{i=0}^{N-1}\sum_{j=0}^{N-1}\alpha^{i+j+2}y_{N-i}y_{N-j}^{2}x_{N-i}\right]b^{\top}C \\ + b\mathbb{E}\left[\sum_{i=0}^{N-1}\sum_{j=0}^{N-1}\alpha^{i+j+2}y_{N-j}y_{N-i}^{2}x_{N-j}\right]B^{\top}C \\ + \mathbb{E}\left[\sum_{i=0}^{N-1}\sum_{j=0}^{N-1}\alpha^{i+j+2}y_{N-i}^{2}y_{N-j}^{2}\right]bb^{\top}C \\ = \left(\frac{\alpha^{2}(1-\alpha^{N})^{2}}{(1-\alpha)^{2}} + \frac{(d+1)\alpha^{2}(1-\alpha^{2N})}{(1-\alpha)(1+\alpha)}\right)BB^{\top}C \\ + \left(\frac{d^{2}\alpha^{2}\left(1-\alpha^{N}\right)^{2}}{(1-\alpha)^{2}} + \frac{(2d^{2}+6d)\alpha^{2}\left(1-\alpha^{2N}\right)}{(1-\alpha)(1+\alpha)}\right)bb^{\top}C$$

The last equality follows from lemma A.3, where we have: $\mathbb{E}\Big[\sum_{i=0}^{N-1}\sum_{j=0}^{N-1}\alpha^{i+j+2}y_{N-i}y_{N-j}\boldsymbol{x}_{N-i}\boldsymbol{x}_{N-j}^{\top}\Big] = \Big(\frac{\alpha^2\left(1-\alpha^N\right)^2}{(1-\alpha)^2} + \frac{(d+1)\alpha^2\left(1-\alpha^{2N}\right)}{(1-\alpha)(1+\alpha)}\Big)\boldsymbol{I},$ $\mathbb{E}\Big[\sum_{i=0}^{N-1}\sum_{j=0}^{N-1}\alpha^{i+j+2}y_{N-i}y_{N-j}^2\boldsymbol{x}_{N-i}\Big] = \mathbf{0} \text{ and } \mathbb{E}\Big[\sum_{i=0}^{N-1}\sum_{j=0}^{N-1}\alpha^{i+j+2}y_{N-i}^2y_{N-j}^2\Big] = \Big(\frac{d^2\alpha^2\left(1-\alpha^N\right)^2}{(1-\alpha)^2} + \frac{(2d^2+6d)\alpha^2\left(1-\alpha^{2N}\right)}{(1-\alpha)(1+\alpha)}\Big).$

$$\mathbb{E}\Big[\underbrace{\boldsymbol{w}^{\top}\boldsymbol{x}_{q}}_{\bullet}\cdot\sum_{i=0}^{N-1}\alpha^{i+1}y_{N-i}(\boldsymbol{B}\boldsymbol{x}_{N-i}+y_{N-i}\boldsymbol{b})\boldsymbol{x}_{q}^{\top}\Big]$$

$$=\mathbb{E}\Big[\sum_{i=0}^{N-1}\alpha^{i+1}y_{N-i}(\boldsymbol{B}\boldsymbol{x}_{N-i}+y_{N-i}\boldsymbol{b})\underbrace{\boldsymbol{w}^{\top}\boldsymbol{x}_{q}}_{\bullet}\boldsymbol{x}_{q}^{\top}\Big]$$

$$= \mathbf{B} \mathbb{E} \Big[\sum_{i=0}^{N-1} \alpha^{i+1} y_{N-i} \mathbf{x}_{N-i} \mathbf{w}^{\top} \Big] \mathbb{E} \Big[\mathbf{x}_{q} \mathbf{x}_{q}^{\top} \Big]$$

$$+ \mathbf{b} \mathbb{E} \Big[\sum_{i=0}^{N-1} \alpha^{i+1} y_{N-i}^{2} \mathbf{w}^{\top} \Big] \mathbb{E} \Big[\mathbf{x}_{q} \mathbf{x}_{q}^{\top} \Big]$$

$$= \alpha \Big(\frac{1 - \alpha^{N}}{1 - \alpha} \Big) \mathbf{B}$$

The last equality follows from lemma A.3, where we have: $\mathbb{E}\Big[\sum_{i=0}^{N-1}\alpha^{i+1}y_{N-i}\boldsymbol{x}_{N-i}\boldsymbol{w}^{\top}\Big] = \alpha\Big(\frac{1-\alpha^N}{1-\alpha}\Big)\boldsymbol{I}, \mathbb{E}\Big[\sum_{i=0}^{N-1}\alpha^{i+1}y_{N-i}^2\boldsymbol{w}\Big] = \mathbf{0}, \text{ and } \mathbb{E}\Big[\boldsymbol{x}_q\boldsymbol{x}_q^{\top}\Big] = \boldsymbol{I}.$

$$\mathbb{E}\left[\underbrace{\boldsymbol{x}_{q}^{\top}\boldsymbol{C}^{\top}\sum_{i=0}^{N-1}\alpha^{i+1}y_{N-i}(\boldsymbol{B}\boldsymbol{x}_{N-i}+y_{N-i}\boldsymbol{b})}_{\bullet}\cdot\boldsymbol{C}\boldsymbol{x}_{q}\sum_{i=0}^{N-1}\alpha^{i+1}y_{N-i}\boldsymbol{x}_{N-i}^{\top}\right]$$

$$=\mathbb{E}\left[\boldsymbol{C}\boldsymbol{x}_{q}\underbrace{\left(\boldsymbol{x}_{q}^{\top}\boldsymbol{C}^{\top}\sum_{i=0}^{N-1}\alpha^{i+1}y_{N-i}(\boldsymbol{B}\boldsymbol{x}_{N-i}+y_{N-i}\boldsymbol{b})\right)}_{\bullet}\sum_{i=0}^{N-1}\alpha^{i+1}y_{N-i}\boldsymbol{x}_{N-i}^{\top}\right]$$

$$=\boldsymbol{C}\mathbb{E}\left[\boldsymbol{x}_{q}\boldsymbol{x}_{q}^{\top}\right]\boldsymbol{C}^{\top}\boldsymbol{B}\mathbb{E}\left[\sum_{i=0}^{N-1}\sum_{j=0}^{N-1}\alpha^{i+j+2}y_{N-i}y_{N-j}\boldsymbol{x}_{N-i}\boldsymbol{x}_{N-j}^{\top}\right]$$

$$=\left(\frac{\alpha^{2}\left(1-\alpha^{N}\right)^{2}}{(1-\alpha)^{2}}+\frac{(d+1)\alpha^{2}\left(1-\alpha^{2N}\right)}{(1-\alpha)(1+\alpha)}\right)\boldsymbol{C}\boldsymbol{C}^{\top}\boldsymbol{B}$$

The last equality follows from lemma A.3, where we have: $\mathbb{E}\Big[\sum_{i=0}^{N-1}\sum_{j=0}^{N-1}\alpha^{i+j+2}y_{N-i}y_{N-j}\boldsymbol{x}_{N-i}\boldsymbol{x}_{N-j}^{\top}\Big] = \Big(\frac{\alpha^2\left(1-\alpha^N\right)^2}{(1-\alpha)^2} + \frac{(d+1)\alpha^2\left(1-\alpha^{2N}\right)}{(1-\alpha)(1+\alpha)}\Big)\boldsymbol{I}, \text{ and } \mathbb{E}\Big[\boldsymbol{x}_q\boldsymbol{x}_q^{\top}\Big] = \boldsymbol{I}.$

$$\mathbb{E}\left[\underbrace{\boldsymbol{w}^{\top}\boldsymbol{x}_{q}}_{\bullet}\cdot\boldsymbol{C}\boldsymbol{x}_{q}\sum_{i=0}^{N-1}\alpha^{i+1}y_{N-i}\boldsymbol{x}_{N-i}^{\top}\right]$$

$$=\boldsymbol{C}\mathbb{E}\left[\boldsymbol{x}_{q}\underbrace{\boldsymbol{x}_{q}^{\top}\boldsymbol{w}}_{\bullet}\sum_{i=0}^{N-1}\alpha^{i+1}y_{N-i}\boldsymbol{x}_{N-i}^{\top}\right]$$

$$=\boldsymbol{C}\mathbb{E}\left[\boldsymbol{x}_{q}\boldsymbol{x}_{q}^{\top}\right]\mathbb{E}\left[\sum_{i=0}^{N-1}\alpha^{i+1}y_{N-i}\boldsymbol{w}\boldsymbol{x}_{N-i}^{\top}\right]$$

$$=\alpha\left(\frac{1-\alpha^{N}}{1-\alpha}\right)\boldsymbol{C}$$

The last equality follows from lemma A.3, where we have: $\mathbb{E}\Big[\sum_{i=0}^{N-1}\alpha^{i+1}y_{N-i}\boldsymbol{x}_{N-i}\boldsymbol{w}^{\top}\Big] = \alpha\Big(\frac{1-\alpha^N}{1-\alpha}\Big)\boldsymbol{I}$, and $\mathbb{E}\Big[\boldsymbol{x}_q\boldsymbol{x}_q^{\top}\Big] = \boldsymbol{I}$.

$$\mathbb{E}\Big[\underbrace{\left(\boldsymbol{x}_{q}^{\top}\boldsymbol{C}^{\top}\sum_{i=0}^{N-1}\alpha^{i+1}y_{N-i}(\boldsymbol{B}\boldsymbol{x}_{N-i}+y_{N-i}\boldsymbol{b})\right)}\cdot\boldsymbol{C}\boldsymbol{x}_{q}\cdot\sum_{i=0}^{N-1}\alpha^{i+1}y_{N-i}^{2}\Big]$$

$$\begin{split} &= \mathbb{E}\Big[\boldsymbol{C}\boldsymbol{x}_{q} \underbrace{\left(\boldsymbol{x}_{q}^{\top} \boldsymbol{C}^{\top} \sum_{i=0}^{N-1} \alpha^{i+1} y_{N-i} (\boldsymbol{B} \boldsymbol{x}_{N-i} + y_{N-i} \boldsymbol{b}) \right) \cdot \sum_{i=0}^{N-1} \alpha^{i+1} y_{N-i}^{2}} \\ &= \boldsymbol{C} \mathbb{E}\Big[\boldsymbol{x}_{q} \boldsymbol{x}_{q}^{\top} \Big] \boldsymbol{C}^{\top} \boldsymbol{B} \mathbb{E}\Big[\sum_{i=0}^{N-1} \sum_{j=0}^{N-1} \alpha^{i+j+2} y_{N-i} y_{N-j}^{2} \boldsymbol{x}_{N-i} \Big] \\ &+ \boldsymbol{C} \mathbb{E}\Big[\boldsymbol{x}_{q} \boldsymbol{x}_{q}^{\top} \Big] \boldsymbol{C}^{\top} \boldsymbol{b} \mathbb{E}\Big[\sum_{i=0}^{N-1} \sum_{j=0}^{N-1} \alpha^{i+j+2} y_{N-i}^{2} y_{N-j}^{2} \Big] \\ &= \Big(\frac{d^{2} \alpha^{2} \Big(1 - \alpha^{N}\Big)^{2}}{(1 - \alpha)^{2}} + \frac{(2d^{2} + 6d) \alpha^{2} \Big(1 - \alpha^{2N}\Big)}{(1 - \alpha)(1 + \alpha)} \Big) \boldsymbol{C} \boldsymbol{C}^{\top} \boldsymbol{b} \\ \text{The last equality follows from lemma A.3, where we have:} \\ \mathbb{E}\Big[\sum_{i=0}^{N-1} \sum_{j=0}^{N-1} \alpha^{i+j+2} y_{N-i} y_{N-j}^{2} \boldsymbol{x}_{N-i} \Big] &= \mathbf{0}, \quad \mathbb{E}\Big[\sum_{i=0}^{N-1} \sum_{j=0}^{N-1} \alpha^{i+j+2} y_{N-i}^{2} y_{N-j}^{2} \Big] \\ &= \Big(\frac{d^{2} \alpha^{2} \Big(1 - \alpha^{N}\Big)^{2}}{(1 - \alpha)^{2}} + \frac{(2d^{2} + 6d) \alpha^{2} \Big(1 - \alpha^{2N}\Big)}{(1 - \alpha)(1 + \alpha)} \Big), \text{ and } \mathbb{E}\Big[\boldsymbol{x}_{q} \boldsymbol{x}_{q}^{\top}\Big] = \boldsymbol{I}. \\ &\mathbb{E}\Big[\underbrace{\boldsymbol{w}^{\top} \boldsymbol{x}_{q} \cdot \boldsymbol{C} \boldsymbol{x}_{q} \cdot \sum_{i=0}^{N-1} \alpha^{i+1} y_{N-i}^{2} \Big]}_{\boldsymbol{I} = \boldsymbol{C}} \\ &= \mathbb{E}\Big[\boldsymbol{C} \boldsymbol{x}_{q} \boldsymbol{x}_{q}^{\top} \boldsymbol{w} \cdot \sum_{i=0}^{N-1} \alpha^{i+1} y_{N-i}^{2} \Big] \\ &= \boldsymbol{C} \mathbb{E}\Big[\boldsymbol{x}_{q} \boldsymbol{x}_{q}^{\top}\Big] \mathbb{E}\Big[\sum_{i=0}^{N-1} \alpha^{i+1} y_{N-i}^{2} \boldsymbol{w} \Big] \\ &= \boldsymbol{C} \mathbb{E}\Big[\boldsymbol{x}_{q} \boldsymbol{x}_{q}^{\top}\Big] \mathbb{E}\Big[\sum_{i=0}^{N-1} \alpha^{i+1} y_{N-i}^{2} \boldsymbol{w} \Big] \\ &= \boldsymbol{0} \end{aligned}$$

The last equality follows from lemma A.3, where we have $\mathbb{E}\Big[\sum_{i=0}^{N-1}\alpha^{i+1}y_{N-i}^2w\Big]=\mathbf{0}$

$$\mathbb{E}\left[\underbrace{\boldsymbol{x}_{q}^{\top}\boldsymbol{C}^{\top}\left(\sum_{i=0}^{N-1}\alpha^{i+1}y_{N-i}(\boldsymbol{B}\boldsymbol{x}_{N-i}+y_{N-i}\boldsymbol{b}+\boldsymbol{b}_{B})\right)}_{\bullet}\cdot\boldsymbol{C}\boldsymbol{x}_{q}\cdot\sum_{i=0}^{N-1}\alpha^{i+1}y_{N-i}\right]$$

$$=\mathbb{E}\left[\boldsymbol{C}\boldsymbol{x}_{q}\cdot\boldsymbol{x}_{q}^{\top}\boldsymbol{C}^{\top}\left(\sum_{i=0}^{N-1}\alpha^{i+1}y_{N-i}(\boldsymbol{B}\boldsymbol{x}_{N-i}+y_{N-i}\boldsymbol{b}+\boldsymbol{b}_{B})\right)\cdot\sum_{j=0}^{N-1}\alpha^{j+1}y_{N-j}\right]$$

$$=\boldsymbol{C}\mathbb{E}\left[\boldsymbol{x}_{q}\boldsymbol{x}_{q}^{\top}\right]\boldsymbol{C}^{\top}\boldsymbol{B}\mathbb{E}\left[\sum_{i=0}^{N-1}\sum_{j=0}^{N-1}\alpha^{i+j+2}y_{N-i}y_{N-j}\boldsymbol{x}_{N-i}\right]$$

$$+\boldsymbol{C}\mathbb{E}\left[\boldsymbol{x}_{q}\boldsymbol{x}_{q}^{\top}\right]\boldsymbol{C}^{\top}\boldsymbol{b}\mathbb{E}\left[\sum_{i=0}^{N-1}\sum_{j=0}^{N-1}\alpha^{i+j+2}y_{N-i}^{2}y_{N-j}\right]$$

$$+\boldsymbol{C}\mathbb{E}\left[\boldsymbol{x}_{q}\boldsymbol{x}_{q}^{\top}\right]\boldsymbol{C}^{\top}\boldsymbol{b}_{B}\mathbb{E}\left[\sum_{i=0}^{N-1}\sum_{j=0}^{N-1}\alpha^{i+j+2}y_{N-i}y_{N-j}\right]$$

$$=\frac{d\alpha^{2}\left(1-\alpha^{2N}\right)}{(1-\alpha)(1+\alpha)}\boldsymbol{C}\boldsymbol{C}^{\top}\boldsymbol{b}_{B}$$

The last equality follows from lemma A.3, where we have:
$$\mathbb{E}\left[\sum_{i=0}^{N-1}\sum_{j=0}^{N-1}\alpha^{i+j+2}\boldsymbol{x}_{N-i}\underbrace{\boldsymbol{x}_{N-i}^{\top}\boldsymbol{w}}_{\boldsymbol{y}_{N-j}}\underbrace{\boldsymbol{x}_{N-j}^{\top}\boldsymbol{w}}_{\boldsymbol{y}_{N-j}}\right] = \mathbf{0}, \mathbb{E}\left[\sum_{i=0}^{N-1}\sum_{j=0}^{N-1}\alpha^{i+j+2}y_{N-i}^{2}y_{N-j}\right] = 0,$$

$$\mathbb{E}\left[\sum_{i=0}^{N-1}\sum_{j=0}^{N-1}\alpha^{i+j+2}y_{N-i}y_{N-j}\right] = \frac{d\alpha^{2}\left(1-\alpha^{2N}\right)}{(1-\alpha)(1+\alpha)}, \text{ and } \mathbb{E}\left[\boldsymbol{x}_{q}\boldsymbol{x}_{q}^{\top}\right] = \boldsymbol{I}.$$

$$\mathbb{E}\left[\underbrace{\boldsymbol{w}^{\top}\boldsymbol{x}_{q}\cdot\boldsymbol{C}\boldsymbol{x}_{q}\cdot\sum_{i=0}^{N-1}\alpha^{i+1}y_{N-i}}\right]$$

$$= \mathbb{E}\left[\boldsymbol{C}\boldsymbol{x}_{q}\,\boldsymbol{x}_{q}^{\top}\boldsymbol{w}\cdot\sum_{i=0}^{N-1}\alpha^{i+1}y_{N-i}\right]$$

$$= \boldsymbol{C}\mathbb{E}\left[\boldsymbol{x}_{q}\boldsymbol{x}_{q}^{\top}\right]\mathbb{E}\left[\sum_{i=0}^{N-1}\alpha^{i+1}y_{N-i}\boldsymbol{w}\right]$$

The last equality follows from lemma A.3, where we have $\mathbb{E}\Big[\sum_{i=0}^{N-1}\alpha^{i+1}y_{N-i}\boldsymbol{w}\Big]=\boldsymbol{0}$.

C.4 Proof of claim B.1

This Section presents the bounds for terms $\boldsymbol{b}_i^{\top}(T+1)\boldsymbol{b}_i(T+1)$, $\boldsymbol{c}_i^{\top}(Tt+1)\boldsymbol{c}_i(T+1)$ and $\boldsymbol{b}^{\top}(T+1)\boldsymbol{b}(T+1)$, establishing the property $\mathcal{A}(T+1)$.

