

# GATING IS WEIGHTING: UNDERSTANDING GATED LINEAR ATTENTION THROUGH IN-CONTEXT LEARNING

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Linear attention methods provide a strong alternative to softmax attention as they allow for efficient recurrent decoding. Recent research has focused on enhancing standard linear attention by incorporating gating while retaining its computational benefits. Such Gated Linear Attention (GLA) architectures include highly competitive models such as Mamba and RWKV. In this work, we examine the in-context learning capabilities of the GLA model and make the following contributions. We show that a multilayer GLA can implement a general class of Weighted Preconditioned Gradient Descent (WPGD) algorithms with data-dependent weights. These weights are induced by the gating and allows the model to control the contribution of individual tokens to prediction. To further understand the mechanics of weighting, we introduce a novel data model with multitask prompts and characterize the optimization landscape of the problem of learning a WPGD algorithm. We identify mild conditions under which there is a unique (global) minimum up to scaling invariance, and the associated WPGD algorithm is unique as well. Finally, we translate these findings to explore the optimization landscape of GLA and shed light on how gating facilitates context-aware learning and when it is provably better than vanilla linear attention.

## 1 INTRODUCTION

The Transformer architecture (Vaswani, 2017) has become the de facto standard for language modeling tasks. The key component of the Transformer is the self-attention mechanism, which computes softmax-based similarities between all token pairs. Despite its success, the self-attention mechanism has quadratic complexity with respect to sequence length, making it computationally expensive for long sequences. To address this issue, a growing body of work has proposed near-linear time approaches to sequence modeling. The initial approaches included linear attention and state-space models, both achieving  $O(1)$  inference complexity per generated token, thanks to their recurrent form. While these initial architectures typically do not match softmax-attention in performance, recent recurrent models such as Mamba (Gu & Dao, 2023; Dao & Gu, 2024), mLSTM (Beck et al., 2024), GLA Transformer (Yang et al., 2023), and RWKV-6 (Peng et al., 2024) achieve highly competitive results with the softmax Transformer. Notably, as highlighted in Yang et al. (2023), these architectures can be viewed as variants of *gated linear attention* (GLA), which incorporates a gating mechanism within the recurrence of linear attention.

Given a sequence of tokens  $(z_i)_{i=1}^n \subset \mathbb{R}^d$  and associated query, key, and value embeddings  $(q_i, k_i, v_i)_{i=1}^n \subset \mathbb{R}^d$ , with  $d$  being the embedding dimension, the GLA recurrence is given by

$$S_i = G_i \odot S_{i-1} + v_i k_i^\top, \quad \text{and} \quad o_i = S_i q_i. \quad (1)$$

Here,  $S_i \in \mathbb{R}^{d \times d}$  represents the *2D state variable*,  $o_i \in \mathbb{R}^d$  represents the  $i$ 'th output token, and the *gating variable*  $G_i := g(z_i) \in \mathbb{R}^{d \times d}$  is applied to the state through the Hadamard product  $\odot$ . When the gating is removed, the model reduces to causal linear attention (Katharopoulos et al., 2020).

The central objective of this work is to enhance the mathematical understanding of the GLA mechanism. In-context learning (ICL), one of the most remarkable features of modern sequence models, provides a powerful framework to achieve this aim. ICL refers to the ability of a sequence model to implicitly infer functional relationships from the demonstrations provided in its context window

(Brown, 2020; Min et al., 2022). It is inherently related to the model’s ability to emulate learning algorithms. Notably, ICL has been a major topic of empirical and theoretical interest in recent years. More specifically, a series of works have examined the approximation and optimization characteristics of linear attention, and have provably connected linear attention to the preconditioned gradient descent algorithm (Von Oswald et al., 2023; Ahn et al., 2024; Zhang et al., 2024). Given that the GLA recurrence in (1) has a richer design space, this leads us to ask:

**Q:** *What are the ICL capabilities of the GLA mechanism? What learning algorithm does it emulate when presented with an ICL task?*

**Contributions:** The GLA recurrence in (1) enables the sequence model to weight past information in a data-dependent manner through the gating mechanism  $(\mathbf{G}_i)_{i=1}^n$ . Building on this observation, we demonstrate that GLA models can implement a *data-dependent Weighted Preconditioned Gradient Descent (WPGD)* algorithm. Specifically, a one-step version of this algorithm with scalar gating, where all entries of  $\mathbf{G}_i$  are identical, is described by the prediction:

$$\hat{\mathbf{y}} = \mathbf{x}^\top \mathbf{P} \mathbf{X}^\top (\mathbf{y} \odot \boldsymbol{\omega}). \quad (2)$$

Here,  $\mathbf{X} \in \mathbb{R}^{n \times d}$  is the input feature matrix;  $\mathbf{y} \in \mathbb{R}^n$  is the associated label vector;  $\mathbf{x} \in \mathbb{R}^d$  represents the test/query input to predict;  $\mathbf{P} \in \mathbb{R}^{d \times d}$  is the preconditioning matrix; and  $\boldsymbol{\omega} \in \mathbb{R}^n$  weights the individual samples. When  $\boldsymbol{\omega}$  is fixed, we drop “data-dependent” and simply refer to this algorithm as the WPGD algorithm. However, for GLA,  $\boldsymbol{\omega} := \boldsymbol{\omega}(\mathbf{X}, \mathbf{y})$  depends on the data through recursive multiplication of the gating variables. Building on this formalism, we make the following specific contributions:

- **ICL capabilities of GLA (§3):** Through constructive arguments, we demonstrate that a multilayer GLA model can implement data-dependent WPGD iterations, with weights induced by the gating function. This construction sheds light on the role of causal masking and the expressivity distinctions between scalar- and vector-valued gating functions.
- **Landscape of 1-step WPGD (§4):** The GLA  $\Leftrightarrow$  WPGD connection motivates us to ask: *How does WPGD weigh demonstrations in terms of their relevance to the query?* To address this, we study the fundamental problem of learning an optimal WPGD algorithm: Given a tuple  $(\mathbf{X}, \mathbf{y}, \mathbf{x}, \mathbf{y}) \sim \mathcal{D}$ , with  $\mathbf{y}$  being the label associated with the query, we investigate the population risk minimization:

$$\mathcal{L}_{\text{WPGD}}^* := \min_{\mathbf{P}, \boldsymbol{\omega}} \mathcal{L}_{\text{WPGD}}(\mathbf{P}, \boldsymbol{\omega}) \quad \text{where} \quad \mathcal{L}_{\text{WPGD}}(\mathbf{P}, \boldsymbol{\omega}) = \mathbb{E}_{\mathcal{D}} \left[ \left( \mathbf{y} - \mathbf{x}^\top \mathbf{P} \mathbf{X} (\boldsymbol{\omega} \odot \mathbf{y}) \right)^2 \right]. \quad (3)$$

As our primary mathematical contribution, we characterize the loss landscape under a general multitask data setting, where the tasks associated with the demonstrations  $(\mathbf{X}, \mathbf{y})$  have varying degrees of correlation to the target task  $(\mathbf{x}, \mathbf{y})$ . We carefully analyze this loss landscape and show that, under mild conditions, there is a unique (global) minimum  $(\mathbf{P}, \boldsymbol{\omega})$  up to scaling invariance, and the associated WPGD algorithm is also unique.

- **Loss landscape of 1-layer GLA (§5):** The landscape is highly intricate due to the recursively multiplied gating variables. We show that learning the optimal GLA layer can be connected to solving (3) with a constraint  $\boldsymbol{\omega} \in \mathcal{C}$ , where the restriction  $\mathcal{C}$  is induced by the choice of gating function and input space. Solidifying this connection, we introduce a multitask prompt model under which we characterize the loss landscape of GLA and the influence of task correlations. Our analysis and experiments reveal insightful distinctions between linear attention, GLA with scalar gating, and GLA with vector-valued gating.

## 1.1 RELATED WORK

We discuss prior literature under two topics.

**Efficient sequence models.** Recent sequence model proposals – such as RetNet (Sun et al., 2023), Mamba (Gu & Dao, 2023), xLSTM (Beck et al., 2024), GLA Transformer (Yang et al., 2023), RWKV-6 (Peng et al., 2024) – admit efficient recurrent forms while being increasingly competitive with the transformer architecture with softmax-attention. However, we have a rather limited theoretical understanding of these architectures, especially, when it comes to their optimization landscape and ICL capabilities. Park et al. (2024); Grazzi et al. (2024) demonstrate that Mamba is effective in competitive with a transformer of similar size in various ICL tasks whereas Arora et al. (2024);

Jelassi et al. (2024) establish theoretical and empirical shortcomings of recurrent models for solving recall tasks. It is worth mentioning that, GLA models also connect to state-space models and linear RNNs (De et al., 2024; Orvieto et al., 2023; Gu et al., 2021; Fu et al., 2022), as they could be viewed as time-varying SSMs (Dao & Gu, 2024; Sieber et al., 2024). Finally, GLA models are also closely related to implicit self-attention frameworks. For example, the work by Zimmerman et al. (2024) on unified implicit attention highlights how models such as Mamba (Gu & Dao, 2023) and RWKV (Peng et al., 2023) can be viewed under a shared attention mechanism. Additionally, Zong et al. (2024) leverage gated cross-attention for robust multimodal fusion, demonstrating another practical application of gated mechanisms. Both approaches align with GLA’s data-dependent gating, suggesting its potential for explainability and stable fusion tasks.