Recurring the *Vector-coupled Dynamics* equations of $\boldsymbol{b}_i^{\top}(t+1)\boldsymbol{b}_i(t+1)$, $\boldsymbol{c}_i^{\top}(t+1)\boldsymbol{c}_i(t+1)$ and $\boldsymbol{b}^{\top}(t+1)\boldsymbol{b}(t+1)$ in lemma A.7, we have:

$$\begin{aligned} & \boldsymbol{b}_{i}^{\top}(T+1)\boldsymbol{b}_{i}(T+1) = \boldsymbol{b}_{i}^{\top}(T)\boldsymbol{b}_{i}(T) + 2\eta\Big(\big(\beta_{3} - \beta_{1}\boldsymbol{c}_{i}^{\top}(T)\boldsymbol{b}_{i}(T)\big)\boldsymbol{c}_{i}^{\top}(T)\boldsymbol{b}_{i}(T) \\ & - \beta_{1}\sum_{k\neq i}^{d}\big(\boldsymbol{c}_{k}^{\top}(T)\boldsymbol{b}_{i}(T)\big)^{2}\Big) + \eta^{2} \left\|\bar{\boldsymbol{b}}_{i}(T)\right\|_{2}^{2} \\ & = \boldsymbol{b}_{i}^{\top}(0)\boldsymbol{b}_{i}(0) + \sum_{s=0}^{T}\Big(2\eta\Big(\big(\beta_{3} - \beta_{1}\boldsymbol{c}_{i}^{\top}(s)\boldsymbol{b}_{i}(s)\big)\boldsymbol{c}_{i}^{\top}(s)\boldsymbol{b}_{i}(s) - \beta_{1}\sum_{k\neq i}^{d}\big(\boldsymbol{c}_{k}^{\top}(s)\boldsymbol{b}_{i}(s)\big)^{2}\Big) \\ & + \eta^{2} \left\|\bar{\boldsymbol{b}}_{i}(s)\right\|_{2}^{2}\Big) \\ & = \boldsymbol{b}_{i}^{\top}(0)\boldsymbol{b}_{i}(0) + 2\eta\underbrace{\sum_{s=0}^{T}\big(\beta_{3} - \beta_{1}\boldsymbol{c}_{i}^{\top}(s)\boldsymbol{b}_{i}(s)\big)\boldsymbol{c}_{i}^{\top}(s)\boldsymbol{b}_{i}(s)}_{\text{term II}} - 2\eta\beta_{1}\sum_{k\neq i}^{d}\underbrace{\sum_{s=0}^{T}\big(\boldsymbol{c}_{k}^{\top}(s)\boldsymbol{b}_{i}(s)\big)^{2}}_{\text{term III}} \\ & + \eta^{2}\underbrace{\sum_{s=0}^{T}\left\|\bar{\boldsymbol{b}}_{i}(s)\right\|_{2}^{2}}_{\text{term III}} \\ & \boldsymbol{c}_{i}^{\top}(T+1)\boldsymbol{c}_{i}(T+1) = \boldsymbol{c}_{i}^{\top}(T)\boldsymbol{c}_{i}(T) + 2\eta\Big(\big(\beta_{3} - \beta_{1}\boldsymbol{c}_{i}^{\top}(T)\boldsymbol{b}_{i}(T)\big)\boldsymbol{c}_{i}^{\top}(T)\boldsymbol{b}_{i}(T) \\ & - \beta_{1}\underbrace{\sum_{k\neq i}}_{k\neq i}^{d}\big(\boldsymbol{c}_{i}^{\top}(T)\boldsymbol{b}_{k}(T)\big)^{2} - \beta_{2}\big(\boldsymbol{c}_{i}^{\top}(T)\boldsymbol{b}(T)\big)^{2}\big) + \eta^{2}\left\|\bar{\boldsymbol{c}}_{i}(T)\right\|_{2}^{2} \end{aligned}$$

 $= \boldsymbol{c}_i^{\top}(T)\boldsymbol{c}_i(T) + \sum_{i=0}^{T} \left(2\eta \left(\left(\beta_3 - \beta_1 \boldsymbol{c}_i^{\top}(s)\boldsymbol{b}_i(s)\right) \boldsymbol{c}_i^{\top}(s)\boldsymbol{b}_i(s) - \beta_1 \sum_{i=0}^{d} \left(\boldsymbol{c}_i^{\top}(s)\boldsymbol{b}_k(s)\right)^2 \right) \right)$

$$-\beta_{2}(\boldsymbol{c}_{i}^{\top}(s)\boldsymbol{b}(s))^{2}) + \eta^{2} \|\bar{\boldsymbol{c}}_{i}(s)\|_{2}^{2})$$

$$= \boldsymbol{c}_{i}^{\top}(0)\boldsymbol{c}_{i}(0) + 2\eta \sum_{s=0}^{T} (\beta_{3} - \beta_{1}\boldsymbol{c}_{i}^{\top}(s)\boldsymbol{b}_{i}(s))\boldsymbol{c}_{i}^{\top}(s)\boldsymbol{b}_{i}(s) - 2\eta\beta_{1} \sum_{k\neq i}^{d} \sum_{s=0}^{T} (\boldsymbol{c}_{i}^{\top}(s)\boldsymbol{b}_{k}(s))^{2}$$

$$= \operatorname{term} \ \Pi$$

$$-2\eta\beta_{2} \sum_{s=0}^{T} (\boldsymbol{c}_{i}^{\top}(s)\boldsymbol{b}(s))^{2} + \eta^{2} \sum_{s=0}^{T} \|\bar{\boldsymbol{c}}_{i}(s)\|_{2}^{2}$$

$$= \operatorname{term} \ \Pi$$

$$\boldsymbol{b}^{\top}(T+1)\boldsymbol{b}(T+1) = \boldsymbol{b}^{\top}(T)\boldsymbol{b}(T) - 2\eta \left(\beta_2 \sum_{k=1}^{d} \left(\boldsymbol{c}_k^{\top}(T)\boldsymbol{b}(T)\right)^2\right) + \eta^2 \left\|\bar{\boldsymbol{b}}(T)\right\|_2^2$$

$$= \boldsymbol{b}^{\top}(0)\boldsymbol{b}(0) + 2\eta\beta_2 \sum_{k=1}^{d} \sum_{s=0}^{T} \left(\boldsymbol{c}_k^{\top}(s)\boldsymbol{b}(s)\right)^2 + \eta^2 \sum_{s=0}^{T} \left\|\bar{\boldsymbol{b}}(s)\right\|_2^2$$

$$= \operatorname{term\ IV} \operatorname{IV} \operatorname{term\ VI}$$

To bound terms I - VI, we will use some inequalities from property $\mathcal{B}(t)$ and lemma C.7 as following with $i, j, k \in [1, d], i \neq j$:

$$|\beta_{3} - \beta_{1} \boldsymbol{c}_{i}^{\top}(s)\boldsymbol{b}_{i}(s)| \leq \delta(t) \exp(-\eta\beta_{1}\gamma t)$$

$$|\boldsymbol{c}_{i}^{\top}(t)\boldsymbol{b}_{j}(t)| \leq 2\delta(t) \exp(-\eta\beta_{1}\gamma t)$$

$$|\boldsymbol{c}_{i}^{\top}(t)\boldsymbol{b}(t)| \leq 2\delta(t) \exp(-\eta\beta_{2}\gamma t) + \frac{\delta(t)}{\beta_{2}} \exp(-\eta\beta_{1}\gamma t)$$

$$\left|\bar{\boldsymbol{b}}_{i}(t)^{\top}\bar{\boldsymbol{b}}_{k}(t)\right| \leq 8d_{h}d^{2}\delta(t)^{2} \exp(-\eta\beta_{1}\gamma t),$$

$$\left|\bar{\boldsymbol{c}}_{i}(t)^{\top}\bar{\boldsymbol{c}}_{k}(t)\right| \leq 8d_{h}d^{2}\delta(t)^{2} \exp(-\eta\beta_{1}\gamma t) + 40\beta_{2}^{2}d_{h}\delta(t)^{2} \exp(-\eta\beta_{2}\gamma t),$$

$$\left\|\bar{\boldsymbol{b}}(t)\right\|_{2}^{2} \leq 16d_{h}\beta_{2}^{2}d^{2}\delta(t)^{2} \exp(-\eta\beta_{2}\gamma t) + 2d_{h}d^{2}\delta(t)^{2} \exp(-\eta\beta_{1}\gamma t),$$

Next we begin bounding terms I - VI.

Bound of term I: By $\left|\beta_3 - \beta_1 c_i^{\top}(s) b_i(s)\right| \leq \delta(s) \exp(-\eta \beta_1 \gamma s)$, we have:

$$\begin{aligned} \left| \boldsymbol{c}_{i}^{\top}(s)\boldsymbol{b}_{i}(s) \right| &\leq \frac{\delta(s)\exp(-\eta\beta_{1}\gamma s) + \beta_{3}}{\beta_{1}} \\ &\leq \frac{4\delta(s) + 2\alpha}{\alpha^{2}} \\ &\leq \frac{5\delta(s)}{\alpha^{2}} \\ &\leq 6\delta(s) \end{aligned}$$

The third inequality is by $\delta(s) \geq 2\sqrt{d_h \log(4d(2d+1)/\delta)} \geq 2\alpha = 2\alpha = 2\exp((-\ln 2)/N)$. For the last inequality, as long as $N \geq \frac{2\ln 2}{\ln 6 - \ln 5}$, we have $\frac{5}{\alpha^2} \leq 6$.

$$\left| \sum_{s=0}^{T} (\beta_3 - \beta_1 \boldsymbol{c}_i^{\top}(s) \boldsymbol{b}_i(s)) \boldsymbol{c}_i^{\top}(s) \boldsymbol{b}_i(s) \right|$$

$$\leq \sum_{s=0}^{T} 6\delta(s)^2 \exp(-\eta \beta_1 \gamma s)$$

$$\leq 6\delta_{\max}^2 \int_{-1}^{\infty} \exp(-\eta \beta_1 \gamma s) ds$$
$$\leq \frac{6\delta_{\max}^2 \exp(\eta \beta_1 \gamma)}{\eta \beta_1 \gamma}$$

The second inequality is due to $\exp(-\eta \beta_1 \gamma s)$ is monotone decreasing.

Bound of term II:

$$\sum_{s=0}^{T} (\boldsymbol{c}_{k}^{\top}(s)\boldsymbol{b}_{i}(s))^{2}$$

$$\leq \sum_{s=0}^{T} (2\delta(s) \exp(-\eta \beta_{1} \gamma s))^{2}$$

$$\leq 4\delta_{\max}^{2} \int_{-1}^{\infty} \exp(-2\eta \beta_{1} \gamma s) ds$$

$$\leq \frac{2\delta_{\max}^{2} \exp(2\eta \beta_{1} \gamma)}{\eta \beta_{1} \gamma}$$

The second inequality is due to $\exp(-2\eta\beta_1\gamma s)$ is monotone decreasing.

Bound of term III:

$$\sum_{s=0}^{T} \left\| \bar{\boldsymbol{b}}_{i}(s) \right\|_{2}^{2}$$

$$\leq \sum_{s=0}^{T} 8d_{h}d^{2}\delta(t)^{2} \exp(-\eta \beta_{1} \gamma t)$$

$$\leq 8d_{h}d^{2}\delta_{\max}^{2} \int_{-1}^{\infty} \exp(-\eta \beta_{1} \gamma s) ds$$

$$\leq \frac{8d_{h}d^{2}\delta_{\max}^{2} \exp(\eta \beta_{1} \gamma)}{\eta \beta_{1} \gamma}$$

The second inequality is due to $\exp(-\eta \beta_1 \gamma s)$ is monotone decreasing.

Bound of term IV:

$$\begin{split} &\sum_{s=0}^{T} \left(\boldsymbol{c}_{i}^{\top}(s)\boldsymbol{b}(s) \right)^{2} \\ &= \sum_{s=0}^{T} \left(2\delta(s) \exp(-\eta\beta_{2}\gamma s) + \frac{\delta(s)}{\beta_{2}} \exp(-\eta\beta_{1}\gamma s) \right)^{2} \\ &\leq \delta_{\max}^{2} \cdot \left(4\sum_{s=0}^{T} \exp(-2\eta\beta_{2}\gamma s) + \frac{4}{\beta_{2}} \sum_{s=0}^{T} \exp(-\eta(\beta_{1}+\beta_{2})\gamma s) + \frac{1}{\beta_{2}^{2}} \sum_{s=0}^{T} \exp(-2\eta\beta_{1}\gamma s) \right) \\ &\leq \delta_{\max}^{2} \cdot \left(4\int_{-1}^{\infty} \exp(-2\eta\beta_{2}\gamma s) ds + \frac{4}{\beta_{2}} \int_{-1}^{\infty} \exp(-\eta(\beta_{1}+\beta_{2})\gamma s) ds \right. \\ &+ \frac{1}{\beta_{2}^{2}} \int_{-1}^{\infty} \exp(-2\eta\beta_{1}\gamma s) ds \right) \\ &= \frac{2\delta_{\max}^{2} \exp(2\eta\beta_{2}\gamma)}{\eta\beta_{2}\gamma} + \frac{4\delta_{\max}^{2} \exp(\eta(\beta_{1}+\beta_{2})\gamma)}{\eta\beta_{2}(\beta_{1}+\beta_{2})\gamma} + \frac{\delta_{\max}^{2} \exp(2\eta\beta_{1}\gamma)}{2\eta\beta_{1}\beta_{2}^{2}\gamma} \\ &\leq \frac{17\delta_{\max}^{2}}{\eta\beta_{2}\gamma} \end{split}$$

The second inequality is due to $\exp(-2\eta\beta_2\gamma s)$, $\exp(-\eta(\beta_1+\beta_2)\gamma s)$ and $\exp(-2\eta\beta_1\gamma s)$ are monotone decreasing. The last inequality is by $\frac{\exp(\eta(\beta_1+\beta_2)\gamma)}{(\beta_1+\beta_2)} \leq 2$ and $\frac{\exp(2\eta\beta_1\gamma)}{2\beta_1\beta_2} \leq 7$ since $\exp(2\eta\beta_1\gamma) \leq \exp(\eta(\beta_1+\beta_2)\gamma) \leq 2$ and $\beta_1+\beta_2 \geq 1$, $\beta_1\beta_2 \geq \frac{1}{7}$.

Bound of term V:

$$\begin{split} &\sum_{s=0}^{T} \left\| \bar{\boldsymbol{c}}_{i}(s) \right\|_{2}^{2} \\ &\leq \sum_{s=0}^{T} \left(8d_{h}d^{2}\delta(t)^{2} \exp(-\eta\beta_{1}\gamma t) + 40\beta_{2}^{2}d_{h}\delta(t)^{2} \exp(-\eta\beta_{2}\gamma t) \right) \\ &\leq 8d_{h}d^{2}\delta_{\max}^{2} \int_{-1}^{\infty} \exp(-\eta\beta_{1}\gamma s) ds + 40\beta_{2}^{2}d_{h}\delta_{\max}^{2} \int_{-1}^{\infty} \exp(-\eta\beta_{2}\gamma s) ds \\ &= \frac{8d_{h}d^{2}\delta_{\max}^{2} \exp(\eta\beta_{1}\gamma)}{\eta\beta_{1}\gamma} + \frac{40\beta_{2}d_{h}\delta_{\max}^{2} \exp(\eta\beta_{2}\gamma)}{\eta\gamma} \\ &\leq \frac{16d_{h}d^{2}\delta_{\max}^{2} \exp(\eta\beta_{1}\gamma)}{\eta\beta_{1}\gamma} + \frac{80\beta_{2}d_{h}\delta_{\max}^{2} \exp(\eta\beta_{2}\gamma)}{\eta\gamma} \end{split}$$

The second inequality is due to $\exp(-\eta \beta_2 \gamma s)$ and $\exp(-\eta \beta_1 \gamma s)$ are monotone decreasing.

Bound of term VI:

$$\sum_{s=0}^{T} \left\| \overline{b}(s) \right\|_{2}^{2}$$

$$\leq \sum_{s=0}^{T} \left(16d_{h}\beta_{2}^{2}d^{2}\delta(t)^{2} \exp(-\eta\beta_{2}\gamma t) + 2d_{h}d^{2}\delta(t)^{2} \exp(-\eta\beta_{1}\gamma t) \right)$$

$$\leq 16d_{h}\beta_{2}^{2}d^{2}\delta_{\max}^{2} \int_{-1}^{\infty} \exp(-\eta\beta_{2}\gamma s) ds + 2d_{h}d^{2}\delta_{\max}^{2} \int_{-1}^{\infty} \exp(-\eta\beta_{1}\gamma s) ds$$

$$\leq \frac{16d_{h}\beta_{2}d^{2}\delta_{\max}^{2} \exp(\eta\beta_{2}\gamma)}{\eta\gamma} + \frac{2d_{h}d^{2}\delta_{\max}^{2} \exp(\eta\beta_{1}\gamma)}{\eta\beta_{1}\gamma}$$

The second inequality is due to $\exp(-\eta \beta_2 \gamma s)$ and $\exp(-\eta \beta_1 \gamma s)$ are monotone decreasing. We next use the bounds of I - VI to bound $\boldsymbol{b}_i^{\top}(T+1)\boldsymbol{b}_i(T+1)$, $\boldsymbol{c}_i^{\top}(T+1)\boldsymbol{c}_i(T+1)$ and $\boldsymbol{b}^{\top}(T+1)\boldsymbol{b}(T+1)$.

Lower bound of $\boldsymbol{b}_i^{\top}(T+1)\boldsymbol{b}_i(T+1)$

$$\begin{split} & \boldsymbol{b}_{i}^{\top}(T+1)\boldsymbol{b}_{i}(T+1) \\ &= \underbrace{\boldsymbol{b}_{i}^{\top}(0)\boldsymbol{b}_{i}(0)}_{\geq \frac{3d_{h}}{4}} + 2\eta \underbrace{\sum_{s=0}^{T} \left(\beta_{3} - \beta_{1}\boldsymbol{c}_{i}^{\top}(s)\boldsymbol{b}_{i}(s)\right)\boldsymbol{c}_{i}^{\top}(s)\boldsymbol{b}_{i}(s) - 2\eta\beta_{1} \sum_{k\neq i}^{d} \underbrace{\sum_{s=0}^{T} \left(\boldsymbol{c}_{k}^{\top}(s)\boldsymbol{b}_{i}(s)\right)^{2}}_{\text{term II}} \\ &+ \eta^{2} \underbrace{\sum_{s=0}^{T} \left\| \overline{\boldsymbol{b}}_{i}(s) \right\|_{2}^{2}}_{\geq 0} \\ &\geq \frac{3d_{h}}{4} - 2\eta \cdot \frac{6\delta_{\max}^{2} \exp(\eta\beta_{1}\gamma)}{\eta\beta_{1}\gamma} - 2\eta\beta_{1}(d-1) \cdot \frac{2\delta_{\max}^{2} \exp(2\eta\beta_{1}\gamma)}{\eta\beta_{1}\gamma} \\ &\geq \frac{3d_{h}}{4} - \frac{12\delta_{\max}^{2} \exp(\eta\beta_{1}\gamma)}{\beta_{1}\gamma} - \frac{4(d-1)\delta_{\max}^{2} \exp(2\eta\beta_{1}\gamma)}{\gamma} \end{split}$$

$$\geq \frac{3d_h}{4} - \frac{2*12*9d_h \log(4d(2d+1)/\delta)}{\beta_1 \frac{1}{2}d_h} - \frac{2*4(d-1)*9d_h \log(4d(2d+1)/\delta)}{\frac{1}{2}d_h} \\ \geq \frac{d_h}{2}$$

The third inequality is by $\delta_{\max} = 3\sqrt{d_h \log(4d(2d+1)/\delta)}$, $\exp(\eta \beta_1 \gamma) \le \exp(2\eta \beta_1 \gamma) \le 2$ and $\gamma = \frac{1}{2}d_h$. The last inequality follows from $d_h = \widetilde{\Omega}(d^2) \ge \left(1728\log(4d(2d+1)/\delta) + 576(d-1)\beta_1\log(4d(2d+1)/\delta)\right)/\beta_1$.