**Theory of in-context learning.** The theoretical aspects of ICL has been studied by a growing body of works during the past few years (Xie et al.; von Oswald et al., 2023; Gatmiry et al.; Li et al., 2023; Collins et al., 2024; Wu et al., 2023; Fu et al.; Lin & Lee, 2024; Akyürek et al., 2023; Zhang et al., 2023). A subset of these follow the setting of Garg et al. (2022) which investigates the ICL ability of transformers by focusing on prompts where each example is labeled by a task function from a specific function class, such as linear models. Akyürek et al. (2023) focuses on linear regression and provide a transformer construction that can perform a single step of GD based on in-context examples. Similarly, Von Oswald et al. (2023) provide a construction of weights in linear attention-only transformers that can replicate GD steps for a linear regression task on in-context examples. Notably, they observe similarities between their constructed networks and those resulting from training on ICL prompts for linear regression tasks. Building on these, Zhang et al. (2024); Mahankali et al. (2023); Ahn et al. (2024) focus on the loss landscape of ICL for linear attention models. For a single-layer model trained on in-context prompts for random linear regression tasks, Mahankali et al. (2023); Ahn et al. (2024) show that the resulting model performs a single preconditioned GD step on in-context examples in a test prompt, aligning with the findings of Von Oswald et al. (2023). More recent work (Ding et al., 2023) analyzes the challenges of causal masking in causal language models (causalLM), showing that their suboptimal convergence dynamics closely resemble those of online gradient descent with non-decaying step sizes. Additionally, Li et al. (2024) analyzes the landscape of the H3 architecture, an SSM, under the same dataset model. They show that H3 can implement WPGD thanks to its convolutional/SSM filter. However, their WPGD theory is limited to the trivial setting where all weights are same as they utilize the standard prompt model with IID examples and shared task. Departing from prior works, we introduce novel multitask dataset and prompt models under which nontrivial weighting is provably optimal. Through this, we both characterize the loss landscape of WPGD and also study more sophisticated GLA models and connect them to data-dependent WPGD algorithms.

## 2 PROBLEM SETUP

*Notations.*  $\mathbb{R}^d$  is the  $d$ -dimensional real space, with  $\mathbb{R}_+^d$  and  $\mathbb{R}_{++}^d$  as its positive and strictly positive orthants.  $[n]$  denotes  $\{1, \dots, n\}$ . Bold letters, e.g.,  $\mathbf{a}$  and  $\mathbf{A}$ , represent vectors and matrices. The identity matrix of size  $n$  is  $\mathbf{I}_n$ .  $\mathbf{1}$  and  $\mathbf{0}$  denote the all-one and all-zero vectors or matrices of proper size.  $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  is the Gaussian distribution with mean  $\boldsymbol{\mu}$  and covariance  $\boldsymbol{\Sigma}$ . The symbol  $\odot$  denotes the Hadamard product and  $\oslash$  denotes Hadamard division. Given  $\mathbf{a}_{i+1}, \dots, \mathbf{a}_j \in \mathbb{R}^d$ , we use  $\mathbf{a}_{i:j}$  to denote  $\mathbf{a}_{i+1} \odot \dots \odot \mathbf{a}_j$  for  $i < j$ , and  $\mathbf{a}_{i:i} = \mathbf{1}_d$  is the  $d$ -dimensional all ones vector.

The objective of this work is to develop a theoretical understanding of GLA through ICL. The optimization landscape of standard linear attention has been a topic of significant interest in the ICL literature (Ahn et al., 2024; Li et al., 2024). Following these works, we consider the input prompt

$$\mathbf{Z} = [\mathbf{z}_1 \ \dots \ \mathbf{z}_n \ \mathbf{z}_{n+1}]^\top = \begin{bmatrix} \mathbf{x}_1 & \dots & \mathbf{x}_n & \mathbf{x}_{n+1} \\ y_1 & \dots & y_n & 0 \end{bmatrix}^\top \in \mathbb{R}^{(n+1) \times (d+1)}, \quad (4)$$

where tokens encode the input-label pairs  $(\mathbf{x}_i, y_i)_{i=1}^n \subset \mathbb{R}^d \times \mathbb{R}$ . We aim to enable ICL by training a sequence model  $F \in \mathbb{R}^{(n+1) \times (d+1)} \rightarrow \mathbb{R}$  that predicts the label  $y := y_{n+1}$  associated with the query  $\mathbf{x} := \mathbf{x}_{n+1}$ . This model will utilize the demonstrations  $(\mathbf{x}_i, y_i)_{i=1}^n$  to infer the mapping between  $\mathbf{x}$  and  $y$ . Assuming that the data is distributed as  $(y, \mathbf{Z}) \sim \mathcal{D}$ , the ICL objective is defined as

$$\mathcal{L}(F) = \mathbb{E}_{\mathcal{D}} \left[ (y - F(\mathbf{Z}))^2 \right]. \quad (5)$$

**Linear attention and shared-task distribution.** Central to our paper is the choice of the function class  $F$ . When  $F$  is a linear attention model, the prediction  $F(\mathbf{Z})$  takes the form  $\hat{\mathbf{y}} = \mathbf{z}_{n+1}^\top \mathbf{W}_q \mathbf{W}_k^\top \mathbf{Z}^\top \mathbf{Z} \mathbf{W}_v \mathbf{h}$  where  $\mathbf{W}_k, \mathbf{W}_q, \mathbf{W}_v \in \mathbb{R}^{(d+1) \times (d+1)}$  are attention parameters, and  $\mathbf{h} \in \mathbb{R}^{d+1}$  is the linear prediction head. We assume that the in-context input-label pairs follow a *shared-task distribution*, where  $\boldsymbol{\beta} \sim \mathcal{N}(0, \boldsymbol{\Sigma}_\beta)$ ,  $\mathbf{x}_i$  are i.i.d. with  $\mathbf{x}_i \sim \mathcal{N}(0, \boldsymbol{\Sigma}_x)$ , and  $y_i \sim \mathcal{N}(\boldsymbol{\beta}^\top \mathbf{x}_i, \sigma^2)$ , where  $\sigma \geq 0$  represents the noise level. Under this shared-task distribution, it is shown (Von Oswald et al., 2023; Ahn et al., 2024; Zhang et al., 2024) that the optimal one-layer linear attention predictor  $\hat{\boldsymbol{\beta}}$  coincides with the one-step optimal preconditioned gradient descent. In particular, we have  $\hat{\boldsymbol{\beta}} = \mathbf{P}^* \mathbf{X}^\top \mathbf{y}$ , where

$$\mathbf{P}^* = \arg \min_{\mathbf{P} \in \mathbb{R}^{d \times d}} \mathbb{E}_{\mathcal{D}} \left[ (\mathbf{y} - \mathbf{x}^\top \mathbf{P} \mathbf{X}^\top \mathbf{y})^2 \right] \quad \text{with} \quad \mathbf{X} := [\mathbf{x}_1 \quad \cdots \quad \mathbf{x}_n]^\top \quad \text{and} \quad \mathbf{y} := \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}. \quad (6)$$

**Linear attention and gating.** Given the input prompt  $\mathbf{Z}$ , let  $\mathbf{Q} = \mathbf{Z} \mathbf{W}_q$ ,  $\mathbf{K} = \mathbf{Z} \mathbf{W}_k$  and  $\mathbf{V} = \mathbf{Z} \mathbf{W}_v$  be the corresponding query, key, and value embedding matrices, respectively. The output of *causal* linear attention at time  $i$  can be computed in a recurrent form as  $\mathbf{S}_i = \mathbf{S}_{i-1} + \mathbf{v}_i \mathbf{k}_i^\top$  and  $\mathbf{o}_i = \mathbf{S}_i \mathbf{q}_i$  where  $\mathbf{q}_i, \mathbf{k}_i, \mathbf{v}_i \in \mathbb{R}^{d+1}$  are the query, key, value embeddings of  $\mathbf{z}_i$  and  $\mathbf{S}_0 = \mathbf{0}$ . This recurrent form implies that linear attention has  $O(d^2)$  cost, that is independent of  $N$ , to generate per-token. As presented in (1), GLA follows the same structure as linear attention but with a gating mechanism, which equips the model with the option to pass or suppress the history. As discussed in Yang et al. (2023), the different choices of the gating function correspond to different popular recurrent architectures such as Mamba (Gu & Dao, 2023), Mamba2 (Dao & Gu, 2024), RWKV (Peng et al., 2024), etc.

We will show that GLA can weigh the context window through gating, thus, its capabilities are linked to the WPGD algorithm described in (7). This will in turn facilitate GLA to effectively learn *multitask prompt distributions* described by  $y_i \sim \mathcal{N}(\boldsymbol{\beta}_i^\top \mathbf{x}_i, \sigma^2)$  with  $\boldsymbol{\beta}_i$ 's not necessarily identical.

### 3 WHAT GRADIENT METHODS CAN GLA EMULATE?

In this section, we investigate the ICL capabilities of gated linear attention (GLA) and show that under suitable instantiations of model weights, GLA can implement *data-dependent* WPGD.

#### 3.1 GLA AS A DATA-DEPENDENT WPGD PREDICTOR

**Data-Dependent WPGD.** Given  $\mathbf{X}$  and  $\mathbf{y}$  as defined in (6), consider the weighted least squares objective  $\mathcal{L}(\boldsymbol{\beta}) = \sum_{i=1}^n \Omega_i \cdot (y_i - \boldsymbol{\beta}^\top \mathbf{x}_i)^2$  with weights  $\boldsymbol{\Omega} \in \mathbb{R}^n$ . To optimize this, we use gradient descent (GD) starting from zero initialization,  $\boldsymbol{\beta}_0 = \mathbf{0}$  with a step size of  $\eta = 1/2$ . One step of standard GD is given by

$$\boldsymbol{\beta}_1 = \boldsymbol{\beta}_0 - \eta \nabla \mathcal{L}(\boldsymbol{\beta}_0) = \sum_{i=1}^n \Omega_i \cdot \mathbf{x}_i y_i = \mathbf{X}^\top (\boldsymbol{\Omega} \odot \mathbf{y}).$$

Given a test/query feature  $\mathbf{x}$ , the corresponding prediction is  $\hat{\mathbf{y}} = \mathbf{x}^\top \hat{\boldsymbol{\beta}}$  where  $\hat{\boldsymbol{\beta}} = \boldsymbol{\beta}_1$ . Additionally, if we were using *preconditioned* GD with a preconditioning/projection matrix  $\mathbf{P} \in \mathbb{R}^{d \times d}$ , one step iteration would take the form

$$\hat{\mathbf{y}} = \mathbf{x}^\top \hat{\boldsymbol{\beta}}, \quad \text{where} \quad \hat{\boldsymbol{\beta}} = \mathbf{P} \boldsymbol{\beta}_1 = \mathbf{P} \mathbf{X}^\top (\boldsymbol{\Omega} \odot \mathbf{y}).$$

Above is the basic *scalar-weighted* WPGD predictor which weights individual datapoints. It turns out, *vector-valued gating* can facilitate a more general estimator which weights individual coordinates. To this aim, we introduce an extension as follows: Let  $\mathbf{P}_1, \mathbf{P}_2 \in \mathbb{R}^{d \times d}$  denote the preconditioning matrices, and let  $\boldsymbol{\Omega} \in \mathbb{R}^{n \times d}$  denote the *vector-valued weighting* matrix. Note that  $\boldsymbol{\Omega}$  is now a matrix rather than vector to facilitate coordinate-wise weighting and will remain consistent throughout the paper. We can similarly define

$$\boldsymbol{\beta}_1^{\text{gd}}(\mathbf{P}_1, \mathbf{P}_2, \boldsymbol{\Omega}) := \mathbf{P}_2 (\mathbf{X} \mathbf{P}_1 \odot \boldsymbol{\Omega})^\top \mathbf{y} \quad (7a)$$

as one-step of (generalized) WPGD. Its corresponding prediction on a test query  $\mathbf{x}$  is:

$$\hat{\mathbf{y}} = \mathbf{x}^\top \hat{\boldsymbol{\beta}}, \quad \text{where} \quad \hat{\boldsymbol{\beta}} = \boldsymbol{\beta}_1^{\text{gd}}(\mathbf{P}_1, \mathbf{P}_2, \boldsymbol{\Omega}). \quad (7b)$$

We note that by removing the preconditioning matrices  $P_1$ ,  $P_2$ , and the weighting matrix  $\Omega$  in (7a), it reduces to standard GD. We also note that Li et al. (2024) demonstrates that H3-like models implement one-step WPGD, where the weighting is example-wise, i.e., setting  $\Omega = \omega \mathbf{1}_d^\top$ , and they focus on the shared-task distribution where  $\beta_i \equiv \beta$ . In contrast, our work considers a more general data setting where tasks within an in-context prompt are not necessarily identical.

We first introduce the following model constructions under which we establish the equivalence between GLA (c.f. (1)) and WPGD (c.f. (7)) with the weighting matrix induced by the input data and the gating function. Inspired by previous works (Von Oswald et al., 2023; Ahn et al., 2024), we consider the following restricted attention matrices:

$$W_k = \begin{bmatrix} P_k^\top & \mathbf{0} \\ \mathbf{0} & 0 \end{bmatrix}, \quad W_q = \begin{bmatrix} P_q^\top & \mathbf{0} \\ \mathbf{0} & 0 \end{bmatrix} \quad \text{and} \quad W_v = \begin{bmatrix} \mathbf{0}_{d \times d} & \mathbf{0} \\ \mathbf{0} & 1 \end{bmatrix}, \quad (8)$$

where  $P_k, P_q \in \mathbb{R}^{d \times d}$ . Here note that we set the  $(d+1, d+1)$ 'th entry of  $W_v$  to be one for simplification. More generally, it can be any nonzero number, e.g.,  $v \in \mathbb{R}$ . Then parameterizing  $W_q$  with  $P_q/v$  returns the same output as from (8).

**Theorem 1.** Recall the GLA from (1) and input sequence  $\mathbf{Z}$  from (4), and suppose that at time  $i$ , gating function has the form of  $g(\mathbf{z}_i) = \mathbf{G}_i \in \mathbb{R}^{(d+1) \times (d+1)}$ . Considering model construction in (8) and prediction head  $\mathbf{h} = \mathbf{1}$ , the single-layer GLA prediction returns

$$f_{\text{GLA}}(\mathbf{Z}) := \mathbf{o}_{n+1}^\top \mathbf{h} = \hat{\beta}^\top \mathbf{x} \quad \text{where} \quad \hat{\beta} = \beta_1^{\text{gd}}(P_k, P_q, \Omega).$$

Here,  $\beta_1^{\text{gd}}(\cdot)$  is a one-step WPGD feature predictor defined in (7a),  $P_k, P_q$  correspond to attention weights following (8), and  $\Omega = [\mathbf{g}_{1:n+1} \quad \mathbf{g}_{2:n+1} \quad \cdots \quad \mathbf{g}_{n:n+1}]^\top \in \mathbb{R}^{n \times d}$  where  $\mathbf{g}_{i:n+1}, i \in [n]$  is given by

$$\mathbf{g}_{i:n+1} := (\mathbf{g}_{i+1} \odot \mathbf{g}_{i+2} \cdots \mathbf{g}_{n+1}) \in \mathbb{R}^d \quad \text{and} \quad \mathbf{G}_i = \begin{bmatrix} * & * \\ \mathbf{g}_i^\top & * \end{bmatrix} \quad (9)$$

Here and throughout, we use  $*$  to fill the entries of the matrices that do not affect the final output, and based on the model construction given in (8), these entries can be assigned any value.

Observe that, crucially, since  $\mathbf{g}_i$  (or  $\mathbf{G}_i$ ) is associated with  $\mathbf{z}_i$ ,  $\mathbf{z}_i$  influences the weighting of all history  $\mathbf{z}_{j < i}$ . We defer the proof of Theorem 1 to the Appendix B.1. It is noticeable that only  $d$  of the total  $(d+1)^2$  entries in each gating matrix  $\mathbf{G}_i$  are useful due to the model construction presented in (8). However, if we relax the weight restriction, e.g.,  $W_v = [\mathbf{0}_{(d+1) \times d} \quad \mathbf{1}_{d+1}]$ , then the weighting matrix  $\Omega$  in Theorem 1 is associated with all rows of the  $\mathbf{G}_i$  matrices. We defer the discussion to Appendix B.1.

### 3.2 CAPABILITIES OF MULTI-LAYER GLA

Ahn et al. (2024) demonstrated that, with appropriate construction, an  $L$ -layer linear attention model performs  $L$ -step preconditioned gradient descent on the dataset  $(\mathbf{x}_i, y_i)_{i=1}^n$  provided within the prompt. In this work, we study multi-layer GLA and analyze the associated algorithm class it can emulate. It is worth mentioning that Ahn et al. (2024) does not consider *causal masking* which is integral to multilayer GLA due to its recurrent nature described in (1). Our analysis will capture the impact of gating and causal mask through  $n$  separate gradient descent trajectories that are coupled.

Consider an  $L$ -layer GLA model. For  $\ell \in [L]$ , let  $\mathbf{Z}_\ell$  and  $\mathbf{O}_\ell$  denote the input and output of the  $\ell$ 'th layer. In practice, residual connections are commonly applied. Hence, we define the updated output of the  $\ell$ 'th layer (after applying the residual connection) as  $\tilde{\mathbf{O}}_\ell := \mathbf{Z}_\ell + \mathbf{O}_\ell$ . Note that  $\tilde{\mathbf{O}}_\ell$  also serves as the input to the  $(\ell+1)$ 'th layer, i.e.,  $\mathbf{Z}_{\ell+1} = \tilde{\mathbf{O}}_\ell$ . In the following, we focus on  $(d+1)$ 'th entries of each token's output at each layer, denoted by  $\tilde{o}_{i,\ell} := (\tilde{\mathbf{O}}_\ell)_{i,d+1}$  for  $i \in [n+1], \ell \in [L]$ .

**Theorem 2.** Consider an  $L$ -layer GLA with residual connections, where  $W_k$  and  $W_q$  in the  $\ell$ 'th layer are parameterized by  $P_{k,\ell}, P_{q,\ell} \in \mathbb{R}^{d \times d}$ , following (8), for  $\ell \in [L]$ . Let the gating be a function of the features, e.g.,  $\mathbf{G}_i = g(\mathbf{x}_i)$ , and let  $\Omega$  be defined as in Theorem 1. Additionally, denote the masking as

$$\mathbf{M}_i = \begin{bmatrix} \mathbf{I}_i & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \in \mathbb{R}^{n \times n}, \quad \text{and let } \hat{\beta}_0, \beta_{i,0} = \mathbf{0} \text{ for } i \in [n].$$

Then the  $(d+1)$ 'th entry of the  $i$ 'th token at the  $\ell$ 'th layer outputs:

- For  $i \leq n$ ,  $\tilde{o}_{i,\ell} = y_i - \mathbf{x}_i^\top \beta_{i,\ell}$  where  $\beta_{i,\ell} = \beta_{i,\ell-1} + P_{q,\ell}(\nabla_{i,\ell} \odot \mathbf{g}_{i:n+1})$ ,



- $\tilde{o}_{n+1,\ell} = -\mathbf{x}^\top \hat{\beta}_\ell$  where  $\hat{\beta}_\ell = (1 + \alpha_\ell) \hat{\beta}_{\ell-1} + \mathbf{P}_{q,\ell} (\nabla_{n,\ell} \odot \mathbf{g}_{n+1})$  and  $\alpha_\ell = \mathbf{x}^\top \mathbf{P}_{q,\ell} \mathbf{P}_{k,\ell}^\top \mathbf{x}$ .