Upper bound of $\boldsymbol{b}_i^{\top}(T+1)\boldsymbol{b}_i(T+1)$

$$\begin{aligned} & \boldsymbol{b}_{i}^{\top}(T+1)\boldsymbol{b}_{i}(T+1) \\ &= \underbrace{\boldsymbol{b}_{i}^{\top}(0)\boldsymbol{b}_{i}(0)}_{\leq \frac{5d_{h}}{4}} - 2\eta \sum_{s=0}^{T} \left(\beta_{3} - \beta_{1}\boldsymbol{c}_{i}^{\top}(s)\boldsymbol{b}_{i}(s)\right)\boldsymbol{c}_{i}^{\top}(s)\boldsymbol{b}_{i}(s) - 2\eta\beta_{1} \sum_{k\neq i}^{d} \sum_{s=0}^{T} \left(\boldsymbol{c}_{k}^{\top}(s)\boldsymbol{b}_{i}(s)\right)^{2} \\ &+ \eta^{2} \sum_{s=0}^{T} \left\| \bar{\boldsymbol{b}}_{i}(s) \right\|_{2}^{2} \\ &\leq \frac{5d_{h}}{4} + 2\eta \cdot \frac{6\delta_{\max}^{2} \exp(\eta\beta_{1}\gamma)}{\eta\beta_{1}\gamma} + \eta^{2} \cdot \frac{8d_{h}d^{2}\delta_{\max}^{2} \exp(\eta\beta_{1}\gamma)}{\eta\beta_{1}\gamma} \\ &\leq \frac{5d_{h}}{4} + \frac{12\delta_{\max}^{2} \exp(\eta\beta_{1}\gamma)}{\beta_{1}\gamma} + \frac{8\eta d_{h}d^{2}\delta_{\max}^{2} \exp(\eta\beta_{1}\gamma)}{\beta_{1}\gamma} \\ &\leq \frac{5d_{h}}{4} + \frac{2*12*9d_{h}\log(4d(2d+1)/\delta)}{\beta_{1}\frac{1}{2}d_{h}} + \frac{2*8\eta d_{h}d^{2}*9d_{h}\log(4d(2d+1)/\delta)}{\beta_{1}\frac{1}{2}d_{h}} \\ &\leq 2d_{h} \end{aligned}$$

The third inequality is by $\delta_{\max} = 3\sqrt{d_h \log(4d(2d+1)/\delta)}$, $\exp(\eta \beta_1 \gamma) \leq 2$ and $\gamma = \frac{1}{2}d_h$. The last inequality follows from

$$d_h = \widetilde{\Omega}(d^2)$$

$$\geq (576 \log(4d(2d+1)/\delta) + 192 \log(4d(2d+1)/\delta))/\beta_1$$

$$\geq (576 \log(4d(2d+1)/\delta) + 384\eta d_h d^2 \log(4d(2d+1)/\delta))/\beta_1$$

Lower bound of $\boldsymbol{c}_i^{\top}(T+1)\boldsymbol{c}_i(T+1)$

$$\begin{split} & \mathbf{c}_{i}^{\top}(T+1)\mathbf{c}_{i}(T+1) \\ & = \underbrace{\mathbf{c}_{i}^{\top}(0)\mathbf{c}_{i}(0)}_{\geq \frac{3d_{h}}{4}} + 2\eta \sum_{s=0}^{T} \left(\beta_{3} - \beta_{1}\mathbf{c}_{i}^{\top}(s)\mathbf{b}_{i}(s)\right)\mathbf{c}_{i}^{\top}(s)\mathbf{b}_{i}(s) - 2\eta\beta_{1} \sum_{k\neq i}^{d} \sum_{s=0}^{T} \left(\mathbf{c}_{i}^{\top}(s)\mathbf{b}_{k}(s)\right)^{2} \\ & = \operatorname{term} \ \mathbf{II} \\ & - 2\eta\beta_{2} \sum_{s=0}^{T} \left(\mathbf{c}_{i}^{\top}(s)\mathbf{b}(s)\right)^{2} + \eta^{2} \sum_{s=0}^{T} \left\| \overline{\mathbf{c}}_{i}(s) \right\|_{2}^{2} \\ & \geq \frac{3d_{h}}{4} - 2\eta \cdot \frac{6\delta_{\max}^{2} \exp(\eta\beta_{1}\gamma)}{\eta\beta_{1}\gamma} - 2\eta\beta_{1}(d-1) \cdot \frac{2\delta_{\max}^{2} \exp(2\eta\beta_{1}\gamma)}{\eta\beta_{1}\gamma} - 2\eta\beta_{2} \cdot \frac{17\delta_{\max}^{2}}{\eta\beta_{2}\gamma} \\ & \geq \frac{3d_{h}}{4} - \frac{2*12*9d_{h} \log(4d(2d+1)/\delta)}{\beta_{1}\frac{1}{2}d_{h}} - \frac{2*4(d-1)*9d_{h} \log(4d(2d+1)/\delta)}{\frac{1}{2}d_{h}} \\ & - \frac{34*9d_{h} \log(4d(2d+1)/\delta)}{\frac{1}{2}d_{h}} \end{split}$$

$$\geq \frac{d_h}{2}$$

The second inequality is by $\delta_{\max} = 3\sqrt{d_h \log(4d(2d+1)/\delta)}$, $\exp(\eta \beta_1 \gamma) \le \exp(2\eta \beta_1 \gamma) \le 2$ and $\gamma = \frac{1}{2}d_h$. The last inequality follows from $d_h = \widetilde{\Omega}(d^2) \ge \left(1728\log(4d(2d+1)/\delta) + (576d+1872)\beta_1\log(4d(2d+1)/\delta)\right)/\beta_1$.

Upper bound of $c_i^{\top}(T+1)c_i(T+1)$

$$\begin{split} & \mathbf{c}_{i}^{\top}(T+1)\mathbf{c}_{i}(T+1) \\ & = \underbrace{\mathbf{c}_{i}^{\top}(0)\mathbf{c}_{i}(0)}_{\leq \frac{5d_{h}}{4}} + 2\eta \underbrace{\sum_{s=0}^{T} \left(\beta_{3} - \beta_{1}\mathbf{c}_{i}^{\top}(s)\mathbf{b}_{i}(s)\right)\mathbf{c}_{i}^{\top}(s)\mathbf{b}_{i}(s)}_{\text{term I}} - 2\eta\beta_{1} \underbrace{\sum_{k\neq i}^{d} \underbrace{\sum_{s=0}^{T} \left(\mathbf{c}_{i}^{\top}(s)\mathbf{b}_{k}(s)\right)^{2}}_{=\text{term II}} + 2\eta\beta_{2} \underbrace{\sum_{s=0}^{T} \left(\mathbf{c}_{i}^{\top}(s)\mathbf{b}(s)\right)^{2} + \eta^{2} \underbrace{\sum_{s=0}^{T} \left\|\bar{\mathbf{c}}_{i}(s)\right\|_{2}^{2}}_{\geq 0} \\ & \leq \frac{5d_{h}}{4} + 2\eta \cdot \frac{6\delta_{\max}^{2} \exp(\eta\beta_{1}\gamma)}{\eta\beta_{1}\gamma} + 2\eta\beta_{1}(d-1) \cdot \frac{2\delta_{\max}^{2} \exp(2\eta\beta_{1}\gamma)}{\eta\beta_{1}\gamma} \\ & + \eta^{2} \cdot \left(\frac{16d_{h}d^{2}\delta_{\max}^{2} \exp(\eta\beta_{1}\gamma)}{\eta\beta_{1}\gamma} + \frac{80\beta_{2}d_{h}\delta_{\max}^{2} \exp(\eta\beta_{2}\gamma)}{\eta\gamma}\right) \\ & \leq \frac{5d_{h}}{4} + \frac{2*12*9d_{h}\log(4d(2d+1)/\delta)}{\beta_{1}\frac{1}{2}d_{h}} + \frac{2*4(d-1)*9d_{h}\log(4d(2d+1)/\delta)}{\frac{1}{2}d_{h}} \\ & + \frac{2*16\eta d_{h}d^{2}*9d_{h}\log(4d(2d+1)/\delta)}{\beta_{1}\frac{1}{2}d_{h}} + \frac{2*80\eta\beta_{2}d_{h}*9d_{h}\log(4d(2d+1)/\delta)}{\frac{1}{2}d_{h}} \\ & \leq 2d_{h} \end{split}$$

The second inequality is by $\delta_{\max} = 3\sqrt{d_h \log(4d(2d+1)/\delta)}$, $\exp(\eta \beta_1 \gamma) \le \exp(2\eta \beta_1 \gamma) \le 2$ and $\gamma = \frac{1}{2}d_h$. The last inequality follows from

$$d_h = \widetilde{\Omega}(d^2)$$

$$\geq 576 \log(4d(2d+1)/\delta)/\beta_1 + 192(d-1)\log(4d(2d+1)/\delta)$$

$$+ 384 \log(4d(2d+1)/\delta)/\beta_1 + 3840 \ln 2 \log(4d(2d+1)/\delta)$$

$$\geq 576 \log(4d(2d+1)/\delta)/\beta_1 + 192(d-1)\log(4d(2d+1)/\delta)$$

$$+ 768\eta d_h d^2 \log(4d(2d+1)/\delta)/\beta_1 + 3840\eta \beta_2 d_h \log(4d(2d+1)/\delta)$$

Lower bound of $\boldsymbol{b}^{\top}(T+1)\boldsymbol{b}(T+1)$

$$\begin{aligned} & \boldsymbol{b}^{\top}(T+1)\boldsymbol{b}(T+1) = \underbrace{\boldsymbol{b}^{\top}(0)\boldsymbol{b}(0)}_{\geq \frac{3d_h}{4}} - 2\eta\beta_2 \sum_{k=1}^{d} \underbrace{\sum_{s=0}^{T} \left(\boldsymbol{c}_k^{\top}(s)\boldsymbol{b}(s)\right)^2}_{=\text{term IV}} + \eta^2 \underbrace{\sum_{s=0}^{T} \left\| \overline{\boldsymbol{b}}(s) \right\|_2^2}_{\geq 0} \\ & \geq \frac{3d_h}{4} - 2\eta\beta_2 d \cdot \frac{17\delta_{\text{max}}^2}{\eta\beta_2 \gamma} \\ & \geq \frac{3d_h}{4} - \frac{34d\delta_{\text{max}}^2}{\gamma} \\ & \geq \frac{3d_h}{4} - \frac{34d * 9d_h \log(4d(2d+1)/\delta)}{\frac{1}{2}d_h} \\ & \geq \frac{d_h}{2} \end{aligned}$$

The third inequality is by $\delta_{\max} = 3\sqrt{d_h \log(4d(2d+1)/\delta)}$ and $\gamma = \frac{1}{2}d_h$. The last inequality follows from $d_h = \widetilde{\Omega}(d^2) \ge 2448d \log(4d(2d+1)/\delta)$.

Upper bound of $\boldsymbol{b}^{\top}(T+1)\boldsymbol{b}(T+1)$

$$\begin{aligned} & \boldsymbol{b}^{\top}(T+1)\boldsymbol{b}(T+1) = \underbrace{\boldsymbol{b}^{\top}(0)\boldsymbol{b}(0)}_{\leq \frac{5d_h}{4}} - 2\eta\beta_2 \sum_{k=1}^{d} \sum_{\frac{s=0}{4}}^{T} \left(\boldsymbol{c}_k^{\top}(s)\boldsymbol{b}(s)\right)^2 + \eta^2 \sum_{\frac{s=0}{4}}^{T} \left\|\bar{\boldsymbol{b}}(s)\right\|_2^2 \\ & \leq \frac{5d_h}{4} + 2\eta\beta_2 d \cdot \frac{17\delta_{\max}^2}{\eta\beta_2\gamma} + \eta^2 \cdot \left(\frac{16d_h\beta_2 d^2\delta_{\max}^2 \exp(\eta\beta_2\gamma)}{\eta\gamma} + \frac{2d_h d^2\delta_{\max}^2 \exp(\eta\beta_1\gamma)}{\eta\beta_1\gamma}\right) \\ & \leq \frac{5d_h}{4} + \frac{34d\delta_{\max}^2}{\gamma} + \frac{2*16\eta d_h\beta_2 d^2\delta_{\max}^2}{\gamma} + \frac{2*2\eta d_h d^2\delta_{\max}^2}{\beta_1\gamma} \\ & \leq \frac{5d_h}{4} + \frac{34d*9d_h \log(4d(2d+1)/\delta)}{\frac{1}{2}d_h} + \frac{32\eta d_h\beta_2 d^2*9d_h \log(4d(2d+1)/\delta)}{\frac{1}{2}d_h} \\ & + \frac{4\eta d_h d^2*9d_h \log(4d(2d+1)/\delta)}{\beta_1 \frac{1}{2}d_h} \\ & \leq 2d_h \end{aligned}$$

The second inequality is by $\exp(\eta \beta_1 \gamma) \leq \exp(\eta \beta_2 \gamma) \leq 2$ The third inequality is by $\delta_{\max} = 3\sqrt{d_h \log(4d(2d+1)/\delta)}$ and $\gamma = \frac{1}{2}d_h$. The last inequality follows from $d_h = \widetilde{\Omega}(d^2)$

$$\geq 816d\log(4d(2d+1)/\delta) + 768\ln 2d^2\log(4d(2d+1)/\delta) + 48\log(4d(2d+1)/\delta)/\beta_1$$

$$\geq 816d \log(4d(2d+1)/\delta) + 768\eta\beta_2 d_h d^2 \log(4d(2d+1)/\delta) + 96\eta d_h d^2 \log(4d(2d+1)/\delta)/\beta_1$$

C.5 Proof of claim B.2

This Section presents the exponential decay bounds for terms $(\beta_3 - \beta_1 c_i^{\top}(T+1)b_i(T+1))$, $c_i^{\top}(T+1)b_i(T+1)$ and $c_i^{\top}(T+1)b(T+1)$, establishing the property $\mathcal{B}(T+1)$.

Bound of
$$(\beta_3 - \beta_1 \boldsymbol{c}_i^\top (T+1) \boldsymbol{b}_i (T+1))$$

Recall the following equation from lemma A.7.

$$\boldsymbol{c}_{i}^{\top}(t+1)\boldsymbol{b}_{i}(t+1) = \boldsymbol{c}_{i}^{\top}(t)\boldsymbol{b}_{i}(t) + \eta \Big((\beta_{3} - \beta_{1}\boldsymbol{c}_{i}^{\top}(t)\boldsymbol{b}_{i}(t))\boldsymbol{b}_{i}^{\top}(t)\boldsymbol{b}_{i}(t) - \beta_{1}\sum_{k\neq i}^{d}\boldsymbol{c}_{i}^{\top}(t)\boldsymbol{b}_{k}(t) \cdot \boldsymbol{b}_{k}^{\top}(t)\boldsymbol{b}_{i}(t) \\ - \beta_{2}\boldsymbol{c}_{i}^{\top}(t)\boldsymbol{b}(t) \cdot \boldsymbol{b}_{i}^{\top}(t)\boldsymbol{b}(t) + (\beta_{3} - \beta_{1}\boldsymbol{c}_{i}^{\top}(t)\boldsymbol{b}_{i}(t))\boldsymbol{c}_{i}^{\top}(t)\boldsymbol{c}_{i}(t) \\ - \beta_{1}\sum_{k\neq i}^{d}\boldsymbol{c}_{k}^{\top}(t)\boldsymbol{b}_{i}(t) \cdot \boldsymbol{c}_{i}^{\top}(t)\boldsymbol{c}_{k}(t) + \eta^{2}\bar{\boldsymbol{c}}_{i}^{\top}(t)\bar{\boldsymbol{b}}_{i}(t)$$

Based on the above equation, we have:

$$\left| \left(\beta_{3} - \beta_{1} \boldsymbol{c}_{i}^{\top} (T+1) \boldsymbol{b}_{i} (T+1) \right) \right| = \left| \underline{\beta_{3} - \beta_{1} \boldsymbol{c}_{i}^{\top} (T) \boldsymbol{b}_{i} (T)} \right|
- \eta \beta_{1} \left(\left(\underline{\beta_{3} - \beta_{1} \boldsymbol{c}_{i}^{\top} (T) \boldsymbol{b}_{i} (T) \right) \boldsymbol{b}_{i}^{\top} (T) \boldsymbol{b}_{i} (T) - \beta_{1} \sum_{k \neq i}^{d} \boldsymbol{c}_{i}^{\top} (T) \boldsymbol{b}_{k} (T) \cdot \boldsymbol{b}_{k}^{\top} (T) \boldsymbol{b}_{i} (T) \right)
- \beta_{2} \boldsymbol{c}_{i}^{\top} (T) \boldsymbol{b} (T) \cdot \boldsymbol{b}_{i}^{\top} (T) \boldsymbol{b} (T) + \left(\underline{\beta_{3} - \beta_{1} \boldsymbol{c}_{i}^{\top} (T) \boldsymbol{b}_{i} (T) \right) \boldsymbol{c}_{i}^{\top} (T) \boldsymbol{c}_{i} (T)
- \beta_{1} \sum_{k \neq i}^{d} \boldsymbol{c}_{k}^{\top} (T) \boldsymbol{b}_{i} (T) \cdot \boldsymbol{c}_{i}^{\top} (T) \boldsymbol{c}_{k} (T) \right) - \eta^{2} \beta_{1} \bar{\boldsymbol{c}}_{i}^{\top} (T) \bar{\boldsymbol{b}}_{i} (T) \right|
= \left| \left(1 - \eta \beta_{1} \left(\boldsymbol{b}_{i}^{\top} (T) \boldsymbol{b}_{i} (T) + \boldsymbol{c}_{i}^{\top} (T) \boldsymbol{c}_{i} (T) \right) \right) \left(\underline{\beta_{3} - \beta_{1} \boldsymbol{c}_{i}^{\top} (T) \boldsymbol{b}_{i} (T) \right) \right|
+ \eta \beta_{1}^{2} \sum_{k \neq i}^{d} \boldsymbol{c}_{i}^{\top} (T) \boldsymbol{b}_{k} (T) \cdot \boldsymbol{b}_{k}^{\top} (T) \boldsymbol{b}_{i} (T) + \eta \beta_{1}^{2} \sum_{k \neq i}^{d} \boldsymbol{c}_{k}^{\top} (T) \boldsymbol{b}_{i} (T) \cdot \boldsymbol{c}_{i}^{\top} (T) \boldsymbol{c}_{k} (T)
+ \eta \beta_{1} \beta_{2} \boldsymbol{c}_{i}^{\top} (T) \boldsymbol{b} (T) \cdot \boldsymbol{b}_{i}^{\top} (T) \boldsymbol{b} (T) - \eta^{2} \beta_{1} \bar{\boldsymbol{c}}_{i}^{\top} (T) \bar{\boldsymbol{b}}_{i} (T) \right|$$

The term $\beta_3 - \beta_1 \boldsymbol{c}_i^\top(T) \boldsymbol{b}_i(T)$ is highlighted with underline, and we collect its *negative feedback* terms together. The factor $\left(1 - \eta \beta_1 \left(\boldsymbol{b}_i^\top(T) \boldsymbol{b}_i(T) + \boldsymbol{c}_i^\top(T) \boldsymbol{c}_i(T)\right)\right) \leq 1$ will drive $\left(\beta_3 - \beta_1 \boldsymbol{c}_i^\top(T + 1) \boldsymbol{b}_i(T + 1)\right)$ to converge to zero.