Here, letting  $\mathbf{B}_\ell = [\beta_{1,\ell} \cdots \beta_{n,\ell}]^\top$ ,  $\bar{\mathbf{X}}_\ell = \mathbf{X} \mathbf{P}_{k,\ell} \odot \boldsymbol{\Omega}$ , and  $\hat{\mathbf{y}}_\ell = (\mathbf{X} \odot \mathbf{B}_{\ell-1}) \mathbf{1}$ , we define

$$\nabla_{i,\ell} = \bar{\mathbf{X}}_\ell^\top \mathbf{M}_i (\hat{\mathbf{y}}_\ell - \mathbf{y}).$$

We defer the proof of Theorem 2 to the Appendix B.2. Theorem 2 states that an  $L$ -layer GLA implements  $L$  steps of WPGD but with gradient in a recurrent form. To recap, given data  $(\mathbf{X}, \mathbf{y})$  and prediction  $\hat{\beta}$ , the gradient with respect to the squared loss takes the form  $\mathbf{X}^\top (\mathbf{X} \hat{\beta} - \mathbf{y})$ , up to some constant  $c$ . In comparison,  $\mathbf{P}_{q,\ell} (\nabla_{i,\ell} \odot \mathbf{g}_{i:n+1})$  similarly acts as a gradient but incorporates layer-wise feature preconditioners  $(\mathbf{P}_{q,\ell}, \mathbf{P}_{k,\ell})$ , data weighting  $(\boldsymbol{\Omega})$ , and causality  $(\mathbf{g}_{i:n+1}, \mathbf{M}_i)$ . Here,  $\mathbf{M}_i$  represents causal masking, ensuring that at time  $i$ , only inputs from  $j \leq i$  are used for prediction. Notably, the recurrent structure of GLA allows the gating mechanism to apply context-dependent weighting strategies. These results are consistent with Ding et al. (2023), which demonstrate that causal masking limits convergence by introducing sequence biases, akin to online gradient descent with non-decaying step sizes.

To simplify the theorem statement, we assume that the gating function depends only on the input feature, e.g.,  $\mathbf{G}_i = g(\mathbf{x}_i)$ , ensuring that the corresponding data-dependent weighting is uniform across all layers. This assumption is included solely for clarity in the theorem statement, and the complete result is provided in Appendix B.2. Note that our inclusion of the additional term  $\alpha_\ell$  captures the influence of the last token's output on the next layer's prediction, which is not addressed by Ahn et al. (2024). Based on the above multi-layer GLA result, we have the following corollary for multi-layer linear attention network with causal mask in each layer.

**Corollary 1.** Consider an  $L$ -layer linear attention model with causal mask and residual connection in each layer. Let  $\ell$ 'th layer be parameterized by  $\mathbf{P}_{q,\ell}, \mathbf{P}_{k,\ell}$  as in (8) and define  $\mathbf{P}_\ell := \mathbf{P}_{q,\ell} \mathbf{P}_{k,\ell}^\top$ ,  $\ell \in [L]$ . Let  $\hat{\beta}_0, \beta_{i,0} = \mathbf{0}$  for  $i \in [n]$ . Then, the  $(d+1)$ 'th entry of the  $i$ 'th token of the  $\ell$ 'th layer outputs satisfies:

- For  $i \leq n$ ,  $\tilde{o}_{i,\ell} = y_i - \mathbf{x}_i^\top \beta_{i,\ell}$  where  $\beta_{i,\ell} = \beta_{i,\ell-1} + \mathbf{P}_\ell \nabla_{i,\ell}$ ,
- $\tilde{o}_{n+1,\ell} = -\mathbf{x}^\top \hat{\beta}_\ell$  where  $\hat{\beta}_\ell = (1 + \alpha_\ell) \hat{\beta}_{\ell-1} + \mathbf{P}_\ell \nabla_{n,\ell}$  and  $\alpha_\ell = \mathbf{x}^\top \mathbf{P}_\ell \mathbf{x}$ .

Here, we define  $\nabla_{i,\ell} = \mathbf{X}^\top \mathbf{M}_i (\hat{\mathbf{y}}_\ell - \mathbf{y})$  with  $\hat{\mathbf{y}}_\ell, \mathbf{M}_i$  following the same definitions as in Theorem 2.

Our theoretical results in Theorem 2 focus on multi-layer GLA without Multi-Layer Perceptron (MLP) layers to isolate and analyze the effects of the gating mechanism. However, MLP layers, a key component of standard Transformers, facilitate further nonlinear feature transformations and interactions, potentially enhancing GLA's expressive power. Future work could explore the theoretical foundations of integrating MLPs into GLA and analyze the optimization landscape of general gated attention models, aligning them more closely with conventional Transformer architectures (Gu & Dao, 2023; Dao & Gu, 2024; Peng et al., 2024).

### 3.3 GLA WITH SCALAR GATING

Theorem 1 establishes a connection between 1-layer GLA (c.f. (1)) and one-step WPGD (c.f. (7)), where the weighting in WPGD corresponds to the gating  $g(\mathbf{z}_i) = \mathbf{G}_i$  in GLA, as detailed in Theorem 1. Now let us consider the widely used types of gating functions, such as  $\mathbf{G}_i = \alpha_i \mathbf{1}_{d+1}^\top$  (Yang et al., 2023; Katsch, 2023; Qin et al., 2024; Peng et al., 2024) or  $\mathbf{G}_i = \gamma_i \mathbf{1}_{d+1} \mathbf{1}_{d+1}^\top$  (Dao & Gu, 2024; Beck et al., 2024; Peng et al., 2021; Sun et al., 2024) where  $\alpha_i \in \mathbb{R}^{d+1}$  and  $\gamma_i \in \mathbb{R}$ . In both cases, the gating matrices in (9) take the form of  $\begin{bmatrix} * & * \\ \mathbf{g}_i \mathbf{1}_d^\top & * \end{bmatrix}$ , thus simplifying the predictor to a sample-weighted PGD, as given by

$$f_{\text{GLA}}(\mathbf{Z}) = \hat{\beta}^\top \mathbf{x}, \quad \text{with} \quad \hat{\beta} = \mathbf{P} \mathbf{X}^\top (\omega \odot \mathbf{y}), \quad (10)$$

where  $\mathbf{P} = \mathbf{P}_q \mathbf{P}_k^\top$  and  $\omega = [g_{1:n+1} \cdots g_{n:n+1}]^\top \in \mathbb{R}^n$ . In the remainder, we will mostly focus on the 1-layer GLA with scalar gating as presented in (10).

#### 4 OPTIMIZATION LANDSCAPE OF WPGD

In this section, we explore the problem of learning the optimal sample-weighted PGD algorithm described in (10), a key step leading to our analysis of GLA. The problem is as follows. Recap from (6) that we are given the tuple  $(x, y, X, y) \sim \mathcal{D}$ , where  $X \in \mathbb{R}^{n \times d}$  is the input matrix,  $y \in \mathbb{R}^n$  is the label vector,  $x \in \mathbb{R}^d$  is the query, and  $y \in \mathbb{R}$  is its associated label. The goal is to use  $X, y$  to predict  $y$  given  $x$  via the 1-step WPGD prediction  $\hat{y} = x^\top \hat{\beta}$ , with  $\hat{\beta}$  as in (10). The algorithm learning problem is given by (3) which minimizes the WPGD risk  $\mathbb{E}_{\mathcal{D}}[(y - x^\top P X(\omega \odot y))^2]$ .

Prior research (Mahankali et al., 2023; Li et al., 2024; Ahn et al., 2024) has studied the problem of learning PGD when input-label pairs follow an IID distribution. It is worth noting that while Li et al. (2024) establishes a connection between H3-like models and (10) similar to ours, their work assumes that the optimal  $\omega$  consists of all ones and does not specifically explore the optimization landscape of  $\omega$  when in-context samples are non-IID. Departing from this, we introduce a realistic model where each input-label pair is allowed to come from a distinct task.

**Definition 1** (Correlated task model). *Suppose  $\beta_i \in \mathbb{R}^d \sim \mathcal{N}(0, I)$  are jointly Gaussian for  $i \in [n + 1]$ . Define the pairwise correlations  $r_{ij} = \mathbb{E}[\beta_i^\top \beta_j]/d$  for  $i, j \in [n + 1]$ , and the task and correlation matrices*

$$\beta := \beta_{n+1}, \quad B = [\beta_1 \dots \beta_n]^\top, \quad R = \frac{1}{d} \mathbb{E}[B B^\top], \quad \text{and} \quad r = \frac{1}{d} \mathbb{E}[B \beta]. \quad (11)$$

Additionally, for any  $i, j \in [n + 1]$ ,  $\beta_i - r_{ij}\beta_j$  is independent of  $\beta_j$ .

Note that in (11), we have  $B \in \mathbb{R}^{n \times d}$ ,  $R \in \mathbb{R}^{n \times n}$ , and  $r \in \mathbb{R}^n$ , with normalization ensuring that the entries of  $R$  and  $r$  lie in the range  $[-1, 1]$ , corresponding to correlation coefficients.

**Definition 2** (Multitask distribution).  *$(\beta_i)_{i=1}^{n+1}$  are drawn according to the correlated task model of Definition 1,  $(x_i)_{i=1}^{n+1} \in \mathbb{R}^d$  are IID following  $x_i \sim \mathcal{N}(0, \Sigma)$  and  $y_i \sim \mathcal{N}(x_i^\top \beta_i, \sigma^2)$  for  $i \in [n + 1]$ .*

**Definition 3.** *Let the eigen decompositions of  $\Sigma$  and  $R$  be denoted by  $\Sigma = U \text{diag}(s) U^\top$  and  $R = E \text{diag}(\lambda) E^\top$ , where  $s = [s_1, \dots, s_d]^\top \in \mathbb{R}_{++}^d$  and  $\lambda = [\lambda_1, \dots, \lambda_n]^\top \in \mathbb{R}_+^n$ . Let  $s_{\min}$  and  $s_{\max}$  denote the smallest and largest eigenvalues of  $\Sigma$ , respectively. Further, let  $\lambda_{\min}$  and  $\lambda_{\max}$  denote the nonzero smallest and largest eigenvalues of  $R$ . Define the effective spectral gap of  $\Sigma$  and  $R$ , respectively, as*

$$\Delta_\Sigma := s_{\max} - s_{\min}, \quad \text{and} \quad \Delta_R := \lambda_{\max} - \lambda_{\min}. \quad (12)$$

**Assumption A.** *For the correlation vector  $r$  from (11), we have  $r = E a$  for some  $a = [a_1, \dots, a_n]^\top \in \mathbb{R}^n$  with at least one nonzero  $a_i$ .*

Assumption A essentially ensures that  $r$  (representing the correlations between in-context tasks) can be expressed as a linear transformation of a vector  $a$  of nonzero values. This guarantees that the correlation structure is non-degenerate, meaning that all elements of  $r$  are influenced by meaningful correlations. Assumption A avoids trivial cases where there are no correlations between tasks. By requiring at least one nonzero element in  $a$ , the assumption ensures that the tasks are interrelated.

The following theorem characterizes the stationary points  $(P, \omega)$  of the WPGD objective in (3).

**Theorem 3.** *Consider independent linear data as described in Definition 2. Suppose Assumption A on the correlation vector  $r$  holds. Let the functions  $h : \mathbb{R}_+ \rightarrow \mathbb{R}_+$  and  $g : \mathbb{R}_+ \rightarrow \mathbb{R}_+$  be defined as*

$$h(\tilde{\gamma}) := \sum_{i=1}^n \frac{\lambda_i a_i^2}{(1 + \lambda_i \tilde{\gamma})^2} \left( \sum_{i=1}^n \frac{a_i^2}{(1 + \lambda_i \tilde{\gamma})^2} \right)^{-1}, \quad (13a)$$

$$g(\gamma) := \left( 1 + M \sum_{i=1}^d \frac{s_i^2}{(M + s_i(\gamma + 1))^2} \left( \sum_{i=1}^d \frac{s_i^3}{(M + s_i(\gamma + 1))^2} \right)^{-1} \right)^{-1}, \quad (13b)$$

where  $\{s_i\}_{i=1}^d$  and  $\{\lambda_i\}_{i=1}^n$  are the eigenvalues of  $\Sigma$  and  $R$ , respectively;  $\{a_i\}_{i=1}^n$  are as given in Assumption A; and  $M = \sigma^2 + \sum_{i=1}^d s_i$ .

The risk function  $\mathcal{L}(P, \omega)$  in (3) has a stationary point  $(P^*, \omega^*)$ , up to rescaling, defined as

$$P^* = \Sigma^{-\frac{1}{2}} \left( \frac{\gamma^* + 1}{\sigma^2 + \text{tr}(\Sigma)} \cdot \Sigma + I \right)^{-1} \Sigma^{-\frac{1}{2}}, \quad \text{and} \quad \omega^* = (g(\gamma^*) \cdot R + I)^{-1} r, \quad (14)$$

where  $\gamma^*$  is a fixed point of composite function  $h(g(\gamma))$ .

Theorem 3 characterizes the stationary points  $(\mathbf{P}^*, \omega^*)$ , which exist up to re-scaling. This result presents the first landscape analysis of GLA for the joint learning of  $(\mathbf{P}, \omega)$ , while also exploring the stationary points  $(\mathbf{P}^*, \omega^*)$ . In the following, we provide mild conditions on effective spectral gaps of  $\mathbf{R}$  and  $\Sigma$  under which a unique (global) minimum  $(\mathbf{P}^*, \omega^*)$  exists.

**Theorem 4** (Uniqueness of the WPGD Predictor). *Consider independent linear data as given in Definition 2. Suppose Assumption A on the correlation vector  $\mathbf{r}$  holds, and*

$$\Delta_{\Sigma} \cdot \Delta_{\mathbf{R}} < M + s_{\min}, \quad (15)$$

where  $\Delta_{\Sigma}$  and  $\Delta_{\mathbf{R}}$  denote the effective spectral gaps of  $\Sigma$  and  $\mathbf{R}$ , respectively, as given in (12);  $s_{\min}$  is the smallest eigenvalue of  $\Sigma$ ; and  $M = \sigma^2 + \sum_{i=1}^d s_i$ .

**T1** The composite function  $h(g(\gamma))$  is a contraction mapping and admits a unique fixed point  $\gamma = \gamma^*$ .

**T2** The function  $\mathcal{L}(\mathbf{P}, \omega)$  has a unique (global) minima  $(\mathbf{P}^*, \omega^*)$ , up to re-scaling, given by (14).

*Proof Sketch.* Let  $\gamma := \frac{\omega^\top \mathbf{R} \omega}{\|\omega\|^2}$ . Note that  $\gamma \geq 0$  since  $\mathbf{R}$  is positive semi-definite. From the first-order optimality condition, the solution to (3) takes the following form:

$$\mathbf{P}(\gamma) = C(\mathbf{r}, \omega, \Sigma) \cdot \Sigma^{-\frac{1}{2}} \left( \frac{\gamma + 1}{\sigma^2 + \text{tr}(\Sigma)} \cdot \Sigma + \mathbf{I} \right)^{-1} \Sigma^{-\frac{1}{2}}, \quad (16a)$$

$$\omega(\gamma) = c(\mathbf{r}, \omega, \Sigma) \cdot (g(\gamma) \cdot \mathbf{R} + \mathbf{I})^{-1} \mathbf{r}, \quad (16b)$$

for some constants  $C(\mathbf{r}, \omega, \Sigma)$  and  $c(\mathbf{r}, \omega, \Sigma)$ .

Substituting the expression for  $\omega(\gamma)$  into  $\gamma = \frac{\omega^\top \mathbf{R} \omega}{\|\omega\|^2}$ , and applying Assumption A, we obtain the equation  $\gamma = h(g(\gamma))$ . We then show that whenever  $\Delta_{\Sigma} \cdot \Delta_{\mathbf{R}} < M + s_{\min}$ , the mapping  $h(g(\gamma))$  is a contraction (see Lemma 1). By the Banach Fixed-Point Theorem, this guarantees the existence of a unique fixed point  $\gamma = \gamma^*$ , where  $\gamma^* = h(g(\gamma^*))$ . Finally, substituting  $\gamma^*$  into (16) implies that  $(\mathbf{P}^*, \omega^*)$ , as given in (14), is a unique (global) minima of (3), up to re-scaling. See Appendix C.2 for the complete proof of Theorem 4.  $\square$

Theorem 4 establishes mild conditions under which a unique (global) minimum  $(\mathbf{P}^*, \omega^*)$  exists, up to scaling invariance, and guarantees the uniqueness of the associated WPGD algorithm. It provides the first global landscape analysis for GLA and generalizes prior work (Li et al., 2024; Ahn et al., 2024) on the global landscape by extending the optimization properties of linear attention to the more complex *nonconvex* GLA with joint  $(\mathbf{P}, \omega)$  optimization.

**Remark 1** An interesting observation about the optimal gating parameter  $\omega^*$  is its connection to the correlation matrix  $\mathbf{R}$ , which captures the task correlations in a multitask learning setting. Specifically, the optimal gating given in (14) highlights how  $\omega^*$  depends directly on both the task correlation matrix  $\mathbf{R}$  and the vector  $\mathbf{r}$ , which encodes the correlations between the tasks and the target task.

**Remark 2** Condition (15) provides a *sufficient* condition for the uniqueness of a fixed point. This implies that whenever  $\Delta_{\Sigma} \cdot \Delta_{\mathbf{R}} < M + s_{\min}$ , the mapping  $h(g(\gamma))$  is a contraction, ensuring the existence of a unique fixed point. However, there may be cases where the mapping  $h(g(\gamma))$  does not satisfy Condition (15), yet a unique fixed point (and a unique global minimum) still exists. This is because the Banach Fixed-Point Theorem does not provide a *necessary* condition.

**Corollary 2.** Suppose  $\Sigma = \mathbf{I}$ . Then,  $\Delta_{\Sigma} = 0$ , satisfying Condition (15), and we have  $g(\gamma^*) = \frac{1}{d + \sigma^2 + 1}$ , which yields

$$\mathbf{P}^* = \mathbf{I}, \quad \text{and} \quad \omega^* = (\mathbf{R} + (d + \sigma^2 + 1)\mathbf{I})^{-1} \mathbf{r}. \quad (17)$$

Thus, the optimal risk  $\mathcal{L}_{\text{WPGD}}^*$  defined in (3) is given by

$$\mathcal{L}_{\text{WPGD}}^* = d + \sigma^2 - d \cdot \mathbf{r}^\top (\mathbf{R} + (d + \sigma^2 + 1)\mathbf{I})^{-1} \mathbf{r}. \quad (18)$$

## 5 OPTIMIZATION LANDSCAPE OF GLA

In Section 3, we demonstrated that GLA implements a data-dependent WPGD algorithm. Building on this, in Section 4, we analyze the optimization landscape for minimizing the 1-step WPGD risk



(c.f. (3)) and show that a unique solution achieves the global minimum of the WPGD algorithm. However, in GLA, the search space for  $\omega$  is restricted and data-dependent, meaning that  $\mathcal{L}_{\text{WPGD}}^*$  in (3) represents the best possible risk a GLA model can achieve. In this section, we analyze the loss landscape for training a 1-layer GLA model and explore the scenarios under which GLA can reach the optimal WPGD risk.