By Recurring (Eq. (31)) from 0 to T, we have:

$$\left| \left(\beta_{3} - \beta_{1} \boldsymbol{c}_{i}^{\top} (T+1) \boldsymbol{b}_{i} (T+1) \right) \right| \\
= \left| \prod_{s=0}^{T} \left(1 - \eta \beta_{1} \left(\boldsymbol{b}_{i}^{\top} (s) \boldsymbol{b}_{i} (s) + \boldsymbol{c}_{i}^{\top} (s) \boldsymbol{c}_{i} (s) \right) \right) \left(\beta_{3} - \beta_{1} \boldsymbol{c}_{i}^{\top} (0) \boldsymbol{b}_{i} (0) \right) \right. \\
+ \sum_{s=0}^{T} \prod_{s'=s+1}^{T} \left(1 - \eta \beta_{1} \left(\boldsymbol{b}_{i}^{\top} (s') \boldsymbol{b}_{i} (s') + \boldsymbol{c}_{i}^{\top} (s') \boldsymbol{c}_{i} (s') \right) \right) \\
\cdot \left(\eta \beta_{1}^{2} \sum_{k \neq i}^{d} \boldsymbol{c}_{i}^{\top} (s) \boldsymbol{b}_{k} (s) \cdot \boldsymbol{b}_{k}^{\top} (s) \boldsymbol{b}_{i} (s) + \eta \beta_{1}^{2} \sum_{k \neq i}^{d} \boldsymbol{c}_{k}^{\top} (s) \boldsymbol{b}_{i} (s) \cdot \boldsymbol{c}_{i}^{\top} (s) \boldsymbol{c}_{k} (s) \right. \\
+ \underbrace{\eta \beta_{1} \beta_{2} \boldsymbol{c}_{i}^{\top} (s) \boldsymbol{b} (s) \cdot \boldsymbol{b}_{i}^{\top} (s) \boldsymbol{b} (s)}_{\bullet} - \underbrace{\eta^{2} \beta_{1} \bar{\boldsymbol{c}}_{i}^{\top} (s) \bar{\boldsymbol{b}}_{i} (s)}_{\diamond} \right) \right|$$

Here $\prod_{s=0}^T \left(1 - \eta \beta_1 \left(\boldsymbol{b}_i^\top(s) \boldsymbol{b}_i(s) + \boldsymbol{c}_i^\top(s) \boldsymbol{c}_i(s)\right)\right) \leq (1 - 2\eta \beta_1 \gamma)^{T+1}$ since $\gamma \leq \boldsymbol{b}_i^\top(s) \boldsymbol{b}_i(s), \boldsymbol{c}_i^\top(s) \boldsymbol{c}_i(s)$. Besides, from property $\mathcal{B}(0), \dots, \mathcal{B}(T)$ and lemma C.7 we know that $\boldsymbol{c}_i^\top(s) \boldsymbol{b}_k(s), \boldsymbol{c}_i^\top(s) \boldsymbol{b}(s)$ and $\bar{\boldsymbol{c}}_i^\top(s) \bar{\boldsymbol{b}}_i(s)$ have bounds with exponential decreasing rate. Therefore, it is easy to derive an exponential decreasing upper bound for $\left|\left(\beta_3 - \beta_1 \boldsymbol{c}_i^\top(T+1) \boldsymbol{b}_i(T+1)\right)\right|$.

By substituting the bounds of $\boldsymbol{c}_i^{\top}(s)\boldsymbol{b}_k(s)$, $\boldsymbol{c}_i^{\top}(s)\boldsymbol{b}(s)$, $\bar{\boldsymbol{c}}_i^{\top}(s)\bar{\boldsymbol{b}}_i(s)$, $\boldsymbol{b}_k^{\top}(s)\boldsymbol{b}_i(s)$, $\boldsymbol{c}_i^{\top}(s)\boldsymbol{c}_k(s)$ and $\boldsymbol{b}_i^{\top}(s)\boldsymbol{b}(s)$, we have:

$$\left| \left(\beta_{3} - \beta_{1} \mathbf{c}_{i}^{\top} (T+1) \mathbf{b}_{i} (T+1) \right) \right|$$

$$\leq \left(1 - 2\eta \beta_{1} \gamma \right)^{T+1} \left| \beta_{3} - \beta_{1} \mathbf{c}_{i}^{\top} (0) \mathbf{b}_{i} (0) \right|$$

$$+ \sum_{s=0}^{T} (1 - 2\eta \beta_{1} \gamma)^{T-s} \cdot \left(2\eta \beta_{1}^{2} (d-1) \cdot 2\delta(s)^{2} \exp(-\eta \beta_{1} \gamma s) \right)$$

$$+ \eta \beta_{1} \beta_{2} \cdot \left(2\delta(s)^{2} \exp(-\eta \beta_{2} \gamma s) + \frac{\delta(s)^{2}}{\beta_{2}} \exp(-\eta \beta_{1} \gamma s) \right)$$

$$+ \eta^{2} \beta_{1} \left(8d_{h} d^{2} \delta(t)^{2} \exp(-\eta \beta_{1} \gamma t) + 8\beta_{2} d_{h} d\delta(t)^{2} \exp(-\eta \beta_{2} \gamma t) \right) \right)$$

$$\downarrow \diamond$$

$$(33)$$

The notations \spadesuit , \clubsuit and \diamondsuit highlight the corresponding terms between (Eq. (32)) and (Eq. (33)) for refference.

We further have the following:

$$\begin{split} &\left|\left(\beta_{3}-\beta_{1}c_{i}^{\top}(T+1)\boldsymbol{b}_{i}(T+1)\right)\right| \\ &\leq (1-2\eta\beta_{1}\gamma)^{T+1}\left|\beta_{3}-\beta_{1}c_{i}^{\top}(0)\boldsymbol{b}_{i}(0)\right| \\ &+\sum_{s=0}^{T}(1-2\eta\beta_{1}\gamma)^{T-s}\cdot\left(\underline{2\eta\beta_{1}^{2}(d-1)\cdot2\delta(s)^{2}\exp(-\eta\beta_{1}\gamma s)}\right) \\ &+\underline{\eta\beta_{1}\beta_{2}\cdot\left(2\delta(s)^{2}\exp(-\eta\beta_{2}\gamma s)+\frac{\delta(s)^{2}}{\beta_{2}}\exp(-\eta\beta_{1}\gamma s)\right)} \\ &+\underline{\eta^{2}\beta_{1}\left(8d_{h}d^{2}\delta(t)^{2}\exp(-\eta\beta_{1}\gamma t)+8\beta_{2}d_{h}d\delta(t)^{2}\exp(-\eta\beta_{2}\gamma t)\right)}\right) \\ &\leq \exp(-2\eta\beta_{1}\gamma(T+1))\left|\beta_{3}-\beta_{1}c_{i}^{\top}(0)\boldsymbol{b}_{i}(0)\right| \\ &+\sum_{s=0}^{T}\exp(2\eta\beta_{1}\gamma(s-T))\cdot\left(\left(4\eta\beta_{1}^{2}(d-1)\delta(s)^{2}+\eta\beta_{1}\delta(s)^{2}+8\eta^{2}\beta_{1}d_{h}d^{2}\delta(t)^{2}\right)\exp(-\eta\beta_{1}\gamma s) \\ &+\left(2\eta\beta_{1}\beta_{2}\delta(s)^{2}+8\eta^{2}\beta_{1}\beta_{2}d_{h}d\delta(s)^{2}\right)\exp(-\eta\beta_{2}\gamma s)\right) \\ &\leq (\beta_{3}+\beta_{1}\delta(0))\exp(-2\eta\beta_{1}\gamma(T+1)) \\ &+\left(4\eta\beta_{1}^{2}(d-1)\delta(T)^{2}+\eta\beta_{1}\delta(T)^{2}+8\eta^{2}\beta_{1}d_{h}d^{2}\delta(T)^{2}\right)\cdot\frac{2}{\eta\beta_{1}\gamma}\exp(-\eta\beta_{1}\gamma(T+1)) \\ &+\left(2\eta\beta_{1}\beta_{2}\delta(T)^{2}+8\eta^{2}\beta_{1}\beta_{2}d_{h}d\delta(T)^{2}\right)\cdot\frac{3}{\eta\beta_{2}\gamma}\exp(-\eta\beta_{1}\gamma(T+1)) \\ &\leq \left(\frac{\beta_{3}}{\delta(0)}+\beta_{1}+\frac{8\beta_{1}(d-1)\delta(T)}{\gamma}+\frac{2\delta(T)}{\gamma}+\frac{16\eta d_{h}d^{2}\delta(T)}{\gamma}+\frac{6\beta_{1}\delta(T)}{\gamma}+\frac{24\eta\beta_{1}d_{h}d\delta(T)}{\gamma}\right) \\ &\cdot\delta(T)\cdot\exp(-\eta\beta_{1}\gamma(T+1)) \end{aligned}$$

The second inequality is derived by factoring out the factors $\exp(-\eta\beta_1\gamma s)$ and $\exp(-\eta\beta_2\gamma s)$. The third inequality is due to $\sum_{s=0}^T \exp(2\eta\beta_1\gamma(s-T))\cdot \exp(-\eta\beta_1\gamma s) \leq \frac{2}{\eta\beta_1\gamma}\exp(-\eta\beta_1\gamma(T+1))$ and $\sum_{s=0}^T \exp(2\eta\beta_1\gamma(s-T))\cdot \exp(-\eta\beta_2\gamma s) \leq \frac{3}{\eta\beta_2\gamma}\exp(-\eta\beta_1\gamma(T+1))$ in lemma C.5. The fourth inequality is by $\delta(0) \leq \delta(T)$, $\exp(-2\eta\beta_1\gamma(T+1)) \leq \exp(-\eta\beta_1\gamma(T+1))$, and we consider $\beta_3 = \frac{\beta_3}{\delta(0)}\cdot\delta(0) \leq \frac{\beta_3}{\delta(0)}\cdot\delta(T)$. The fifth inequality is by proving $\left(\frac{\beta_3}{\delta(0)}+\beta_1+\frac{8\beta_1(d-1)\delta(T)}{\gamma}+\frac{2\delta(T)}{\gamma}+\frac{16\eta d_h d^2\delta(T)}{\gamma}+\frac{6\beta_1\delta(T)}{\gamma}+\frac{24\eta\beta_1 d_h d\delta(T)}{\gamma}\right) \leq 1$ as follows:

$$\frac{\beta_3}{\delta(0)} + \beta_1 + \frac{8\beta_1(d-1)\delta(T)}{\gamma} + \frac{2\delta(T)}{\gamma} + \frac{16\eta d_h d^2 \delta(T)}{\gamma} + \frac{6\beta_1 \delta(T)}{\gamma} + \frac{24\eta \beta_1 d_h d\delta(T)}{\gamma} \\
\leq \frac{\beta_3}{\delta(0)} + \frac{3}{4} + \frac{\delta(T)}{\gamma} \cdot \left(8\beta_1(d-1) + 2 + 16\eta d_h d^2 + 6\beta_1 + 24\eta \beta_1 d_h d\right) \\
\leq \frac{\beta_3}{2\sqrt{d_h \log(4d(2d+1)/\delta)}} + \frac{3}{4} \\
+ \frac{3\sqrt{d_h \log(4d(2d+1)/\delta)}}{\frac{1}{2}d_h} \cdot \left(8\beta_1(d-1) + 2 + 16\eta d_h d^2 + 6\beta_1 + 24\eta \beta_1 d_h d\right) \\
\leq 1$$

The first inequality is by $\beta_1 \leq \frac{3}{4}$. The second inequality is by $\delta(0) \geq 2\sqrt{d_h \log(4d(2d+1)/\delta)}$, $\delta(T) \leq 3\sqrt{d_h \log(4d(2d+1)/\delta)}$ and $\gamma = \frac{1}{2}d_h$. The last inequality hold as long as $d_h = \frac{1}{2}d_h$.

$$\begin{split} \widetilde{\Omega}(d^2) &\geq \Big(\frac{1}{\sqrt{\log(4d(2d+1)/\delta)}} + 24\sqrt{\log(4d(2d+1)/\delta)} \Big(8\beta_1(d-1) + 2 + 8 + 6\beta_1 + 12\eta\beta_1/d\Big)\Big)^2 \geq \\ &\Big(\frac{1}{\sqrt{\log(4d(2d+1)/\delta)}} + 24\sqrt{\log(4d(2d+1)/\delta)} \Big(8\beta_1(d-1) + 2 + 16\eta d_h d^2 + 6\beta_1 + 24\eta\beta_1 d_h d\Big)\Big)^2. \end{split}$$
 Therefore, we have

$$\left| \left(\beta_3 - \beta_1 \boldsymbol{c}_i^\top (T+1) \boldsymbol{b}_i (T+1) \right) \right| \le \delta(T) \exp(-\eta \beta_1 \gamma (T+1)) \le \delta(T+1) \exp(-\eta \beta_1 \gamma (T+1))$$
 (35)

where the last inequality is by $\delta(T) \leq \delta(T+1)$.

The proof for the bounds of $c_i^{\top}(T+1)b_j(T+1)$ and $c_i^{\top}(T+1)b(T+1)$ are similar to that of $(\beta_3 - \beta_1 c_i^{\top}(T+1)b_i(T+1))$. We presents the calculation as follows.

$$\left(\beta_3 - \beta_1 \boldsymbol{c}_i^{\top}(T+1)\boldsymbol{b}_i(T+1)\right)$$
. We presents the calculation as follows.
Bound of $\boldsymbol{c}_i^{\top}(T+1)\boldsymbol{b}_j(T+1)$

$$\left|\boldsymbol{c}_i^{\top}(T+1)\boldsymbol{b}_j(T+1)\right|$$

$$\begin{vmatrix} c_i^\top (T+1)b_j(T+1) \end{vmatrix} = \begin{vmatrix} c_i^\top (T)b_j(T) + \eta \Big((\beta_3 - \beta_1 c_i^\top (T)b_i(T))b_i^\top (T)b_j(T) - \beta_1 \sum_{k \neq i}^d c_i^\top (T)b_k(T) \cdot b_k^\top (T)b_j(T) \\ - \beta_2 c_i^\top (T)b(T) \cdot b_j^\top (T)b(T) + (\beta_3 - \beta_1 c_j^\top (T)b_j(T))c_i^\top (T)c_j(T) \\ - \beta_1 \sum_{k \neq j}^d c_k^\top (T)b_j(T) \cdot c_i^\top (T)c_k(T) \Big) + \eta^2 \overline{c}_i^\top (T)\overline{b}_j(T) \end{vmatrix} \\ = \begin{vmatrix} \left(1 - \eta\beta_1 (c_i^\top (T)c_i(T) + b_j^\top (T)b_j(T))\right) c_i^\top (T)b_j(T) \\ + \eta (\beta_3 - \beta_1 c_i^\top (T)b_i(T))b_i^\top (T)b_j(T) - \eta\beta_1 \sum_{k \neq i, k \neq j}^d c_i^\top (T)b_k(T) \cdot b_k^\top (T)b_j(T) \\ - \eta\beta_2 c_i^\top (T)b(T) \cdot b_j^\top (T)b(T) + \eta (\beta_3 - \beta_1 c_j^\top (T)b_j(T))c_i^\top (T)c_j(T) \\ - \eta\beta_1 \sum_{k \neq i, k \neq j}^d c_k^\top (T)b_j(T) \cdot c_i^\top (T)c_k(T) + \eta^2 \overline{c}_i^\top (T)\overline{b}_j(T) \Big| \\ = \begin{vmatrix} \prod_{s=0}^T \left(1 - \eta\beta_1 (c_i^\top (s)c_i(s) + b_j^\top (s)b_j(s))\right) c_i^\top (0)b_j(0) \\ + \sum_{s=0}^T \prod_{s'=s+1}^T \left(1 - \eta\beta_1 (c_i^\top (s')c_i(s') + b_j^\top (s')b_j(s'))\right) \\ \cdot \left(\eta (\beta_3 - \beta_1 c_i^\top (s)b_i(s))b_i^\top (s)b_j(s) + \eta (\beta_3 - \beta_1 c_j^\top (s)b_j(s))c_i^\top (s)c_j(s) \right) \\ - \eta\beta_1 \sum_{k \neq i, k \neq j}^d c_i^\top (s)b_k(s) \cdot b_k^\top (s)b_j(s) - \eta\beta_1 \sum_{k \neq i, k \neq j}^d c_k^\top (s)b_j(s) \cdot c_i^\top (s)c_k(s) \\ - \eta\beta_2 c_i^\top (s)b(s) \cdot b_j^\top (s)b(s) + \eta^2 \overline{c}_i^\top (s)\overline{b}_j(s) \Big) \Big|$$

$$-\underbrace{\eta \beta_{2} \boldsymbol{c}_{i}^{\top}(s) \boldsymbol{b}(s) \cdot \boldsymbol{b}_{j}^{\top}(s) \boldsymbol{b}(s)}_{\diamondsuit} + \underbrace{\eta^{2} \bar{\boldsymbol{c}}_{i}^{\top}(s) \bar{\boldsymbol{b}}_{j}(s)}_{\heartsuit} \Big) \Big|$$

$$\leq (1 - 2\eta \beta_{1} \gamma)^{T+1} \Big| \boldsymbol{c}_{i}^{\top}(0) \boldsymbol{b}_{j}(0) \Big|$$

$$+ \sum_{s=0}^{T} (1 - 2\eta \beta_{1} \gamma)^{T-s} \Big(\underbrace{2\eta \delta(s)^{2} \exp(-\eta \beta_{1} \gamma s)}_{\blacktriangle} \Big)$$

$$+\underbrace{2\eta\beta_1(d-2)\cdot 2\delta(s)^2\exp(-\eta\beta_1\gamma s)}_{\bullet}$$

$$+ \eta \beta_{2} \cdot \left(2\delta(s)^{2} \exp(-\eta \beta_{2} \gamma s) + \frac{\delta(s)^{2}}{\beta_{2}} \exp(-\eta \beta_{1} \gamma s)\right)$$

$$+ \eta^{2} \left(8d_{h}d^{2}\delta(t)^{2} \exp(-\eta \beta_{1} \gamma t) + 8\beta_{2}d_{h}d\delta(t)^{2} \exp(-\eta \beta_{2} \gamma t)\right)$$

$$\leq \exp(-2\eta \beta_{1} \gamma (T+1)) \Big| c_{i}^{\top}(0) b_{j}(0) \Big|$$

$$+ \sum_{s=0}^{T} \exp(2\eta \beta_{1} \gamma (s-T)) \Big(\left(2\eta \delta(s)^{2} + 4\eta \beta_{1}(d-2)\delta(s)^{2} + \eta \delta(s)^{2} + 8\eta^{2}d_{h}d^{2}\delta(s)^{2} \right) \exp(-\eta \beta_{1} \gamma s)$$

$$+ \left(2\eta \beta_{2}\delta(s)^{2} + 8\eta^{2}\beta_{2}d_{h}d\delta(s)^{2}\right) \exp(-\eta \beta_{2} \gamma s) \Big)$$

$$\leq \delta(T) \exp(-\eta \beta_{1} \gamma (T+1))$$

$$+ \left(2\eta \delta(T)^{2} + 4\eta \beta_{1}(d-2)\delta(T)^{2} + \eta \delta(T)^{2} + 8\eta^{2}d_{h}d^{2}\delta(T)^{2}\right) \cdot \frac{2}{\eta \beta_{1} \gamma} \exp(-\eta \beta_{1} \gamma (T+1))$$

$$+ \left(2\eta \beta_{2}\delta(T)^{2} + 8\eta^{2}\beta_{2}d_{h}d\delta(T)^{2}\right) \cdot \frac{3}{\eta \beta_{2} \gamma} \exp(-\eta \beta_{1} \gamma (T+1))$$

$$= \left(1 + \frac{4\delta(T)}{\beta_{1} \gamma} + \frac{8(d-2)\delta(T)}{\gamma} + \frac{2\delta(T)}{\beta_{1} \gamma} + \frac{16\eta d_{h}d^{2}\delta(T)}{\beta_{1} \gamma} + \frac{6\delta(T)}{\gamma} + \frac{24\eta d_{h}d\delta(T)}{\gamma}\right)$$

$$\cdot \delta(T) \exp(-\eta \beta_{1} \gamma (T+1))$$

$$\leq 2\delta(T) \exp(-\eta \beta_{1} \gamma (T+1))$$

$$\leq 2\delta(T+1) \exp(-\eta \beta_{1} \gamma (T+1))$$

$$\leq 2\delta(T+1) \exp(-\eta \beta_{1} \gamma (T+1))$$

$$(36)$$

This bound requires $\frac{\delta(T)}{\gamma} \left(\frac{4}{\beta_1} + 8(d-2) + \frac{2}{\beta_1} + \frac{16\eta d_h d^2}{\beta_1} + 6 + 24\eta d_h d \right) \leq 1$, which can be verified by $d_h = \widetilde{\Omega}(d^2) \geq 36 \log(4d(2d+1)/\delta) \left(\frac{4}{\beta_1} + 8(d-2) + \frac{2}{\beta_1} + \frac{2}{\beta_1} + 6 + \frac{12}{d} \right)^2 \geq 36 \log(4d(2d+1)/\delta) \left(\frac{4}{\beta_1} + 8(d-2) + \frac{2}{\beta_1} + \frac{16\eta d_h d^2}{\beta_1} + 6 + 24\eta d_h d \right)^2$.