### 5.1 MULTI-TASK PROMPT MODEL

We consider the following multi-task prompts setting with  $K$  correlated tasks  $(\beta_k)_{k=1}^K$ , and 1 query task  $\beta$ . For each correlated task, draw a length  $n_k$  prompt with IID input-label pairs  $\{(\mathbf{x}_i^{(k)}, y_i^{(k)})\}_{i=1}^{n_k}$  to obtain sequences  $(\mathbf{Z}_k)_{k=1}^K$  and the query example is given by  $\mathbf{z} := (\mathbf{x}, y \sim \mathcal{N}(\mathbf{x}^\top \beta, \sigma^2))$ . Let  $n := \sum_{k=1}^K n_k$ . These sequences  $(\mathbf{Z}_k)_{k=1}^K$  as well as query token  $\mathbf{z}$  are concatenated to form a single prompt  $\mathbf{Z}$ . Recap the GLA prediction from (1) and let  $f_{\text{GLA}}(\mathbf{Z})$  be the GLA prediction as defined in Theorem 1. Additionally, consider the model construction as presented in (8) with  $\mathbf{P}_q, \mathbf{P}_k \in \mathbb{R}^{d \times d}$  being the trainable parameters. Then the GLA optimization problem is described as follows:

$$\mathcal{L}_{\text{GLA}}^* := \min_{\mathbf{P}_k, \mathbf{P}_q, g} \mathcal{L}_{\text{GLA}}(\mathbf{P}_k, \mathbf{P}_q, g) \quad \text{where} \quad \mathcal{L}_{\text{GLA}}(\mathbf{P}_k, \mathbf{P}_q, g) = \mathbb{E}_{\mathcal{D}} \left[ (y - f_{\text{GLA}}(\mathbf{Z}))^2 \right]. \quad (19)$$

Here,  $g \in \mathcal{G}$  represents the gating function.

Note that 1) the task vectors  $(\beta_k)_{k=1}^K$  are not explicitly shown in the prompt, 2) examples  $(\mathbf{x}_i^{(k)}, y_i^{(k)})$  are randomly drawn, and 3) the gating function is applied to the tokens/input samples  $(\mathbf{Z}_k)_{k=1}^K$ . Given the above three evidences, the implicit weighting induced by the GLA model varies across different prompts, and it prevents the GLA from learning the optimal weighting.

To address this, we introduce delimiters to mark the boundary of each task. Let  $(\mathbf{d}_k)_{k=1}^K$  be the delimiters that determine stop of the tasks. Specifically, the final prompt is given by

$$\mathbf{Z} = \left[ \mathbf{Z}_1^\top \quad \mathbf{d}_1 \quad \cdots \quad \mathbf{Z}_K^\top \quad \mathbf{d}_K \quad \mathbf{z} \right]^\top. \quad (20)$$

Additionally, to decouple the influence of gating and data, we envision that each token is  $\mathbf{z}_i = [\mathbf{x}_i, y_i, \mathbf{c}_i]$  where  $\mathbf{c}_i \neq \mathbf{0} \in \mathbb{R}^p$  is the contextual features with  $p$  being its dimension and  $(\mathbf{x}_i, y_i)$  are the data features.

- For task prompts  $\mathbf{Z}_k$ : Contextual features are set to a fixed vector  $\tilde{\mathbf{d}}_0 \neq \mathbf{0}$ .
- For delimiters  $\mathbf{d}_k$ : Data features are set to zero (e.g.,  $\mathbf{x}_i = \mathbf{0}$  and  $y_i = 0$ ) so that  $\mathbf{d}_k = [\mathbf{0}_{d+1} \quad \tilde{\mathbf{d}}_k]$  where  $\tilde{\mathbf{d}}_k$  denotes the context vector.

Note that explicit delimiters have been utilized to address real-world problems (Wang et al., 2024; Asai et al., 2022; Dun et al., 2023) due to their ability to improve efficiency and enhance generalization, particularly in task-mixture or multi-document scenarios. To further verify our claim and motivate the introduction of  $(\mathbf{d}_k)_{k=1}^K$ , in Figure 1, we present the results of GLA training with and without delimiters, shown by the red and green curves, respectively. The black dashed curves represent the optimal WPGD loss  $\mathcal{L}_{\text{WPGD}}^*$  under different scenarios, and training GLA without delimiters (the green solid curve) performs strictly worse. In contrast, training with delimiters can achieve optimal performance under certain scenarios (see Figures 1a, 1b, and 1c). Theorem 5 in the next section provides a theoretical explanation for these observations, as well as the misalignment seen in Figure 1d. Further discussion and experimental details are provided in Section 5.2 and Appendix A.

### 5.2 LOSS LANDSCAPE OF 1-LAYER GLA

Given the input tokens with extended dimension, to ensure that GLA still implements WPGD as in Theorem 1, we propose the following model construction.

$$\tilde{\mathbf{W}}_k = \begin{bmatrix} \mathbf{W}_k & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}, \quad \tilde{\mathbf{W}}_q = \begin{bmatrix} \mathbf{W}_q & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \quad \text{and} \quad \tilde{\mathbf{W}}_v = \begin{bmatrix} \mathbf{W}_v & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}. \quad (21)$$

Here,  $\tilde{\mathbf{W}}_{k,q,v} \in \mathbb{R}^{(d+p+1) \times (d+p+1)}$  and  $\mathbf{W}_{k,q,v} \in \mathbb{R}^{(d+1) \times (d+1)}$  are constructed via (8). The main idea is to set the last  $p$  rows and columns of attention matrices to zeros, ensuring that the delimiters do not affect the final prediction.

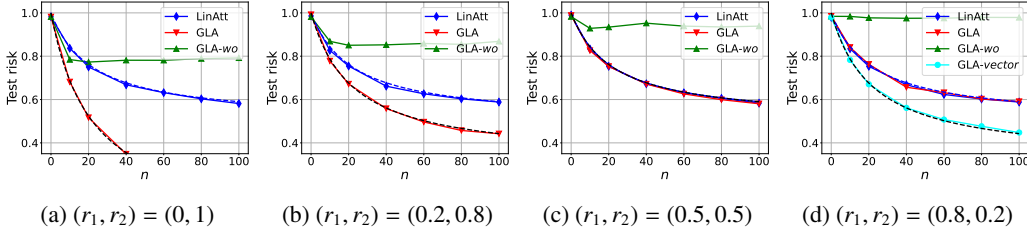


Figure 1: We consider four different types of model training: **LinAtt** (blue solid): Standard linear attention training. **GLA** (red solid): GLA training using prompts with delimiters (see (20)) and scalar gating. **GLA-wo** (green solid): GLA training using prompts without delimiters and with scalar gating. **GLA-vector** (cyan solid): GLA training using prompts with delimiters and vector gating. The blue and black dashed curves represent the optimal linear attention and WPGD risks from (25) and (18), respectively, as the number of in-context examples  $n$  increases. Implementation details are provided in Appendix A.

**Assumption B.** Delimiters  $\bar{\mathbf{d}}_0, \dots, \bar{\mathbf{d}}_K$  are linearly independent, and activation function  $\phi(z) : \mathbb{R} \rightarrow [0, 1]$  is continuous, satisfying  $\phi(-\infty) = 0$  and  $\phi(+\infty) = 1$ .

**Assumption C.** The correlation between context tasks  $(\beta_k)_{k=1}^K$  and query task  $\beta$  satisfies  $\mathbb{E}[\beta_i^\top \beta_j] = 0$  and  $\mathbb{E}[\beta_i^\top \beta] \leq \mathbb{E}[\beta_j^\top \beta]$  for  $1 \leq i \leq j \leq K$ .

Given context examples  $\{(X_k, \mathbf{y}_k) := (\mathbf{x}_i^{(k)}, \mathbf{y}_i^{(k)})_{i=1}^{n_k}\}_{k=1}^K$ , define the concatenated data  $(X, \mathbf{y})$  as follows:

$$X = [X_1^\top \dots X_K^\top]^\top \in \mathbb{R}^{n \times d} \quad \text{and} \quad \mathbf{y} = [\mathbf{y}_1^\top \dots \mathbf{y}_K^\top]^\top \in \mathbb{R}^n. \quad (22)$$

Based on the assumptions above, we are able to establish the equivalence between optimizing 1-layer GLA and optimizing 1-step WPGD predictor under scalar gating.

**Theorem 5 (Scalar Gating).** Recap the loss function  $\mathcal{L}_{\text{WPGD}}(\mathbf{P}, \omega)$  from (3) with dataset  $(X, \mathbf{y})$  defined in (22). Suppose Assumption B holds and consider GLA with scalar gating  $g(z) = \phi(\mathbf{w}_g^\top z) \mathbf{1} \mathbf{1}^\top$  where  $\mathbf{w}_g$  is the trainable parameter. Consider input prompt  $\mathbf{Z}$  defined in (20) and model constructions described in (21). Then the optimal risk  $\mathcal{L}_{\text{GLA}}^*$  defined in (19) obeys

$$\mathcal{L}_{\text{GLA}}^* = \mathcal{L}_{\text{WPGD}}^{\star, \mathcal{W}} \quad \text{where} \quad \mathcal{L}_{\text{WPGD}}^{\star, \mathcal{W}} := \min_{\mathbf{P} \in \mathbb{R}^{d \times d}, \omega \in \mathcal{W}} \mathcal{L}_{\text{WPGD}}(\mathbf{P}, \omega). \quad (23)$$

Here,  $\mathcal{W} := \left\{ [\omega_1 \mathbf{1}_{n_1}^\top \dots \omega_K \mathbf{1}_{n_K}^\top]^\top \in \mathbb{R}^n \mid 0 \leq \omega_i \leq \omega_j \leq 1, \forall 1 \leq i \leq j \leq K \right\}$ . Additionally, suppose Assumption C holds and  $n_i = n_j$ , for any  $i, j \in [K]$ . Let  $\mathcal{L}_{\text{WPGD}}^*$  be the optimal WPGD risk (c.f. (3)). Then  $\mathcal{L}_{\text{GLA}}^*$  satisfies

$$\mathcal{L}_{\text{GLA}}^* = \mathcal{L}_{\text{WPGD}}^*. \quad (24)$$

Assumption B ensures that any  $\omega$  in  $\mathcal{W}$  can be achieved by an appropriate choice of gating parameters. Furthermore, Assumption C guarantees that the optimal choice of  $\omega$  under the WPGD objective lies within the search space  $\mathcal{W}$ . The proof is provided in Appendix D.1.

In Figure 1, we conduct model training to validate our findings. Consider the setting where  $K = 2$  and let  $(r_1, r_2) = (\mathbb{E}[\beta_1^\top \beta]/d, \mathbb{E}[\beta_2^\top \beta]/d)$ . In Figures 1a, 1b, and 1c, Assumption C holds, and the GLA results (shown in solid red) align with the optimal WPGD risk (represented by the dashed black curves), validating (24). However, in Figure 1d, since  $r_1 > r_2$ , Assumption C does not hold, and as a result, the optimal GLA loss  $\mathcal{L}_{\text{GLA}}^*$  obtained from (23) is lower than the optimal WPGD loss  $\mathcal{L}_{\text{WPGD}}^*$ . Further experimental details are deferred to Appendix A.

**Loss landscape of vector gating.** Till now, much of our discussion has focused on the scalar gating setting. It is important to highlight that, even in the scalar-weighting context, analyzing the WPGD problem remains non-trivial due to the joint optimization over  $(\mathbf{P}, \omega)$ . However, as demonstrated in Theorem 5, scalar gating can only express weightings within the set  $\mathcal{W}$ . If Assumption C does not hold,  $\mathcal{L}_{\text{GLA}}^*$  cannot achieve the optimal WPGD loss (see the misalignment between red solid curve, presenting  $\mathcal{L}_{\text{GLA}}^*$ , and black dashed curve, presenting  $\mathcal{L}_{\text{WPGD}}^*$  in Figure 1d). We argue that vector gating overcomes this limitation by applying distinct weighting mechanisms across different dimensions, facilitating stronger expressivity.

**Theorem 6 (Vector Gating).** Recall input prompt  $\mathbf{Z}$  from (20) and model constructions from (21) but with  $\mathbf{W}_v = [\mathbf{0}_{(d+1) \times d} \ \mathbf{u}]$ . Suppose Assumption B holds and consider GLA with vector gating  $g(\mathbf{z}) = \phi(\mathbf{W}_g \mathbf{z}) \mathbf{1}^\top$ . Here,  $\mathbf{u}$  and  $\mathbf{W}_g$  are trainable parameters. Consider Problem (19), where we employ a vector gating  $g(\mathbf{z}) = \phi(\mathbf{W}_g \mathbf{z}) \mathbf{1}^\top$ . Let  $\mathcal{L}_{\text{GLA-v}}^*$  denote its optimal risk, and  $\mathcal{L}_{\text{WPGD}}^*$  be defined as in (3). Then, the optimal risk obeys  $\mathcal{L}_{\text{GLA-v}}^* = \mathcal{L}_{\text{WPGD}}^*$ .

In Theorem 5, the equivalence between  $\mathcal{L}_{\text{GLA}}^*$  and  $\mathcal{L}_{\text{WPGD}}^*$  is established only when both Assumptions B and C are satisfied. In contrast, Theorem 6 demonstrates that applying vector gating requires only Assumption B to establish  $\mathcal{L}_{\text{GLA-v}}^* = \mathcal{L}_{\text{WPGD}}^*$ . Specifically, under the bounded activation model of Assumption B, scalar gating is unable to express non-monotonic weighting schemes. For instance, suppose there are two tasks: Even if Task 1 is more relevant to the query, Assumption B will assign a higher weight to examples in Task 2 resulting in sub-optimal prediction. Theorem 6 shows that vector gating can avoid such bottlenecks by potentially encoding tasks in distinct subspaces. To verify these intuitions, in Figure 1d, we train a GLA model with vector gating and results are presented in cyan curve, which outperform the scalar gating results (red solid) and align with the optimal WPGD loss (black dashed).

**Loss landscape of 1-layer linear attention.** Inspired by the fact that linear attention implements all ones gating, that is,  $\mathbf{G}_i \equiv \mathbf{1}$ . Consider training a single-layer linear attention and let  $f_{\text{ATT}}(\mathbf{Z}) := f_{\text{GLA}}(\mathbf{Z}, \mathbf{G}_i \equiv \mathbf{1})$  be its prediction. Let  $\mathcal{L}_{\text{ATT}}^*$  be the corresponding optimal risk following (19).

**Corollary 3.** Consider a single-layer linear attention following model construction in (8) and consider linear data as given in Definition 2. Let  $\mathbf{R}, \mathbf{r}$  be the corresponding correlation matrix and vector as defined in Definition 1. Suppose  $\Sigma = \mathbf{I}$ . Then the optimal risk obeys

$$\mathcal{L}_{\text{ATT}}^* := \min_{\mathbf{P} \in \mathbb{R}^{d \times d}} \mathcal{L}_{\text{WPGD}}(\mathbf{P}, \omega = \mathbf{1}) = d + \sigma^2 - \frac{d(\mathbf{1}^\top \mathbf{r})^2}{n(d + \sigma^2 + 1) + \mathbf{1}^\top \mathbf{R} \mathbf{1}}. \quad (25)$$

**Corollary 4 (Benefit of Gating).** Consider the same setting as discussed in Corollary 3, and suppose Assumption B holds. Then, we have that  $\mathcal{L}_{\text{ATT}}^* \geq \mathcal{L}_{\text{GLA}}^*$ . Additionally, if Assumption C holds, we obtain

$$\mathcal{L}_{\text{ATT}}^* - \mathcal{L}_{\text{GLA}}^* = d \cdot \mathbf{r}^\top \left( \mathbf{R}_+^{-1} - \frac{\mathbf{1} \mathbf{1}^\top}{\mathbf{1}^\top \mathbf{R}_+ \mathbf{1}} \right) \mathbf{r} \geq 0, \quad \text{where} \quad \mathbf{R}_+ := \mathbf{R} + (d + \sigma^2 + 1) \mathbf{I}.$$

The proof of this corollary is directly from (18), (24) and (25). In the Figure 1, blue solid curves represent the linear attention results and blue dashed are the theory curves following (25). The two curves are aligned in all the subfigures, which validate our Corollary 3. More implementation details are deferred to Appendix A.

## 6 DISCUSSION

To summarize, this work offers a fresh theoretical perspective on gated linear attention models through in-context learning by showing that they can emulate data-dependent weighted preconditioned gradient descent (WPGD) algorithms. Our work also reveals how gating is crucial for achieving ICL with stronger data/context adaptivity by demonstrating clear separations between linear attention, scalar-valued gating, and vector-valued gating. We study the optimization landscape of GLA through a connection to the WPGD formulation (3). We have advocated that (3) is a problem of fundamental mathematical interest in its own right, developed the first characterization of its optimization landscape, and showed that it enjoys unique global minima and no other stationary point under mild conditions.

**Limitations and Future Work.** Our analysis is currently limited to characterizing the landscape of scalar gating in GLA models. Extending this framework to vector-valued gating and exploring when delimiters are necessary for learning, as well as investigating the GLA landscape where gates depend on input features, are promising directions for future research.

## REFERENCES

Kwangjun Ahn, Xiang Cheng, Hadi Daneshmand, and Suvrit Sra. Transformers learn to implement preconditioned gradient descent for in-context learning. *Advances in Neural Information Processing Systems*, 36, 2024.

- Ekin Akyürek, Dale Schuurmans, Jacob Andreas, Tengyu Ma, and Denny Zhou. What learning algorithm is in-context learning? investigations with linear models. In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=0g0X4H8yN4I>.
- Simran Arora, Sabri Eyuboglu, Michael Zhang, Aman Timalina, Silas Alberti, Dylan Zinsley, James Zou, Atri Rudra, and Christopher Ré. Simple linear attention language models balance the recall-throughput tradeoff. *arXiv preprint arXiv:2402.18668*, 2024.
- Akari Asai, Mohammadreza Salehi, Matthew E Peters, and Hannaneh Hajishirzi. Attempt: Parameter-efficient multi-task tuning via attentional mixtures of soft prompts. *arXiv preprint arXiv:2205.11961*, 2022.
- Maximilian Beck, Korbinian Pöppel, Markus Spanring, Andreas Auer, Oleksandra Prudnikova, Michael Kopp, Günter Klambauer, Johannes Brandstetter, and Sepp Hochreiter. xLSTM: Extended long short-term memory. *arXiv preprint arXiv:2405.04517*, 2024.
- Tom B Brown. Language models are few-shot learners. *arXiv preprint arXiv:2005.14165*, 2020.
- Liam Collins, Advait Parulekar, Aryan Mokhtari, Sujay Sanghavi, and Sanjay Shakkottai. In-context learning with transformers: Softmax attention adapts to function lipschitzness. *arXiv preprint arXiv:2402.11639*, 2024.
- Tri Dao and Albert Gu. Transformers are ssms: Generalized models and efficient algorithms through structured state space duality. *arXiv preprint arXiv:2405.21060*, 2024.
- Soham De, Samuel L Smith, Anushan Fernando, Aleksandar Botev, George Cristian-Muraru, Albert Gu, Ruba Haroun, Leonard Berrada, Yutian Chen, Srivatsan Srinivasan, et al. Griffin: Mixing gated linear recurrences with local attention for efficient language models. *arXiv preprint arXiv:2402.19427*, 2024.
- Nan Ding, Tomer Levinboim, Jialin Wu, Sebastian Goodman, and Radu Soricut. Causallm is not optimal for in-context learning. *arXiv preprint arXiv:2308.06912*, 2023.
- Chen Dun, Mirian Hipolito Garcia, Guoqing Zheng, Ahmed Hassan Awadallah, Anastasios Kyrillidis, and Robert Sim. Sweeping heterogeneity with smart mops: Mixture of prompts for llm task adaptation. *arXiv preprint arXiv:2310.02842*, 2023.
- Daniel Y Fu, Tri Dao, Khaled K Saab, Armin W Thomas, Atri Rudra, and Christopher Ré. Hungry hungry hippos: Towards language modeling with state space models. *arXiv preprint arXiv:2212.14052*, 2022.
- Deqing Fu, Tianqi Chen, Robin Jia, and Vatsal Sharan. Transformers learn higher-order optimization methods for in-context learning: A study with linear models. In *NeurIPS 2023 Workshop on Mathematics of Modern Machine Learning*.
- Shivam Garg, Dimitris Tsipras, Percy S Liang, and Gregory Valiant. What can transformers learn in-context? a case study of simple function classes. *Advances in Neural Information Processing Systems*, 35:30583–30598, 2022.
- Khashayar Gatmiry, Nikunj Saunshi, Sashank J Reddi, Stefanie Jegelka, and Sanjiv Kumar. Can looped transformers learn to implement multi-step gradient descent for in-context learning? In *Forty-first International Conference on Machine Learning*.
- Riccardo Grazi, Julien Siems, Simon Schrod, Thomas Brox, and Frank Hutter. Is mamba capable of in-context learning? *arXiv preprint arXiv:2402.03170*, 2024.
- Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces. *arXiv preprint arXiv:2312.00752*, 2023.
- Albert Gu, Karan Goel, and Christopher Ré. Efficiently modeling long sequences with structured state spaces. *arXiv preprint arXiv:2111.00396*, 2021.