Bound of $c_i^{\top}(T+1)b(T+1)$

$$\begin{aligned} & \left| \boldsymbol{c}_{i}^{\top}(T+1)\boldsymbol{b}(T+1) \right| \\ &= \left| \boldsymbol{c}_{i}^{\top}(T)\boldsymbol{b}(T) + \eta \left(\left(\beta_{3} - \beta_{1}\boldsymbol{c}_{i}^{\top}(T)\boldsymbol{b}_{i}(T) \right) \boldsymbol{b}_{i}^{\top}(T)\boldsymbol{b}(T) - \beta_{1} \sum_{k \neq i}^{d} \boldsymbol{c}_{i}^{\top}(T)\boldsymbol{b}_{k}(T) \cdot \boldsymbol{b}_{k}^{\top}(T)\boldsymbol{b}(T) \right. \\ & - \beta_{2}\boldsymbol{c}_{i}^{\top}(T)\boldsymbol{b}(T) \cdot \boldsymbol{b}^{\top}(T)\boldsymbol{b}(T) - \beta_{2} \sum_{k=1}^{d} \boldsymbol{c}_{k}^{\top}(T)\boldsymbol{b}(T) \cdot \boldsymbol{c}_{k}^{\top}(T)\boldsymbol{c}_{i}(T) \right) + \eta^{2}\bar{\boldsymbol{c}}_{i}^{\top}(T)\bar{\boldsymbol{b}}(T) \\ &= \left| \left(1 - \eta\beta_{2} \left(\boldsymbol{b}^{\top}(T)\boldsymbol{b}(T) + \boldsymbol{c}_{i}^{\top}(T)\boldsymbol{c}_{i}(T) \right) \right) \boldsymbol{c}_{i}^{\top}(T)\boldsymbol{b}(T) \right. \\ &+ \eta \left(\beta_{3} - \beta_{1}\boldsymbol{c}_{i}^{\top}(T)\boldsymbol{b}_{i}(T) \right) \boldsymbol{b}_{i}^{\top}(T)\boldsymbol{b}(T) - \eta\beta_{1} \sum_{k \neq i}^{d} \boldsymbol{c}_{i}^{\top}(T)\boldsymbol{b}_{k}(T) \cdot \boldsymbol{b}_{k}^{\top}(T)\boldsymbol{b}(T) \right. \\ &- \eta\beta_{2} \sum_{k \neq i}^{d} \boldsymbol{c}_{k}^{\top}(T)\boldsymbol{b}(T) \cdot \boldsymbol{c}_{k}^{\top}(T)\boldsymbol{c}_{i}(T) + \eta^{2}\bar{\boldsymbol{c}}_{i}^{\top}(T)\bar{\boldsymbol{b}}(T) \right| \\ &= \left| \prod_{s=0}^{T} \left(1 - \eta\beta_{2} \left(\boldsymbol{b}^{\top}(s)\boldsymbol{b}(s) + \boldsymbol{c}_{i}^{\top}(s)\boldsymbol{c}_{i}(s) \right) \right) \boldsymbol{c}_{i}^{\top}(0)\boldsymbol{b}(0) \right. \\ &+ \sum_{s=0}^{T} \prod_{s'=s+1}^{T} \left(1 - \eta\beta_{2} \left(\boldsymbol{b}^{\top}(s')\boldsymbol{b}(s') + \boldsymbol{c}_{i}^{\top}(s')\boldsymbol{c}_{i}(s') \right) \right) \end{aligned}$$

$$\begin{split} &\cdot \left(\underbrace{\eta(\beta_3 - \beta_1 c_i^\top(s)b_i(s))b_i^\top(s)b(s)}_{i} - \eta\beta_1 \sum_{k\neq i}^d c_i^\top(s)b_k(s) \cdot b_k^\top(s)b(s) \right. \\ &- \eta\beta_2 \sum_{k\neq i}^d c_k^\top(s)b(s) \cdot c_k^\top(s)c_i(s) + \underbrace{\eta^2 \overline{c}_i^\top(s)\overline{b}(s)}_{\diamond} \right) \Big| \\ &\leq (1 - 2\eta\beta_2\gamma)^{T+1} \Big| c_i^\top(0)b(0) \Big| \\ &+ \sum_{s=0}^T (1 - 2\eta\beta_2\gamma)^{T-s} \left(\eta\delta(s)^2 \exp(-\eta\beta_1\gamma s) \right. \\ &+ \eta\beta_1(d-1) \cdot 2\delta(s)^2 \exp(-\eta\beta_1\gamma s) \right. \\ &+ \eta\beta_2(d-1) \cdot \left(2\delta(s)^2 \exp(-\eta\beta_1\gamma s) \right. \\ &+ \underbrace{\eta^2 (4d_h a^2\delta(s)^2 \exp(-\eta\beta_1\gamma s) + 28\beta_2^2 d_h d\delta(s)^2 \exp(-\eta\beta_2\gamma s))}_{\diamond} \right) \\ &+ \underbrace{\eta^2 (4d_h a^2\delta(s)^2 \exp(-\eta\beta_1\gamma s) + 28\beta_2^2 d_h d\delta(s)^2 \exp(-\eta\beta_2\gamma s))}_{\diamond} \right) \\ &\leq \exp(-2\eta\beta_2\gamma(T+1)) \Big| c_i^\top(0)b(0) \Big| \\ &+ \sum_{s=0}^T \exp(2\eta\beta_2\gamma(s-T)) \Big((\eta\delta(T)^2 + 2\eta\beta_1(d-1)\delta(T)^2 \\ &+ \eta(d-1)\delta(T)^2 + 4\eta^2 d_h d^2\delta(T)^2 \Big) \exp(-\eta\beta_1\gamma s) \\ &+ (2\eta\beta_2(d-1)\delta(T)^2 + 28\eta^2\beta_2^2 d_h d\delta(s)^2) \exp(-\eta\beta_2\gamma s) \Big) \\ &\leq \delta(T) \exp(-\eta\beta_2\gamma(T+1)) \\ &+ \Big(\eta\delta(T)^2 + 2\eta\beta_1(d-1)\delta(T)^2 + \eta(d-1)\delta(T)^2 + 4\eta^2 d_h d^2\delta(T)^2 \Big) \cdot \frac{2}{\eta\beta_2\gamma} \exp(-\eta\beta_1\gamma(T+1)) \\ &+ \Big(2\eta\beta_2(d-1)\delta(T)^2 + 28\eta^2\beta_2^2 d_h d\delta(T)^2 \Big) \cdot \frac{2}{\eta\beta_2\gamma} \exp(-\eta\beta_2\gamma(T+1)) \\ &+ \frac{\delta(T)}{\beta_2} \Big(\frac{2\delta(T)}{\gamma} + \frac{4\beta_1(d-1)\delta(T)}{\gamma} + \frac{56\eta\beta_2 d_h d\delta(T)}{\gamma} \Big) \exp(-\eta\beta_2\gamma(T+1)) \\ &+ \frac{\delta(T)}{\beta_2} \Big(\frac{2\delta(T)}{\gamma} + \frac{4\beta_1(d-1)\delta(T)}{\gamma} + \frac{\delta(T+1)}{\beta_2} \exp(-\eta\beta_1\gamma(T+1)) \Big) \\ &\leq 2\delta(T) \exp(-\eta\beta_2\gamma(T+1)) + \frac{\delta(T)}{\beta_2} \exp(-\eta\beta_1\gamma(T+1)) \\ &\leq 2\delta(T) \exp(-\eta\beta_2\gamma(T+1)) + \frac{\delta(T+1)}{\beta_2} \exp(-\eta\beta_1\gamma(T+1)) \\ &\leq 2\delta(T+1) \exp(-\eta\beta_2\gamma(T+1) + \frac{\delta(T+1)}{\beta_2} \exp(-\eta\beta_1\gamma(T+1$$

Property $\mathcal{B}(T+1)$ is established by (Eq. (35)), (Eq. (36)) and (Eq. (37)).

C.5.1 Auxiliary lemma

Lemma C.5 As long as $2\eta\beta_1\gamma \leq \ln 2$, and $2\eta\beta_2\gamma \leq \ln 2$, we have

$$\sum_{s=0}^{T} \exp(2\eta\beta_1\gamma(s-T)) \cdot \exp(-\eta\beta_1\gamma s) \le \frac{2}{\eta\beta_1\gamma} \exp(-\eta\beta_1\gamma(T+1))$$

$$\sum_{s=0}^{T} \exp(2\eta\beta_1\gamma(s-T)) \cdot \exp(-\eta\beta_2\gamma s) \le \frac{3}{\eta\beta_2\gamma} \exp(-\eta\beta_1\gamma(T+1))$$

$$\sum_{s=0}^{T} \exp(2\eta\beta_2\gamma(s-T)) \cdot \exp(-\eta\beta_1\gamma s) \le \frac{2}{\eta\beta_2\gamma} \exp(-\eta\beta_1\gamma(T+1))$$

$$\sum_{s=0}^{T} \exp(2\eta\beta_2\gamma(s-T)) \cdot \exp(-\eta\beta_2\gamma s) \le \frac{2}{\eta\beta_2\gamma} \exp(-\eta\beta_2\gamma(T+1))$$

Proof of lemma C.5.

$$\sum_{s=0}^{T} \exp(2\eta\beta_1\gamma(s-T)) \cdot \exp(-\eta\beta_1\gamma s)$$

$$= \sum_{s=0}^{T} \exp(\eta\beta_1\gamma s - 2\eta\beta_1\gamma T)$$

$$\leq \int_{0}^{T+1} \exp(\eta\beta_1\gamma s - 2\eta\beta_1\gamma T) ds$$

$$\leq \frac{1}{\eta\beta_1\gamma} \left(\exp(-\eta\beta_1\gamma(T-1)) - \exp(-2\eta\beta_1\gamma T)\right)$$

$$\leq \frac{1}{\eta\beta_1\gamma} \exp(-\eta\beta_1\gamma(T-1))$$

$$\leq \frac{\exp(2\eta\beta_1\gamma)}{\eta\beta_1\gamma} \exp(-\eta\beta_1\gamma(T+1))$$

$$\leq \frac{2}{\eta\beta_1\gamma} \exp(-\eta\beta_1\gamma(T+1))$$

The first inequality is due to $\exp(\eta \beta_1 \gamma s)$ is monotone increasing. The last inequality is due to $2\eta \beta_1 \gamma \leq \ln 2$.

$$\sum_{s=0}^{T} \exp(2\eta\beta_1\gamma(s-T)) \cdot \exp(-\eta\beta_2\gamma s)$$

$$= \sum_{s=0}^{T} \exp(-\eta(\beta_2 - 2\beta_1)\gamma s - 2\eta\beta_1\gamma T)$$

$$\leq \int_{-1}^{T} \exp(-\eta(\beta_2 - 2\beta_1)\gamma s - 2\eta\beta_1\gamma T) ds$$

$$\leq \frac{1}{\eta(\beta_2 - 2\beta_1)\gamma} \left(\exp(\eta(\beta_2 - 2\beta_1)\gamma - 2\eta\beta_1\gamma T) - \exp(-\eta\beta_2\gamma T) \right)$$

$$\leq \frac{1}{\eta(\beta_2 - 2\beta_1)\gamma} \exp(\eta(\beta_2 - 2\beta_1)\gamma - 2\eta\beta_1\gamma T)$$

$$\leq \frac{\exp(\eta\beta_2\gamma)}{\eta(\beta_2 - 2\beta_1)\gamma} \exp(-2\eta\beta_1\gamma (T+1))$$

$$\leq \frac{2\exp(\eta\beta_2\gamma)}{\eta\beta_2\gamma}\exp(-2\eta\beta_1\gamma(T+1))$$

$$\leq \frac{3}{\eta\beta_2\gamma}\exp(-\eta\beta_1\gamma(T+1))$$

The first inequality is due to $\exp(-\eta(\beta_2-2\beta_1)\gamma s)$ is monotone decreasing. The fourth inequality is due to $\beta_2 \geq 4\beta_1$, thus $\frac{1}{\beta_2-2\beta_1} \leq \frac{2}{\beta_2}$. The last inequality is due to $\eta\beta_2\gamma \leq (\ln 2)/2 \leq \ln(3/2)$.

$$\sum_{s=0}^{T} \exp(2\eta \beta_2 \gamma(s-T)) \cdot \exp(-\eta \beta_1 \gamma s)$$

$$= \sum_{s=0}^{T} \exp(\eta(2\beta_2 - \beta_1)\gamma s - 2\eta \beta_2 \gamma T)$$

$$\leq \int_{0}^{T+1} \exp(\eta(2\beta_2 - \beta_1)\gamma s - 2\eta \beta_2 \gamma T) ds$$

$$\leq \frac{1}{\eta(2\beta_2 - \beta_1)\gamma} \left(\exp(2\eta \beta_2 \gamma - \eta \beta_1 \gamma (T+1)) - \exp(-2\eta \beta_2 \gamma T) \right)$$

$$\leq \frac{1}{\eta(2\beta_2 - \beta_1)\gamma} \exp(2\eta \beta_2 \gamma - \eta \beta_1 \gamma (T+1))$$

$$\leq \frac{\exp(2\eta \beta_2 \gamma)}{\eta(2\beta_2 - \beta_1)\gamma} \exp(-\eta \beta_1 \gamma (T+1))$$

$$\leq \frac{2}{\eta \beta_2 \gamma} \exp(-\eta \beta_1 \gamma (T+1))$$

The first inequality is due to $\exp(\eta(2\beta_2 - \beta_1)\gamma s)$ is monotone increasing. The last inequality is due to $\beta_2 \ge \beta_1$ and $2\eta\beta_2\gamma \le \ln 2$.

The proof of $\sum_{s=0}^T \exp(2\eta\beta_2\gamma(s-T))\cdot \exp(-\eta\beta_2\gamma s) \leq \frac{2}{\eta\beta_2\gamma}\exp(-\eta\beta_2\gamma(T+1))$ is similar to $\sum_{s=0}^T \exp(2\eta\beta_1\gamma(s-T))\cdot \exp(-\eta\beta_1\gamma s) \leq \frac{2}{\eta\beta_1\gamma}\exp(-\eta\beta_1\gamma(T+1))$. Just replace β_1 with β_2 , and consider $2\eta\beta_2\gamma \leq \ln 2$.

C.6 Proof of claim B.3

This Section presents the bounds for terms $\boldsymbol{b}_i^{\top}(T+1)\boldsymbol{b}_j(T+1)$, $\boldsymbol{c}_i^{\top}(T+1)\boldsymbol{c}_j(T+1)$ and $\boldsymbol{b}_i^{\top}(T+1)\boldsymbol{b}(T+1)$ with $i,j\in[1,d],i\neq j$, establishing the property $\mathcal{C}(T+1)$.

Recall the *Vector-coupled Dynamics* equations of $\boldsymbol{b}_i^{\top}(t+1)\boldsymbol{b}_j(t+1)$, $\boldsymbol{c}_i^{\top}(t+1)\boldsymbol{c}_j(t+1)$ and $\boldsymbol{b}_i^{\top}(t+1)\boldsymbol{b}(t+1)$ in lemma A.7:

$$\mathbf{b}_{i}^{\top}(t+1)\mathbf{b}_{j}(t+1) \\
= \mathbf{b}_{i}^{\top}(t)\mathbf{b}_{j}(t) + \eta \Big(2(\beta_{3} - \beta_{1}\mathbf{c}_{i}^{\top}(t)\mathbf{b}_{i}(t))\mathbf{c}_{i}^{\top}(t)\mathbf{b}_{j}(t) + 2(\beta_{3} - \beta_{1}\mathbf{c}_{j}^{\top}(t)\mathbf{b}_{j}(t))\mathbf{c}_{j}^{\top}(t)\mathbf{b}_{i}(t) \\
- \beta_{3}(\mathbf{c}_{i}^{\top}(t)\mathbf{b}_{j}(t) + \mathbf{c}_{j}^{\top}(t)\mathbf{b}_{i}(t)) - 2\beta_{1}\sum_{k\neq i, k\neq j}^{d} \mathbf{c}_{k}^{\top}(t)\mathbf{b}_{i}(t) \cdot \mathbf{c}_{k}^{\top}(t)\mathbf{b}_{j}(t) + \eta^{2}\bar{\mathbf{b}}_{i}^{\top}(t)\bar{\mathbf{b}}_{j}(t)$$
(38)

$$\mathbf{c}_{i}^{\top}(t+1)\mathbf{c}_{j}(t+1) \\
= \mathbf{c}_{i}^{\top}(t)\mathbf{c}_{j}(t) + \eta \Big(2(\beta_{3} - \beta_{1}\mathbf{c}_{i}^{\top}(t)\mathbf{b}_{i}(t))\mathbf{c}_{j}^{\top}(t)\mathbf{b}_{i}(t) + 2(\beta_{3} - \beta_{1}\mathbf{c}_{j}^{\top}(t)\mathbf{b}_{j}(t))\mathbf{c}_{i}^{\top}(t)\mathbf{b}_{j}(t) \\
- \beta_{3}(\mathbf{c}_{i}^{\top}(t)\mathbf{b}_{j}(t) + \mathbf{c}_{j}^{\top}(t)\mathbf{b}_{i}(t)) - 2\beta_{1}\sum_{k\neq i, k\neq j}^{d} \mathbf{c}_{i}^{\top}(t)\mathbf{b}_{k}(t) \cdot \mathbf{c}_{j}^{\top}(t)\mathbf{b}_{k}(t) \\
- 2\beta_{2}\mathbf{c}_{i}^{\top}(t)\mathbf{b}(t) \cdot \mathbf{c}_{j}^{\top}(t)\mathbf{b}(t) + \eta^{2}\bar{\mathbf{c}}_{i}^{\top}(t)\bar{\mathbf{c}}_{j}(t)$$
(39)

$$\mathbf{b}_{i}^{\top}(t+1)\mathbf{b}(t+1) \\
= \mathbf{b}_{i}^{\top}(t)\mathbf{b}(t) + \eta \Big((\beta_{3} - \beta_{1}\mathbf{c}_{i}^{\top}(t)\mathbf{b}_{i}(t))\mathbf{c}_{i}^{\top}(t)\mathbf{b}(t) - \beta_{1} \sum_{k \neq i}^{d} \mathbf{c}_{k}^{\top}(t)\mathbf{b}_{i}(t) \cdot \mathbf{c}_{k}^{\top}(t)\mathbf{b}(t) \\
- \beta_{2} \sum_{k=1}^{d} \mathbf{c}_{k}^{\top}(t)\mathbf{b}(t) \cdot \mathbf{c}_{k}^{\top}(t)\mathbf{b}_{i}(t) \Big) + \eta^{2} \bar{\mathbf{b}}_{i}^{\top}(t)\bar{\mathbf{b}}(t) \tag{40}$$

To give bounds for the above three terms, we will use the following bounds from property $\mathcal{B}(t)$ and lemma C.7:

$$|\beta_3 - \beta_1 \boldsymbol{c}_i^{\top}(t)\boldsymbol{b}_i(t)| \leq \delta(t) \exp(-\eta \beta_1 \gamma t)$$

$$|\boldsymbol{c}_i^{\top}(t)\boldsymbol{b}_j(t)| \le 2\delta(t)\exp(-\eta\beta_1\gamma t)$$

$$|\boldsymbol{c}_i^{\top}(t)\boldsymbol{b}(t)| \le 2\delta(t)\exp(-\eta\beta_2\gamma t) + \frac{\delta(t)}{\beta_2}\exp(-\eta\beta_1\gamma t)$$

$$\left| \bar{\boldsymbol{b}}_i(t)^{\top} \bar{\boldsymbol{b}}_j(t) \right| \le 8d_h d^2 \delta(t)^2 \exp(-\eta \beta_1 \gamma t),$$

$$\left| \bar{\boldsymbol{c}}_i(t)^{\top} \bar{\boldsymbol{c}}_j(t) \right| \le 8d_h d^2 \delta(t)^2 \exp(-\eta \beta_1 \gamma t) + 40\beta_2^2 d_h \delta(t)^2 \exp(-\eta \beta_2 \gamma t),$$

$$\left| \overline{\boldsymbol{b}}_{i}^{\top}(t)\overline{\boldsymbol{b}}(t) \right| \leq 8\beta_{2}d_{h}d^{2}\delta(t)^{2}\exp(-\eta\beta_{2}\gamma t) + 4d_{h}d^{2}\delta(t)^{2}\exp(-\eta\beta_{1}\gamma t)$$

Besides, by $\left|\beta_3 - \beta_1 \boldsymbol{c}_i^\top(t) \boldsymbol{b}_i(t)\right| \leq \delta(t) \exp(-\eta \beta_1 \gamma t)$, we have:

$$\begin{aligned} \left| \boldsymbol{c}_{i}^{\top}(t)\boldsymbol{b}_{i}(t) \right| &\leq \frac{\delta(t) \exp(-\eta \beta_{1} \gamma t) + \beta_{3}}{\beta_{1}} \\ &\leq \frac{4\delta(t) + 2\alpha}{\alpha^{2}} \\ &\leq \frac{5\delta(t)}{\alpha^{2}} \\ &\leq 6\delta(t) \end{aligned}$$

The third inequality is by $\delta(s) \geq 2\sqrt{d_h \log(4d(2d+1)/\delta)} \geq 2\alpha = 2\alpha = 2\exp((-\ln 2)/N)$. For the last inequality, as long as $N \geq \frac{2\ln 2}{\ln 6 - \ln 5}$, we have $\frac{5}{\alpha^2} \leq 6$.

We will provide the upper bounds for $\left| \boldsymbol{b}_i^{\top}(T+1)\boldsymbol{b}_j(T+1) \right|$, $\left| \boldsymbol{c}_i^{\top}(T+1)\boldsymbol{c}_j(T+1) \right|$ and $\left| \boldsymbol{b}_i^{\top}(T+1)\boldsymbol{b}_j(T+1) \right|$ by substituting the above bounds into (Eq. 38), (Eq. 39) and (Eq. 40) respectively.

Bound of $\left| \boldsymbol{b}_i^{\top}(T+1)\boldsymbol{b}_j(T+1) \right|$

$$\begin{vmatrix} \mathbf{b}_{i}^{\top}(T+1)\mathbf{b}_{j}(T+1) \end{vmatrix} = \begin{vmatrix} \mathbf{b}_{i}^{\top}(T)\mathbf{b}_{j}(T) + \eta \left(2\left(\beta_{3} - \beta_{1}\mathbf{c}_{i}^{\top}(T)\mathbf{b}_{i}(T) \right)\mathbf{c}_{i}^{\top}(T)\mathbf{b}_{j}(T) + 2\left(\beta_{3} - \beta_{1}\mathbf{c}_{j}^{\top}(T)\mathbf{b}_{j}(T) \right)\mathbf{c}_{j}^{\top}(T)\mathbf{b}_{i}(T) \right) - 2\beta_{1} \sum_{k \neq i, k \neq j}^{d} \mathbf{c}_{k}^{\top}(T)\mathbf{b}_{i}(T) \cdot \mathbf{c}_{k}^{\top}(T)\mathbf{b}_{j}(T) \right) + \eta^{2} \bar{\mathbf{b}}_{i}^{\top}(T)\bar{\mathbf{b}}_{j}(T) \begin{vmatrix} \mathbf{c}_{i}^{\top}(T)\mathbf{b}_{i}(T) - 2\beta_{1} \sum_{k \neq i, k \neq j}^{d} \mathbf{c}_{k}^{\top}(T)\mathbf{b}_{i}(T) \cdot \mathbf{c}_{k}^{\top}(T)\mathbf{b}_{j}(T) + \eta^{2} \bar{\mathbf{b}}_{i}^{\top}(T)\bar{\mathbf{b}}_{j}(T) \end{vmatrix}$$

$$\leq \left| \mathbf{b}_{i}^{\top}(T)\mathbf{b}_{j}(T) \right| + 2\eta \left| \left(\beta_{3} - \beta_{1}\mathbf{c}_{i}^{\top}(T)\mathbf{b}_{i}(T)\right)\mathbf{c}_{i}^{\top}(T)\mathbf{b}_{j}(T) \right| + 2\eta \left| \left(\beta_{3} - \beta_{1}\mathbf{c}_{j}^{\top}(T)\mathbf{b}_{j}(T)\right)\mathbf{c}_{j}^{\top}(T)\mathbf{b}_{i}(T) \right| + \eta\beta_{3} \left| \left(\mathbf{c}_{i}^{\top}(T)\mathbf{b}_{j}(T) + \mathbf{c}_{j}^{\top}(T)\mathbf{b}_{i}(T)\right) \right| + 2\eta\beta_{1} \sum_{k \neq i, k \neq j}^{d} \left| \mathbf{c}_{k}^{\top}(T)\mathbf{b}_{i}(T) \cdot \mathbf{c}_{k}^{\top}(T)\mathbf{b}_{j}(T) \right| + \eta^{2} \left| \bar{\mathbf{b}}_{i}^{\top}(T)\bar{\mathbf{b}}_{j}(T) \right|$$

$$\leq \delta(T) + 4\eta \cdot \delta(T) \exp(-\eta\beta_{1}\gamma T) \cdot 2\delta(T) \exp(-\eta\beta_{1}\gamma T) + 2\eta\beta_{3} \cdot 2\delta(T) \exp(-\eta\beta_{1}\gamma T) + 2\eta\beta_{1}(d-2) \cdot \left(2\delta(T) \exp(-\eta\beta_{1}\gamma T) + 4\eta\beta_{3} + 8\eta\beta_{1}(d-2)\delta(T) \exp(-\eta\beta_{1}\gamma T) \right)$$

$$\leq \delta(T) + \delta(T) \left(8\eta\delta(T) \exp(-\eta\beta_{1}\gamma T) + 4\eta\beta_{3} + 8\eta\beta_{1}(d-2)\delta(T) \exp(-\eta\beta_{1}\gamma T) + 8\eta^{2}d_{h}d^{2}\delta_{max} \right) \exp(-\eta\beta_{1}\gamma T)$$

$$\leq \delta(T) + \delta(T) \left(8\eta\delta_{max} + 4\eta\beta_{3} + 8\eta\beta_{1}(d-2)\delta_{max} + 8\eta^{2}d_{h}d^{2}\delta_{max} \right) \exp(-\eta\beta_{1}\gamma T)$$

$$\leq \delta(T) + \delta(T) \left(8\eta\delta_{max} + 4\eta\beta_{3} + 8\eta\beta_{1}(d-2)\delta_{max} + 8\eta^{2}d_{h}d^{2}\delta_{max} \right) \exp(-\eta\beta_{1}\gamma T)$$

$$(41)$$

The first inequality is derived by triangle inequality. The second inequality is derived by $\left| \boldsymbol{b}_i^\top (T) \boldsymbol{b}_j (T) \right| \leq \delta(T)$ and substituting the bounds of $|\beta_3 - \beta_1 \boldsymbol{c}_i^\top (t) \boldsymbol{b}_i (t)|$, $|\boldsymbol{c}_i^\top (t) \boldsymbol{b}_j (t)|$ and $\left| \overline{\boldsymbol{b}}_i (t)^\top \overline{\boldsymbol{b}}_j (t) \right|$. The third inequality is derived by factoring out the common factor $\delta(T)$. The last inequality is derived by $\delta(T) \leq \delta_{\max}$ and $\exp(-\eta \beta_1 \gamma T) \leq 1$.

Bound of
$$\left| \boldsymbol{c}_i^{\top} (T+1) \boldsymbol{c}_j (T+1) \right|$$

$$\begin{vmatrix} \mathbf{c}_{i}^{\top}(T+1)\mathbf{c}_{j}(T+1) \end{vmatrix} = \begin{vmatrix} \mathbf{c}_{i}^{\top}(T)\mathbf{c}_{j}(T) + \eta \left(2(\beta_{3} - \beta_{1}\mathbf{c}_{i}^{\top}(T)\mathbf{b}_{i}(T))\mathbf{c}_{j}^{\top}(T)\mathbf{b}_{i}(T) + 2(\beta_{3} - \beta_{1}\mathbf{c}_{j}^{\top}(T)\mathbf{b}_{j}(T))\mathbf{c}_{i}^{\top}(T)\mathbf{b}_{j}(T) \right) - \beta_{3} \left(\mathbf{c}_{i}^{\top}(T)\mathbf{b}_{j}(T) + \mathbf{c}_{j}^{\top}(T)\mathbf{b}_{i}(T) \right) - 2\beta_{1} \sum_{k \neq i, k \neq j}^{d} \mathbf{c}_{i}^{\top}(T)\mathbf{b}_{k}(T) \cdot \mathbf{c}_{j}^{\top}(T)\mathbf{b}_{k}(T) \\ - 2\beta_{2}\mathbf{c}_{i}^{\top}(T)\mathbf{b}(T) \cdot \mathbf{c}_{j}^{\top}(T)\mathbf{b}(T) \right) + \eta^{2}\bar{\mathbf{c}}_{i}^{\top}(T)\bar{\mathbf{c}}_{j}(T) \Big| \\ \leq \left| \mathbf{c}_{i}^{\top}(T)\mathbf{c}_{j}(T) \right| + 2\eta \left| \left(\beta_{3} - \beta_{1}\mathbf{c}_{i}^{\top}(T)\mathbf{b}_{i}(T) \right) \mathbf{c}_{j}^{\top}(T)\mathbf{b}_{i}(T) \right| + 2\eta \left| \left(\beta_{3} - \beta_{1}\mathbf{c}_{j}^{\top}(T)\mathbf{b}_{j}(T) \right) \mathbf{c}_{i}^{\top}(T)\mathbf{b}_{j}(T) \right| \\ + \eta\beta_{3} \left| \left(\mathbf{c}_{i}^{\top}(T)\mathbf{b}_{j}(T) + \mathbf{c}_{j}^{\top}(T)\mathbf{b}_{i}(T) \right) \right| + 2\eta\beta_{1} \sum_{k \neq i, k \neq j}^{d} \mathbf{c}_{i}^{\top}(T)\mathbf{b}_{k}(T) \cdot \mathbf{c}_{j}^{\top}(T)\mathbf{b}_{k}(T) \Big| \\ + 2\eta\beta_{2} \left| \mathbf{c}_{i}^{\top}(T)\mathbf{b}(T) \cdot \mathbf{c}_{j}^{\top}(T)\mathbf{b}(T) \right| + \eta^{2} \left| \mathbf{c}_{i}^{\top}(T)\bar{\mathbf{c}}_{j}(T) \right| \\ \leq \delta(T) + 4\eta \cdot \delta(T) \exp(-\eta\beta_{1}\gamma T) \cdot 2\delta(T) \exp(-\eta\beta_{1}\gamma T) + 2\eta\beta_{3} \cdot 2\delta(T) \exp(-\eta\beta_{1}\gamma T) \\ + 2\eta\beta_{1}(d-2) \cdot \left(2\delta(T) \exp(-\eta\beta_{1}\gamma T) \right) \cdot 2\delta(T) \exp(-\eta\beta_{1}\gamma T) + 2\eta\beta_{3} \cdot 2\delta(T) \exp(-\eta\beta_{1}\gamma T) \\ + 2\eta\beta_{2} \left(2\delta(T) \exp(-\eta\beta_{1}\gamma T) + 40\beta_{2}^{2}d_{h}\delta(T)^{2} \exp(-\eta\beta_{2}\gamma T) \right) \\ \leq \delta(T) + \delta(T) \left(8\eta\delta(T) \exp(-\eta\beta_{1}\gamma T) + 4\eta\beta_{3} + 8\eta\beta_{1}(d-2)\delta(T) \exp(-\eta\beta_{1}\gamma T) + 8\eta^{2}d_{h}d^{2}\delta(T) \right) \\ + \frac{2\eta\delta(T)}{\beta_{2}} \exp(-\eta\beta_{1}\gamma T) + 8\eta\delta(T) \exp(-\eta\beta_{2}\gamma T) \right) \cdot \exp(-\eta\beta_{1}\gamma T) \\ \leq \delta(T) + \delta(T) \left(8\eta\beta_{2}\delta(T) \exp(-\eta\beta_{2}\gamma T) + 40\eta^{2}\beta_{2}^{2}d_{h}\delta(T) \right) \cdot \exp(-\eta\beta_{2}\gamma T) \\ \leq \delta(T) + \delta(T) \left(16\eta\delta_{\max} + 4\eta\beta_{3} + 8\eta\beta_{1}(d-2)\delta_{\max} + 8\eta^{2}d_{h}d^{2}\delta_{\max} + \frac{2\eta\delta_{\max}}{\beta_{2}} \right) \cdot \exp(-\eta\beta_{1}\gamma T) \\ + \delta(T) \left(8\eta\beta_{2}\delta_{\max} + 40\eta^{2}\beta_{2}^{2}d_{h}\delta_{\max} \right) \cdot \exp(-\eta\beta_{2}\gamma T) \end{aligned}$$

The first inequality is derived by triangle inequality. The second inequality is derived by $\left| \boldsymbol{b}_i^\top(T) \boldsymbol{b}_j(T) \right| \leq \delta(T)$ and substituting the bounds of $|\beta_3 - \beta_1 \boldsymbol{c}_i^\top(t) \boldsymbol{b}_i(t)|$, $|\boldsymbol{c}_i^\top(t) \boldsymbol{b}_j(t)|$ and $\left| \overline{\boldsymbol{c}}_i(t)^\top \overline{\boldsymbol{c}}_j(t) \right|$. The third inequality is derived by factoring out the common factor $\delta(T)$. The last inequality is derived by $\delta(T) \leq \delta_{\max}$, $\exp(-\eta \beta_1 \gamma T) \leq 1$ and $\exp(-\eta \beta_2 \gamma T) \leq 1$.

Bound of $\left| \boldsymbol{b}_i^{\top} (T+1) \boldsymbol{b} (T+1) \right|$

$$\begin{aligned} & \left| b_{i}^{\top}(T+1)b(T+1) \right| \\ & = \left| b_{i}^{\top}(T)b(T) + \eta \left(\left(\beta_{3} - \beta_{1}c_{i}^{\top}(T)b_{i}(T) \right)c_{i}^{\top}(T)b(T) - \beta_{1} \sum_{k \neq i}^{d} c_{k}^{\top}(T)b_{i}(T) \cdot c_{k}^{\top}(T)b(T) \right) \\ & - \beta_{2} \sum_{k=1}^{d} c_{k}^{\top}(T)b(T) \cdot c_{k}^{\top}(T)b_{i}(T) \right) + \eta^{2} \overline{b}_{i}^{\top}(T) \overline{b}(T) \right| \\ & = \left| b_{i}^{\top}(T)b(T) + \eta \left(\left(\beta_{3} - \beta_{1}c_{i}^{\top}(T)b_{i}(T) \right)c_{i}^{\top}(T)b(T) - \beta_{1} \sum_{k \neq i}^{d} c_{k}^{\top}(T)b_{i}(T) \cdot c_{k}^{\top}(T)b(T) \right) \\ & - \beta_{2} \sum_{k \neq i}^{d} c_{k}^{\top}(T)b(T) \cdot c_{k}^{\top}(T)b_{i}(T) - \beta_{2}c_{i}^{\top}(T)b(T) \cdot c_{i}^{\top}(T)b_{i}(T) \right) + \eta^{2} \overline{b}_{i}^{\top}(T) \overline{b}(T) \right| \\ & \leq \left| b_{i}^{\top}(T)b(T) \right| + \eta \left| \left(\beta_{3} - \beta_{1}c_{i}^{\top}(T)b_{i}(T) \right)c_{i}^{\top}(T)b(T) \right| + \eta \beta_{1} \sum_{k \neq i}^{d} \left| c_{k}^{\top}(T)b_{i}(T) \cdot c_{k}^{\top}(T)b(T) \right| \\ & + \eta \beta_{2} \sum_{k \neq i}^{d} \left| c_{k}^{\top}(T)b(T) \cdot c_{k}^{\top}(T)b_{i}(T) \right| + \eta \beta_{2} \left| c_{i}^{\top}(T)b(T) \cdot c_{i}^{\top}(T)b_{i}(T) \right| + \eta^{2} \left| \overline{b}_{i}^{\top}(T) \overline{b}(T) \right| \\ & \leq \delta(T) + \eta \cdot \delta(T) \exp(-\eta \beta_{1} \gamma T) \cdot \left(2\delta(T) \exp(-\eta \beta_{2} \gamma T) + \frac{\delta(T)}{\beta_{2}} \exp(-\eta \beta_{1} \gamma T) \right) \\ & + \eta \beta_{2} \cdot 6\delta(T) \cdot \left(2\delta(T) \exp(-\eta \beta_{2} \gamma T) + \frac{\delta(T)}{\beta_{2}} \exp(-\eta \beta_{1} \gamma T) \right) \\ & + \eta \beta_{2} \cdot 6\delta(T) \cdot \left(2\delta(T) \exp(-\eta \beta_{2} \gamma T) + 4d_{n}d^{2}\delta(T)^{2} \exp(-\eta \beta_{1} \gamma T) \right) \\ & \leq \delta(T) + \delta(T) \left(\eta(2(\beta_{1} + \beta_{2})(d-1) + 1 \right) \cdot \frac{\delta(T)}{\beta_{2}} \exp(-\eta \beta_{1} \gamma T) \\ & + \delta(T) \left(2\eta(2(\beta_{1} + \beta_{2})(d-1) + 1 \right) \delta(T) \exp(-\eta \beta_{1} \gamma T) \\ & + \delta(T) \left(2\eta(2(\beta_{1} + \beta_{2})(d-1) + 1 \right) \delta(T) \exp(-\eta \beta_{1} \gamma T) \\ & \leq \delta(T) + \delta(T) \left(4\eta d\delta_{\max} + 4\eta^{2} d_{n}d^{2}\delta_{\max} \right) \cdot \exp(-\eta \beta_{1} \gamma T) \\ & + \delta(T) \left(8\eta \beta_{2} d\delta_{\max} + 12\eta \beta_{2} \delta_{\max} + 8\eta^{2} \beta_{2} d_{n}d^{2}\delta_{\max} \right) \cdot \exp(-\eta \beta_{2} \gamma T) \end{aligned}$$

The first inequality is derived by triangle inequality. The second inequality is derived by $\left| \boldsymbol{b}_i^\top(T) \boldsymbol{b}_j(T) \right| \leq \delta(T)$ and substituting the bounds of $|\beta_3 - \beta_1 \boldsymbol{c}_i^\top(t) \boldsymbol{b}_i(t)|, |\boldsymbol{c}_i^\top(t) \boldsymbol{b}_j(t)|, |\boldsymbol{c}_i^\top(t) \boldsymbol{b}_i(t)|$ and $\left| \bar{\boldsymbol{b}}_i(t)^\top \bar{\boldsymbol{b}}(t) \right|$. The third inequality is derived by factoring out the common factor $\delta(T)$. The last inequality is derived by $\delta(T) \leq \delta_{\max}, \exp(-\eta \beta_1 \gamma T) \leq 1$ and $2(\beta_1 + \beta_2)(d-1) + 1 \leq 4\beta_2 d$ since $\beta_2 \geq \beta_1$ and $\beta_2 \geq 1$.

We next provide the upper bound for $\delta(T+1)$.

$$\delta(T+1) = \max\{|\boldsymbol{b}_{i}^{\top}(T+1)\boldsymbol{b}_{j}(T+1)|, |\boldsymbol{c}_{i}^{\top}(T+1)\boldsymbol{c}_{j}(T+1)|, |\boldsymbol{b}_{i}^{\top}(T+1)\boldsymbol{b}(T+1)|\}
\leq \delta(T) + \delta(T) \left(16\eta d\delta_{\max} + 4\eta\beta_{3} + 8\eta\beta_{1}(d-2)\delta_{\max} + 8\eta^{2}d_{h}d^{2}\delta_{\max} + \frac{2\eta\delta_{\max}}{\beta_{2}}\right) \cdot \exp(-\eta\beta_{1}\gamma T)
+ \delta(T) \left(8\eta\beta_{2}d\delta_{\max} + 40\eta^{2}\beta_{2}^{2}d_{h}\delta_{\max} + 12\eta\beta_{2}\delta_{\max} + 8\eta^{2}\beta_{2}d_{h}d^{2}\delta_{\max}\right) \cdot \exp(-\eta\beta_{2}\gamma T)$$
(44)

This inequality can be verified by comparing with (Eq. (41)), (Eq. (42)), (Eq. (43)). To give more precise bound, we introduce the following lemma:

Lemma C.6 If $y(t+1) \le y(t) + cy(t) \exp(-at) + dy(t) \exp(-bt)$, with $a, b, c, d > 0, t \ge 0$ and $a, b \le \ln 2$, then y(t) satisfies:

$$y(t) \le y(0) \exp(\frac{2c}{a} + \frac{2d}{b})$$

Proof of lemma C.6.

$$y(t+1) \leq y(t) + cy(t) \exp(-at) + dy(t) \exp(-bt)$$

$$\Rightarrow y(t+1) \leq y(t)(1 + c \exp(-at) + d \exp(-bt))$$

$$\Rightarrow y(t+1) \leq y(0) \prod_{s=0}^{t} (1 + c \exp(-as) + d \exp(-bs))$$

$$\Rightarrow \ln y(t+1) \leq \ln y(0) + \sum_{s=0}^{t} \ln(1 + c \exp(-as) + d \exp(-bs))$$

$$\Rightarrow \ln y(t+1) \leq \ln y(0) + \sum_{s=0}^{t} (c \exp(-as) + d \exp(-bs))$$

$$\Rightarrow \ln y(t+1) \leq \ln y(0) + \int_{-1}^{t} (c \exp(-as) + d \exp(-bs)) ds$$

$$\Rightarrow \ln y(t+1) \leq \ln y(0) + (\frac{c}{a}(\exp(a) - \exp(-at)) + \frac{d}{b}(\exp(b) - \exp(-bt)))$$

$$\Rightarrow \ln y(t+1) \leq \ln y(0) + (\frac{2c}{a} + \frac{2d}{b})$$

$$\Rightarrow y(t+1) \leq y(0) \exp(\frac{2c}{a} + \frac{2d}{b})$$

The fourth arrow is due to $\ln(1+x) \le x$ for $x \ge 0$. The fifth arrow is due to $\exp(-as), \exp(-bs)$ are monotone decreasing. the 7-th arrow is due to $a, b \le \ln 2$ and $-\exp(-at) \le 0, -\exp(-bt) \le 0$.

Lemma C.6 presents the core idea of establishing property $\mathcal{C}(T+1)$. If $a\gg c$ and $b\gg d$ in the above lemma, we will have $y(t+1)\leq y(0)\cdot O(1)$. Similarly, as Mamba converges quickly ($\mathbf{C}^{\top}\mathbf{B}\to \frac{\beta_3}{\beta_1}\mathbf{I}$, $\mathbf{C}^{\top}\mathbf{b}\to \mathbf{0}$), we can prove that $|\boldsymbol{b}_i^{\top}(t)\boldsymbol{b}_j(t)|, |\boldsymbol{c}_i^{\top}(t)\boldsymbol{c}_j(t)|, |\boldsymbol{b}_i^{\top}(t)\boldsymbol{b}(t)|$ hold their magnitudes around their initial states.

We next combine (Eq. (44)) and lemma C.6 to give bound for $\delta(T+1)$.

$$\begin{split} &\delta(T+1) \\ &\leq \delta(T) + \delta(T) \Big(16 \eta d\delta_{\max} + 4 \eta \beta_3 + 8 \eta \beta_1 (d-2) \delta_{\max} + 8 \eta^2 d_h d^2 \delta_{\max} + \frac{2 \eta \delta_{\max}}{\beta_2} \Big) \cdot \exp(-\eta \beta_1 \gamma T) \\ &+ \delta(T) \Big(8 \eta \beta_2 d\delta_{\max} + 40 \eta^2 \beta_2^2 d_h \delta_{\max} + 12 \eta \beta_2 \delta_{\max} + 8 \eta^2 \beta_2 d_h d^2 \delta_{\max} \Big) \cdot \exp(-\eta \beta_2 \gamma T) \\ &\leq \delta(0) \cdot \exp\Big(\frac{2 \Big(16 \eta d\delta_{\max} + 4 \eta \beta_3 + 8 \eta \beta_1 (d-2) \delta_{\max} + 8 \eta^2 d_h d^2 \delta_{\max} + \frac{2 \eta \delta_{\max}}{\beta_2} \Big)}{\eta \beta_1 \gamma} \end{split}$$

$$\begin{split} & + \frac{2 \left(8 \eta \beta_2 d \delta_{\max} + 40 \eta^2 \beta_2^2 d_h \delta_{\max} + 12 \eta \beta_2 \delta_{\max} + 8 \eta^2 \beta_2 d_h d^2 \delta_{\max} \right)}{\eta \beta_2 \gamma} \\ & \leq \delta(0) \cdot \exp \left(\frac{3 \sqrt{d_h \log(4d(2d+1)/\delta)}}{\frac{1}{2} d_h} \cdot \left(\frac{32d}{\beta_1} + \frac{8\beta_3}{\beta_1} + 16(d-2) + \frac{16 \eta d_h d^2}{\beta_1} + \frac{4}{\beta_1 \beta_2} \right) \\ & + 16d + 80 \eta \beta_2 d_h + 24 + 16 \eta d_h d^2 \right) \right) \\ & \leq \frac{3}{2} \cdot \delta(0) \\ & \leq 3 \sqrt{d_h \log(4d(2d+1)/\delta)} \end{split}$$

The first inequality is derived by (Eq. (44)). The second inequality is derived by lemma C.6. The third inequality is derived by $\gamma = \frac{1}{2}d_h$. The last inequality is derived by

$$\begin{split} &d_h = \widetilde{\Omega}(d^2) \\ &\geq \frac{36}{(\ln(3/2))^2} \log(4d(2d+1)/\delta) \Big(\frac{32d}{\beta_1} + \frac{8\beta_3}{\beta_1} \\ &+ 16(d-2) + \frac{8}{\beta_1} + \frac{4}{\beta_1\beta_2} + 16d + 80 \ln 2 + 24 + 8 \Big)^2 \\ &\geq \frac{36}{(\ln(3/2))^2} \log(4d(2d+1)/\delta) \Big(\frac{32d}{\beta_1} + \frac{8\beta_3}{\beta_1} \\ &+ 16(d-2) + \frac{16\eta d_h d^2}{\beta_1} + \frac{4}{\beta_1\beta_2} + 16d + 80\eta\beta_2 d_h + 24 + 16\eta d_h d^2 \Big)^2 \end{split}$$

C.7 Bounds of η^2 terms

This Section presents the bounds for $\bar{\boldsymbol{b}}_i(t)^{\top}\bar{\boldsymbol{b}}_j(t)$, $\bar{\boldsymbol{c}}_i(t)^{\top}\bar{\boldsymbol{c}}_j(t)$, $\left\|\bar{\boldsymbol{b}}(t)\right\|_2^2$, $\bar{\boldsymbol{c}}_i^{\top}(t)\bar{\boldsymbol{b}}_j(t)$, $\bar{\boldsymbol{c}}_i^{\top}(t)\bar{\boldsymbol{b}}(t)$, $\bar{\boldsymbol{b}}_i^{\top}(t)\bar{\boldsymbol{b}}(t)$ (these terms usually appear in the *Vector-coupled Dynamics* equations with a η^2 factor) with $i,j\in[1,d]$ under the assumption that $\mathcal{A}(t)$, $\mathcal{B}(t)$, and $\mathcal{C}(t)$ hold.

Lemma C.7 *Under the assumption that* A(t), B(t), and C(t) hold, we have the following bounds:

$$\left| \bar{\boldsymbol{b}}_{i}(t)^{\top} \bar{\boldsymbol{b}}_{j}(t) \right| \leq 8d_{h}d^{2}\delta(t)^{2} \exp(-\eta\beta_{1}\gamma t),$$

$$\left| \bar{\boldsymbol{c}}_{i}(t)^{\top} \bar{\boldsymbol{c}}_{j}(t) \right| \leq 8d_{h}d^{2}\delta(t)^{2} \exp(-\eta\beta_{1}\gamma t) + 24\beta_{2}^{2}d_{h}\delta(t)^{2} \exp(-\eta\beta_{2}\gamma t),$$

$$\left\| \bar{\boldsymbol{b}}(t) \right\|_{2}^{2} \leq 16d_{h}\beta_{2}^{2}d^{2}\delta(t)^{2} \exp(-\eta\beta_{2}\gamma t) + 2d_{h}d^{2}\delta(t)^{2} \exp(-\eta\beta_{1}\gamma t),$$

$$\left| \bar{\boldsymbol{c}}_{i}^{\top}(t)\bar{\boldsymbol{b}}_{j}(t) \right| \leq 8d_{h}d^{2}\delta(t)^{2} \exp(-\eta\beta_{1}\gamma t) + 8\beta_{2}d_{h}d\delta(t)^{2} \exp(-\eta\beta_{2}\gamma t),$$

$$\left| \bar{\boldsymbol{c}}_{i}^{\top}(t)\bar{\boldsymbol{b}}(t) \right| \leq 4d_{h}d^{2}\delta(t)^{2} \exp(-\eta\beta_{1}\gamma t) + 28\beta_{2}^{2}d_{h}d\delta(t)^{2} \exp(-2\eta\beta_{2}\gamma t),$$

$$\left| \bar{\boldsymbol{b}}_{i}^{\top}(t)\bar{\boldsymbol{b}}(t) \right| \leq 8\beta_{2}d_{h}d^{2}\delta(t)^{2} \exp(-\eta\beta_{2}\gamma t) + 4d_{h}d^{2}\delta(t)^{2} \exp(-\eta\beta_{1}\gamma t)$$

where $i, j \in [1, d]$. Note that this lemma does not require $i \neq j$.

Firstly, recall the following dynamics equation in lemma A.6:

$$\mathbf{b}_{i}(t+1) = \mathbf{b}_{i}(t) + \eta \Big((\beta_{3} - \beta_{1} \mathbf{c}_{i}^{\top}(t) \mathbf{b}_{i}(t)) \mathbf{c}_{i}(t) - \beta_{1} \sum_{k \neq i}^{d} \mathbf{c}_{k}^{\top}(t) \mathbf{b}_{i}(t) \cdot \mathbf{c}_{k}(t) \Big)$$

$$=: \mathbf{b}_{i}(t) + \eta \bar{\mathbf{b}}_{i}(t)$$

$$\boldsymbol{c}_i(t+1) = \boldsymbol{c}_i(t) + \eta \Big(\big(\beta_3 - \beta_1 \boldsymbol{c}_i^\top(t) \boldsymbol{b}_i(t) \big) \boldsymbol{b}_i(t) - \beta_1 \sum_{k \neq i}^d \boldsymbol{c}_i^\top(t) \boldsymbol{b}_k(t) \cdot \boldsymbol{b}_k(t) \Big)$$

$$-\beta_2 \boldsymbol{c}_i^{\top}(t) \boldsymbol{b}(t) \cdot \boldsymbol{b}(t)$$

$$=: \boldsymbol{c}_i(t) + \eta \bar{\boldsymbol{c}}_i(t)$$

$$\boldsymbol{b}(t+1) = \boldsymbol{b}(t) - \eta \left(\beta_2 \sum_{i=1}^{d} \boldsymbol{c}_k^{\top}(t) \boldsymbol{b}(t) \cdot \boldsymbol{c}_k(t) \right) =: \boldsymbol{b}(t) + \eta \bar{\boldsymbol{b}}(t)$$

Thus we have

$$\bar{\boldsymbol{b}}_{i}(t) = (\beta_{3} - \beta_{1}\boldsymbol{c}_{i}^{\top}(t)\boldsymbol{b}_{i}(t))\boldsymbol{c}_{i}(t) - \beta_{1}\sum_{k\neq i}^{d}\boldsymbol{c}_{k}^{\top}(t)\boldsymbol{b}_{i}(t) \cdot \boldsymbol{c}_{k}(t)$$

$$\bar{\boldsymbol{c}}_{i}(t) = (\beta_{3} - \beta_{1}\boldsymbol{c}_{i}^{\top}(t)\boldsymbol{b}_{i}(t))\boldsymbol{b}_{i}(t) - \beta_{1}\sum_{k\neq i}^{d}\boldsymbol{c}_{i}^{\top}(t)\boldsymbol{b}_{k}(t) \cdot \boldsymbol{b}_{k}(t) - \beta_{2}\boldsymbol{c}_{i}^{\top}(t)\boldsymbol{b}(t) \cdot \boldsymbol{b}(t)$$

$$\bar{\boldsymbol{b}}(t) = \beta_{2}\sum_{k=1}^{d}\boldsymbol{c}_{k}^{\top}(t)\boldsymbol{b}(t) \cdot \boldsymbol{c}_{k}(t)$$

Recalling the properties A(t) and B(t):

$$\begin{aligned} \mathcal{A}(t): \\ d_h/2 &\leq \boldsymbol{b}_i^\top(t)\boldsymbol{b}_i(t), \boldsymbol{c}_i^\top(t)\boldsymbol{c}_i(t), \boldsymbol{b}^\top(t)\boldsymbol{b}(t) \leq 2d_h \\ \mathcal{B}(t): \\ |\beta_3/\beta_1 - \boldsymbol{c}_i^\top(t)\boldsymbol{b}_i(t)| &\leq \delta(t)\exp(-\eta\beta_1\gamma t) \\ |\boldsymbol{c}_i^\top(t)\boldsymbol{b}_j(t)| &\leq 2\delta(t)\exp(-\eta\beta_1\gamma t) \\ |\boldsymbol{c}_i^\top(t)\boldsymbol{b}(t)| &\leq 2\delta(t)\exp(-\eta\beta_2\gamma t) + \frac{\delta(t)}{\beta_2}\exp(-\eta\beta_1\gamma t) \end{aligned}$$

We can derive the follow bounds for the norm of $\bar{b}_i(t)$, $\bar{c}_i(t)$ and $\bar{b}(t)$:

$$\begin{aligned} & \left\| \bar{\boldsymbol{b}}_{i}(t) \right\| = \left\| \left(\beta_{3} - \beta_{1} \boldsymbol{c}_{i}^{\top}(t) \boldsymbol{b}_{i}(t) \right) \boldsymbol{c}_{i}(t) - \beta_{1} \sum_{k \neq i}^{d} \boldsymbol{c}_{k}^{\top}(t) \boldsymbol{b}_{i}(t) \cdot \boldsymbol{c}_{k}(t) \right\| \\ & \leq \left\| \boldsymbol{c}_{i}(t) \right\| \cdot \left| \left(\beta_{3} - \beta_{1} \boldsymbol{c}_{i}^{\top}(t) \boldsymbol{b}_{i}(t) \right) \right| + \beta_{1} \sum_{k \neq i}^{d} \left\| \boldsymbol{c}_{k}(t) \right\| \cdot \left| \boldsymbol{c}_{k}^{\top}(t) \boldsymbol{b}_{i}(t) \right| \\ & \leq \sqrt{2d_{h}} \delta(t) \exp(-\eta \beta_{1} \gamma t) + \sqrt{2d_{h}} \beta_{1}(d-1) \cdot 2\delta(t) \exp(-\eta \beta_{1} \gamma t) \\ & \leq 2\sqrt{2d_{h}} d\delta(t) \exp(-\eta \beta_{1} \gamma t) \end{aligned}$$

The last inequality is by $\beta_1 \leq 1$.

$$\begin{aligned} & \left\| \bar{\boldsymbol{c}}_{i}(t) \right\| = \left\| \left(\beta_{3} - \beta_{1} \boldsymbol{c}_{i}^{\top}(t) \boldsymbol{b}_{i}(t) \right) \boldsymbol{b}_{i}(t) - \beta_{1} \sum_{k \neq i}^{d} \boldsymbol{c}_{i}^{\top}(t) \boldsymbol{b}_{k}(t) \cdot \boldsymbol{b}_{k}(t) - \beta_{2} \boldsymbol{c}_{i}^{\top}(t) \boldsymbol{b}(t) \cdot \boldsymbol{b}(t) \right\| \\ & \leq \left\| \boldsymbol{b}_{i}(t) \right\| \cdot \left| \left(\beta_{3} - \beta_{1} \boldsymbol{c}_{i}^{\top}(t) \boldsymbol{b}_{i}(t) \right) \right| + \beta_{1} \sum_{k \neq i}^{d} \left\| \boldsymbol{b}_{k}(t) \right\| \cdot \left| \boldsymbol{c}_{i}^{\top}(t) \boldsymbol{b}_{k}(t) \right| + \beta_{2} \left\| \boldsymbol{b}(t) \right\| \cdot \left| \boldsymbol{c}_{i}^{\top}(t) \boldsymbol{b}(t) \right| \\ & \leq \sqrt{2d_{h}} \delta(t) \exp(-\eta \beta_{1} \gamma t) + \sqrt{2d_{h}} \beta_{1}(d-1) \cdot 2\delta(t) \exp(-\eta \beta_{1} \gamma t) \\ & + \sqrt{2d_{h}} \beta_{2} \cdot \left(2\delta(t) \exp(-\eta \beta_{2} \gamma t) + \frac{\delta(t)}{\beta_{2}} \exp(-\eta \beta_{1} \gamma t) \right) \\ & \leq 2\sqrt{2d_{h}} d\delta(t) \exp(-\eta \beta_{1} \gamma t) + 2\sqrt{2d_{h}} \beta_{2} \delta(t) \exp(-\eta \beta_{2} \gamma t) \end{aligned}$$

The last inequality is by $\beta_1 \leq 1$.

$$\left\| \bar{\boldsymbol{b}}(t) \right\| = \left\| \beta_2 \sum_{k=1}^{d} \boldsymbol{c}_k^{\top}(t) \boldsymbol{b}(t) \cdot \boldsymbol{c}_k(t) \right\|$$

$$\leq \beta_2 \sum_{k=1}^{d} \left\| \boldsymbol{c}_k(t) \right\| \cdot \left| \boldsymbol{c}_k^{\top}(t) \boldsymbol{b}(t) \right|$$

$$\leq \sqrt{2d_h} \beta_2 d \cdot \left(2\delta(t) \exp(-\eta \beta_2 \gamma t) + \frac{\delta(t)}{\beta_2} \exp(-\eta \beta_1 \gamma t) \right)$$

$$\leq 2\sqrt{2d_h} \beta_2 d\delta(t) \exp(-\eta \beta_2 \gamma t) + \sqrt{2d_h} d\delta(t) \exp(-\eta \beta_1 \gamma t)$$

By multiplying them pairwise, we obtain:

$$\left| \bar{\boldsymbol{b}}_i(t)^\top \bar{\boldsymbol{b}}_j(t) \right| \leq \left(2\sqrt{2d_h} d\delta(t) \exp(-\eta \beta_1 \gamma t) \right)^2 \leq 8d_h d^2 \delta(t)^2 \exp(-\eta \beta_1 \gamma t)$$

The last inequality is by $\exp(-2\eta\beta_1\gamma t) \leq \exp(-\eta\beta_1\gamma t)$.

$$\begin{aligned} &\left| \bar{\boldsymbol{c}}_{i}(t)^{\top} \bar{\boldsymbol{c}}_{j}(t) \right| \leq \left(2\sqrt{2d_{h}}d\delta(t) \exp(-\eta\beta_{1}\gamma t) + 2\sqrt{2d_{h}}\beta_{2}\delta(t) \exp(-\eta\beta_{2}\gamma t) \right)^{2} \\ &= 8d_{h}d^{2}\delta(t)^{2} \exp(-2\eta\beta_{1}\gamma t) + 8\beta_{2}^{2}d_{h}\delta(t)^{2} \exp(-2\eta\beta_{2}\gamma t) \\ &+ 16d_{h}\beta_{2}d\delta(t)^{2} \exp(-\eta(\beta_{1}+\beta_{2})\gamma t) \\ &\leq 8d_{h}d^{2}\delta(t)^{2} \exp(-\eta\beta_{1}\gamma t) + 40\beta_{2}^{2}d_{h}\delta(t)^{2} \exp(-\eta\beta_{2}\gamma t) \end{aligned}$$

The last inequality is by $\beta_2 d \leq 2\beta_2^2$ because $\beta_2 = \Omega(d^2) \geq d$, and $\exp(-2\eta\beta_1\gamma t) \leq \exp(-\eta\beta_1\gamma t)$, $\exp(-2\eta\beta_2\gamma t) \leq \exp(-\eta\beta_2\gamma t)$, $\exp(-\eta(\beta_1+\beta_2)\gamma t) \leq \exp(-\eta\beta_2\gamma t)$.

$$\begin{split} & \left\| \bar{\boldsymbol{b}}(t) \right\|_{2}^{2} \leq \left(2\sqrt{2d_{h}}\beta_{2}d\delta(t) \exp(-\eta\beta_{2}\gamma t) + \sqrt{2d_{h}}d\delta(t) \exp(-\eta\beta_{1}\gamma t) \right)^{2} \\ &= 8d_{h}\beta_{2}^{2}d^{2}\delta(t)^{2} \exp(-2\eta\beta_{2}\gamma t) + 2d_{h}d^{2}\delta(t)^{2} \exp(-2\eta\beta_{1}\gamma t) \\ &+ 8d_{h}\beta_{2}d^{2}\delta(t)^{2} \exp(-\eta(\beta_{1} + \beta_{2})\gamma t) \\ &\leq 16d_{h}\beta_{2}^{2}d^{2}\delta(t)^{2} \exp(-\eta\beta_{2}\gamma t) + 2d_{h}d^{2}\delta(t)^{2} \exp(-\eta\beta_{1}\gamma t) \end{split}$$

The last inequality is by $\beta_2 \leq \beta_2^2$ because $\beta_2 \geq 1$, and $\exp(-2\eta\beta_1\gamma t) \leq \exp(-\eta\beta_1\gamma t)$, $\exp(-\eta\beta_2\gamma t) \leq \exp(-\eta\beta_2\gamma t)$, $\exp(-\eta(\beta_1+\beta_2)\gamma t) \leq \exp(-\eta\beta_2\gamma t)$.

$$\begin{aligned} & \left| \bar{\boldsymbol{c}}_{i}^{\top}(t)\bar{\boldsymbol{b}}_{j}(t) \right| \leq \left\| \bar{\boldsymbol{c}}_{i}(t) \right\| \cdot \left\| \bar{\boldsymbol{b}}_{j}(t) \right\| \\ & \leq \left(2\sqrt{2d_{h}}d\delta(t)\exp(-\eta\beta_{1}\gamma t) + 2\sqrt{2d_{h}}\beta_{2}\delta(t)\exp(-\eta\beta_{2}\gamma t) \right) \cdot 2\sqrt{2d_{h}}d\delta(t)\exp(-\eta\beta_{1}\gamma t) \\ & \leq 8d_{h}d^{2}\delta(t)^{2}\exp(-\eta\beta_{1}\gamma t) + 8\beta_{2}d_{h}d\delta(t)^{2}\exp(-\eta\beta_{2}\gamma t) \end{aligned}$$

The last inequality is by $\exp(-2\eta\beta_1\gamma t) \le \exp(-\eta\beta_1\gamma t)$, $\exp(-\eta(\beta_1+\beta_2)\gamma t) \le \exp(-\eta\beta_2\gamma t)$.

$$\begin{split} & \left| \bar{\boldsymbol{c}}_{i}^{\top}(t)\bar{\boldsymbol{b}}(t) \right| \leq \left\| \bar{\boldsymbol{c}}_{i}(t) \right\| \cdot \left\| \bar{\boldsymbol{b}}(t) \right\| \\ & \leq \left(2\sqrt{2d_{h}}d\delta(t)\exp(-\eta\beta_{1}\gamma t) + 2\sqrt{2d_{h}}\beta_{2}\delta(t)\exp(-\eta\beta_{2}\gamma t) \right) \\ & \cdot \left(2\sqrt{2d_{h}}\beta_{2}d\delta(t)\exp(-\eta\beta_{2}\gamma t) + \sqrt{2d_{h}}d\delta(t)\exp(-\eta\beta_{1}\gamma t) \right) \\ & = 4d_{h}d^{2}\delta(t)^{2}\exp(-2\eta\beta_{1}\gamma t) + 8\beta_{2}^{2}d_{h}d\delta(t)^{2}\exp(-2\eta\beta_{2}\gamma t) \\ & + 8\beta_{2}d_{h}d^{2}\delta(t)^{2}\exp(-\eta(\beta_{1}+\beta_{2})\gamma t) + 4\beta_{2}d_{h}d\delta(t)^{2}\exp(-\eta(\beta_{1}+\beta_{2})\gamma t) \\ & \leq 4d_{h}d^{2}\delta(t)^{2}\exp(-\eta\beta_{1}\gamma t) + 28\beta_{2}^{2}d_{h}d\delta(t)^{2}\exp(-2\eta\beta_{2}\gamma t) \end{split}$$

The last inequality is by $\beta_2 d \leq 2\beta_2^2$, $\beta_2 \leq \beta_2^2$, and $\exp(-2\eta\beta_1\gamma t) \leq \exp(-\eta\beta_1\gamma t)$, $\exp(-\eta\beta_2\gamma t) \leq \exp(-\eta\beta_2\gamma t)$, $\exp(-\eta(\beta_1+\beta_2)\gamma t) \leq \exp(-\eta\beta_2\gamma t)$.

$$\begin{split} \left| \overline{\boldsymbol{b}}_{i}^{\top}(t) \overline{\boldsymbol{b}}(t) \right| &\leq \left\| \overline{\boldsymbol{b}}_{i}(t) \right\| \cdot \left\| \overline{\boldsymbol{b}}(t) \right\| \\ &\leq 2 \sqrt{2 d_{h}} d\delta(t) \exp(-\eta \beta_{1} \gamma t) \cdot \left(2 \sqrt{2 d_{h}} \beta_{2} d\delta(t) \exp(-\eta \beta_{2} \gamma t) + \sqrt{2 d_{h}} d\delta(t) \exp(-\eta \beta_{1} \gamma t) \right) \\ &\leq 8 \beta_{2} d_{h} d^{2} \delta(t)^{2} \exp(-\eta \beta_{2} \gamma t) + 4 d_{h} d^{2} \delta(t)^{2} \exp(-\eta \beta_{1} \gamma t) \\ \text{The last inequality is by } \exp(-2 \eta \beta_{1} \gamma t) \leq \exp(-\eta \beta_{1} \gamma t), \quad \exp(-\eta (\beta_{1} + \beta_{2}) \gamma t) \leq \exp(-\eta \beta_{2} \gamma t). \end{split}$$

D Discussion

In this section, we show that orthogonal initialization Mamba can be trained to ICL solution, and compare our method with some previous works.

Orthogonal Initialization Now we assume that each column of W_B and W_C are initialized with orthogonal columns of unit norm. Then we have

$$C^{\top}(0)C(0) = B^{\top}(0)B(0) = I, \quad B^{\top}(0)b(0) = C^{\top}(0)b(0) = 0.$$

Consider the following update rule as part of lemma 5.1.

$$\boldsymbol{B}(t+1) = \boldsymbol{B}(t) + \eta \beta_3 \boldsymbol{C}(t) - \eta \beta_1 \boldsymbol{C}(t) \boldsymbol{C}(t)^{\mathsf{T}} \boldsymbol{B}(t), \tag{45}$$

$$\boldsymbol{C}(t+1) = \boldsymbol{C}(t) + \eta \beta_3 \boldsymbol{B}(t) - \eta \beta_1 \boldsymbol{B}(t) \boldsymbol{B}(t)^{\mathsf{T}} \boldsymbol{C}(t) - \eta \beta_2 \boldsymbol{b}(t) \boldsymbol{b}(t)^{\mathsf{T}} \boldsymbol{C}(t), \tag{46}$$

$$\boldsymbol{b}(t+1) = \boldsymbol{b}(t) - \eta \beta_2 \boldsymbol{C}(t) \boldsymbol{C}(t)^{\mathsf{T}} \boldsymbol{b}(t). \tag{47}$$

By (Eq. (45)), (Eq. (46)) and (Eq. (47)), we have:

$$\boldsymbol{B}^{\top}(t+1)\boldsymbol{b}(t+1) = \boldsymbol{B}(t)^{\top}\boldsymbol{b}(t) + \eta\beta_{3}\boldsymbol{C}(t)^{\top}\boldsymbol{b}(t) - \eta\beta_{1}\boldsymbol{B}(t)^{\top}\boldsymbol{C}(t)\boldsymbol{C}(t)^{\top}\boldsymbol{b}(t) - \eta\beta_{2}\boldsymbol{B}(t)^{\top}\boldsymbol{C}(t)\boldsymbol{C}(t)^{\top}\boldsymbol{b}(t) - \eta^{2}\beta_{2}\beta_{3}\boldsymbol{C}(t)^{\top}\boldsymbol{C}(t)\boldsymbol{C}(t)^{\top}\boldsymbol{b}(t) + \eta^{2}\beta_{1}\beta_{2}\boldsymbol{B}(t)^{\top}\boldsymbol{C}(t)\boldsymbol{C}(t)^{\top}\boldsymbol{C}(t)\boldsymbol{C}(t)^{\top}\boldsymbol{b}(t)$$

$$C(t+1)^{\top} \boldsymbol{b}(t+1) = \boldsymbol{C}(t)^{\top} \boldsymbol{b}(t) + \eta \beta_3 \boldsymbol{B}(t)^{\top} \boldsymbol{b}(t) - \eta \beta_1 \boldsymbol{C}(t)^{\top} \boldsymbol{B}(t) \boldsymbol{B}(t)^{\top} \boldsymbol{b}(t) - \eta \beta_2 \boldsymbol{C}(t)^{\top} \boldsymbol{b}(t) \boldsymbol{b}(t)^{\top} \boldsymbol{b}(t)$$
$$- \eta \beta_2 \boldsymbol{C}(t)^{\top} \boldsymbol{C}(t) \boldsymbol{C}(t)^{\top} \boldsymbol{b}(t) - \eta^2 \beta_2 \beta_3 \boldsymbol{B}(t)^{\top} \boldsymbol{C}(t) \boldsymbol{C}(t)^{\top} \boldsymbol{b}(t)$$
$$+ \eta^2 \beta_1 \beta_2 \boldsymbol{C}(t)^{\top} \boldsymbol{B}(t) \boldsymbol{B}(t)^{\top} \boldsymbol{C}(t) \boldsymbol{C}(t)^{\top} \boldsymbol{b}(t) + \eta^2 \beta_2^2 \boldsymbol{C}(t)^{\top} \boldsymbol{b}(t) \boldsymbol{b}(t)^{\top} \boldsymbol{C}(t) \boldsymbol{C}(t)^{\top} \boldsymbol{b}(t)$$

Combining $B^{\top}(0)b(0) = C^{\top}(0)b(0) = 0$ with induction, we can derive that $B^{\top}(t)b(t) = C^{\top}(t)b(t) = 0$ for $t \ge 0$. Thus we only need to consider the following dynamics.

$$\boldsymbol{B}(t+1) = \boldsymbol{B}(t) + \eta \beta_3 \boldsymbol{C}(t) - \eta \beta_1 \boldsymbol{C}(t) \boldsymbol{C}(t)^{\top} \boldsymbol{B}(t), \tag{48}$$

$$C(t+1) = C(t) + \eta \beta_3 B(t) - \eta \beta_1 B(t) B(t)^{\mathsf{T}} C(t)$$
(49)

Denote $\boldsymbol{B}(t)^{\top}\boldsymbol{B}(t) = D(t)$, $\boldsymbol{C}(t)^{\top}\boldsymbol{C}(t) = E(t)$ and $\boldsymbol{C}(t)^{\top}\boldsymbol{B}(t) = F(t)$ then by (Eq. (48)) and (Eq. (49)), we have

$$F(t+1) = F(t) + \eta \beta_3 (\mathbf{D}(t) + \mathbf{E}(t)) - \eta \beta_1 F(t) \mathbf{D}(t)$$

$$+ \eta^2 \beta_3^2 F(t)^{\mathsf{T}} - 2\eta^2 \beta_1 \beta_3 F(t) F(t)^{\mathsf{T}} - \eta \beta_1 E(t) F(t)$$

$$+ \eta^2 \beta_1^2 F(t) F(t)^{\mathsf{T}} F(t)$$
(50)

$$D(t+1) = D(t) + \eta \beta_3 (\mathbf{F}(t) + \mathbf{F}(t)^{\top}) - \eta \beta_1 (\mathbf{F}(t)^{\top} \mathbf{F}(t) + \mathbf{F}(t)^{\top} \mathbf{F}(t))$$

$$+ \eta^2 \beta_3^2 \mathbf{E}(t) - \eta^2 \beta_1 \beta_3 \mathbf{F}(t)^{\top} \mathbf{E}(t)$$

$$- \eta^2 \beta_1 \beta_3 \mathbf{E}(t) \mathbf{F}(t) + \eta^2 \beta_1^2 \mathbf{F}(t)^{\top} \mathbf{E}(t) \mathbf{F}(t)$$
(51)

$$E(t+1) = E(t) + \eta \beta_3 (\mathbf{F}(t) + \mathbf{F}(t)^{\top}) - \eta \beta_1 (\mathbf{F}(t)^{\top} \mathbf{F}(t) + \mathbf{F}(t)^{\top} \mathbf{F}(t))$$

$$+ \eta^2 \beta_3^2 \mathbf{D}(t) - \eta^2 \beta_1 \beta_3 \mathbf{F}(t)^{\top} \mathbf{D}(t)$$

$$- \eta^2 \beta_1 \beta_3 \mathbf{D}(t) \mathbf{F}(t) + \eta^2 \beta_1^2 \mathbf{F}(t)^{\top} \mathbf{D}(t) \mathbf{F}(t)$$
(52)

Note that $D(0)=E(0)={\bf I}$ and $F(0)={\bf 0}$. By induction we can see that D(t), E(t) and F(t) are diagonal matrix for t>0. Because of the symmetry, we have D(t)=E(t). Now we denote $D(t)=E(t)=g(t){\bf I}$ and $F(t)=h(t){\bf I}$. Then based on (Eq. (50)), (Eq. (51)) and (Eq. (52)), we have:

$$q(t+1) = q(t) + \eta(2h(t) + \eta\beta_3 q(t) - \eta\beta_1 q(t)h(t))(\beta_3 - \beta_1 h(t))$$
(53)

$$h(t+1) = h(t) + \eta g(t)(\beta_3 - \beta_1 h(t)) + \eta^2 h(t)(\beta_3 - \beta_1 h(t))^2$$
(54)

Since g(0)=1 and h(0)=0 at initialization, h(t) will converge to $\frac{\beta_3}{\beta_1}$ (i.e. $C^{\top}B\to \frac{\beta_3}{\beta_1}I$).

Compare with Other Techniques (Eq. (45)), (Eq. (46)) and (Eq. (47)) can be viewed as the gradient descent that minimize the following target:

$$\frac{1}{2} \| \boldsymbol{C}^{\top} \boldsymbol{W}_{B} \boldsymbol{X} - \boldsymbol{Y} \|_{F}^{2} \tag{55}$$

where X and Y satisfy:

$$\boldsymbol{X}\boldsymbol{X}^{\top} = \begin{bmatrix} \beta_1 & & & \\ & \ddots & & \\ & & \beta_1 & \\ & & & \beta_2 \end{bmatrix} \in \mathbb{R}^{(d+1)\times(d+1)}, \boldsymbol{X}\boldsymbol{Y}^{\top} = \begin{bmatrix} \beta_3 & & & \\ & \ddots & & \\ & & \beta_3 & \\ & & \boldsymbol{0}_{1\times d} \end{bmatrix} \in \mathbb{R}^{(d+1)\times(d)}.$$

To establish convergence for this problem under gaussian initialization, Arora et al. (2019) require the standard deviation to be small enough, while Du and Hu (2019) require larger dimension d_h because their method relies on the condition number of X. Our method balances the requirements on initialization and dimension. The fine-grained nature of our analysis (particularly the *Vector-coupled Dynamics*) enables extension to various problem beyong (Eq. (55)).

E Additional Experimental Results

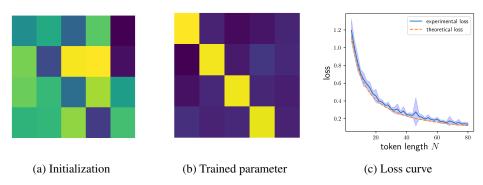


Figure 2: (a) Visualization of matrix product $C^{\top}W_B$ before training; (b) Post-training visualization of matrix product $C^{\top}W_B$; (c) Test loss versus token sequence length N. Blue curve: experimental loss; orange dashed line: theoretical loss $\frac{d}{2}\left(1-\frac{\beta_3^2}{\beta_1}\right)$.

Experiments Setting We follow Section 3 to generate the dateset and initialize the model. Specifically, we set dimension d=4, $d_h=80$, prompt token length N=50, and train the Mamba model on 3000 sequences by gradient descent. Moreover, we vary the length of the prompt token N from 4 to 80 and compare the test loss with the theoretical loss. For each N, we conduct 10 independent experiments and report the averaged results. All experiments are performed on an NVIDIA A800 GPU.

Experiment Result Figure 2a and Figure 2b show that $C^{\top}B$ can be trained to diagonal matrix from random initialization. Figure 2c show that the experimental loss aligns with the theoretical loss $\mathcal{L}(\theta) = \frac{d}{2} \left(1 - \frac{\beta_3^2}{\beta_1}\right)$, noting that the theoretical loss $\left(1 - \frac{\beta_3^2}{\beta_1}\right)$ has an upper bound $\frac{3d(d+1)}{2N}$ that decays linearly with N. These experimental results further verified our theoretical proof.

Mamba vs Linear Attention Optimal linear attention outperforms Mamba under our construction, and they have O(1/N) error upper bound with different constant factors. We provide a comprarison of loss between optimal Mamba (under our Assumption 4.1) with optimal linear attention as in Table 2 with setting $d=10, N=10, 20, \ldots, 80$.

When N is smaller than d We also test the case when $N \le d$ in Table 3 with setting $d = 20, N = 4, 6, \dots, 20$.

Convergence of w_{Δ} We set $w_{\Delta}=0$ in the assumption. Now we show that random initializd $w_{\Delta}=0$ can converge to 0 experimental. The results is in Table 4.

Different d_h Table 5 shows the mean value and standard deviation of the loss for smaller d_h (in 10 repeated experiments). We set d=4, N=30, and the theoretical loss is 0.2954.

Table 2: Comparison of Mamba and Linear Attention

N	10	20	30	40	50	60	70	80
Mamba							0.7022	
Linear Attention	2.6190	1.7742	1.3415	1.0784	0.9016	0.7746	0.6790	0.6044

Table 3: Experiment for $N \leq d$

N	4	6	8	10	12	14	16	18	20
Experimental Loss Theoretical Loss								5.6193 5.4810	

Table 4: Convergence of w_{Δ}

Epoch	0	10	20	30	40	50	60	70	80
$egin{array}{c} \ oldsymbol{w}_\Delta\ _2 \ \ oldsymbol{w}_\Delta\ _2^2 \end{array}$	0.8883	0.7513	0.4821	0.3331	0.2444	0.2026	0.1868	0.1799	0.1773
	0.7891	0.5645	0.2324	0.1109	0.0597	0.0410	0.0349	0.0324	0.0314

Table 5: Different d_h

d_h	6	8	10	12	14	16	18	20
mean(loss) std(loss)		0.2933 0.0055						0.2959 0.0142