- Samy Jelassi, David Brandfonbrener, Sham M Kakade, and Eran Malach. Repeat after me: Transformers are better than state space models at copying. *arXiv preprint arXiv:2402.01032*, 2024.
- Angelos Katharopoulos, Apoorv Vyas, Nikolaos Pappas, and François Fleuret. Transformers are rnns: Fast autoregressive transformers with linear attention. In *International conference on machine learning*, pp. 5156–5165. PMLR, 2020.
- Tobias Katsch. Gateloop: Fully data-controlled linear recurrence for sequence modeling. *arXiv preprint arXiv:2311.01927*, 2023.
- Yingcong Li, Muhammed Emrullah Ildiz, Dimitris Papailiopoulos, and Samet Oymak. Transformers as algorithms: Generalization and stability in in-context learning. In *International Conference on Machine Learning*, pp. 19565–19594. PMLR, 2023.
- Yingcong Li, Ankit Singh Rawat, and Samet Oymak. Fine-grained analysis of in-context linear estimation: Data, architecture, and beyond. *arXiv preprint arXiv:2407.10005*, 2024.
- Ziqian Lin and Kangwook Lee. Dual operating modes of in-context learning. *arXiv preprint arXiv:2402.18819*, 2024.
- Arvind Mahankali, Tatsunori B Hashimoto, and Tengyu Ma. One step of gradient descent is provably the optimal in-context learner with one layer of linear self-attention. *arXiv preprint arXiv:2307.03576*, 2023.
- Sewon Min, Xinxu Lyu, Ari Holtzman, Mikel Artetxe, Mike Lewis, Hannaneh Hajishirzi, and Luke Zettlemoyer. Rethinking the role of demonstrations: What makes in-context learning work? *arXiv preprint arXiv:2202.12837*, 2022.
- Antonio Orvieto, Samuel L Smith, Albert Gu, Anushan Fernando, Caglar Gulcehre, Razvan Pascanu, and Soham De. Resurrecting recurrent neural networks for long sequences. In *International Conference on Machine Learning*, pp. 26670–26698. PMLR, 2023.
- Jongho Park, Jaeseung Park, Zheyang Xiong, Nayoung Lee, Jaewoong Cho, Samet Oymak, Kangwook Lee, and Dimitris Papailiopoulos. Can mamba learn how to learn? a comparative study on in-context learning tasks. *arXiv preprint arXiv:2402.04248*, 2024.
- Bo Peng, Eric Alcaide, Quentin Anthony, Alon Albalak, Samuel Arcadinho, Stella Biderman, Huanqi Cao, Xin Cheng, Michael Chung, Matteo Grella, et al. Rwkv: Reinventing rnns for the transformer era. *arXiv preprint arXiv:2305.13048*, 2023.
- Bo Peng, Daniel Goldstein, Quentin Anthony, Alon Albalak, Eric Alcaide, Stella Biderman, Eugene Cheah, Teddy Ferdinan, Haowen Hou, Przemysław Kazienko, et al. Eagle and finch: RWKV with matrix-valued states and dynamic recurrence. *arXiv preprint arXiv:2404.05892*, 2024.
- Hao Peng, Nikolaos Pappas, Dani Yogatama, Roy Schwartz, Noah A Smith, and Lingpeng Kong. Random feature attention. *arXiv preprint arXiv:2103.02143*, 2021.
- Zhen Qin, Songlin Yang, Weixuan Sun, Xuyang Shen, Dong Li, Weigao Sun, and Yiran Zhong. Hgrn2: Gated linear rnns with state expansion. *arXiv preprint arXiv:2404.07904*, 2024.
- Jerome Sieber, Carmen Amo Alonso, Alexandre Didier, Melanie N Zeilinger, and Antonio Orvieto. Understanding the differences in foundation models: Attention, state space models, and recurrent neural networks. *arXiv preprint arXiv:2405.15731*, 2024.
- Yutao Sun, Li Dong, Shaohan Huang, Shuming Ma, Yuqing Xia, Jilong Xue, Jianyong Wang, and Furu Wei. Retentive network: A successor to transformer for large language models. *arXiv preprint arXiv:2307.08621*, 2023.
- Yutao Sun, Li Dong, Yi Zhu, Shaohan Huang, Wenhui Wang, Shuming Ma, Quanlu Zhang, Jianyong Wang, and Furu Wei. You only cache once: Decoder-decoder architectures for language models. *arXiv preprint arXiv:2405.05254*, 2024.
- A Vaswani. Attention is all you need. *Advances in Neural Information Processing Systems*, 2017.



- Johannes Von Oswald, Eyvind Niklasson, Ettore Randazzo, João Sacramento, Alexander Mordvintsev, Andrey Zhmoginov, and Max Vladymyrov. Transformers learn in-context by gradient descent. In *International Conference on Machine Learning*, pp. 35151–35174. PMLR, 2023.
- Johannes von Oswald, Eyvind Niklasson, Maximilian Schlegel, Seijin Kobayashi, Nicolas Zucchet, Nino Scherrer, Nolan Miller, Mark Sandler, Max Vladymyrov, Razvan Pascanu, et al. Uncovering mesa-optimization algorithms in transformers. *arXiv preprint arXiv:2309.05858*, 2023.
- Ruochen Wang, Sohyun An, Minhao Cheng, Tianyi Zhou, Sung Ju Hwang, and Cho-Jui Hsieh. One prompt is not enough: Automated construction of a mixture-of-expert prompts. *arXiv preprint arXiv:2407.00256*, 2024.
- Jingfeng Wu, Difan Zou, Zixiang Chen, Vladimir Braverman, Quanquan Gu, and Peter L Bartlett. How many pretraining tasks are needed for in-context learning of linear regression? *arXiv preprint arXiv:2310.08391*, 2023.
- Sang Michael Xie, Aditi Raghunathan, Percy Liang, and Tengyu Ma. An explanation of in-context learning as implicit bayesian inference. In *International Conference on Learning Representations*.
- Songlin Yang, Bailin Wang, Yikang Shen, Rameswar Panda, and Yoon Kim. Gated linear attention transformers with hardware-efficient training. *arXiv preprint arXiv:2312.06635*, 2023.
- Ruiqi Zhang, Spencer Frei, and Peter L Bartlett. Trained transformers learn linear models in-context. *arXiv preprint arXiv:2306.09927*, 2023.
- Ruiqi Zhang, Spencer Frei, and Peter L Bartlett. Trained transformers learn linear models in-context. *Journal of Machine Learning Research*, 25(49):1–55, 2024.
- Itamar Zimmerman, Ameen Ali, and Lior Wolf. A unified implicit attention formulation for gated-linear recurrent sequence models. *arXiv preprint arXiv:2405.16504*, 2024.
- Chang Zong, Jian Shao, Weiming Lu, and Yueting Zhuang. Stock movement prediction with multimodal stable fusion via gated cross-attention mechanism. *arXiv preprint arXiv:2406.06594*, 2024.

## CONTENTS

<b>A Implementation Detail</b>	<b>15</b>
A.1 Multi-layer Experiments	16
<b>B GLA <math>\Leftrightarrow</math> WPGD</b>	<b>16</b>
B.1 Proof of Theroem 1	16
B.2 Proof of Theorem 2	18
<b>C Optimization Landscape of WPGD</b>	<b>20</b>
C.1 Proof of Theorem 3	20
C.2 Proof of Theorem 4	23
C.3 Proof of Corollary 2	26
<b>D Loss Landscape of 1-layer GLA</b>	<b>27</b>
D.1 Proof of Theorem 5	27
D.2 Proof of Theorem 6	28
D.3 Proof of Corollary 3	28

## A IMPLEMENTATION DETAIL

**Data generation.** Consider ICL problem with input in the form of multi-task prompt as described in Section 5.1. In the experiments, we set  $K = 2$ , dimensions  $d = 10$  and  $p = 5$ , uniform context length  $n_1 = n_2 = \bar{n}$ , and vary  $\bar{n}$  from 0 to 50. Let  $(r_1, r_2) := (\mathbb{E}[\beta_1^\top \beta]/d, \mathbb{E}[\beta_2^\top \beta]/d)$  denote the correlations between in-context tasks  $\beta_1, \beta_2$  and query task  $\beta$ . We generate task vectors as follows:

$$\beta_1, \beta_2 \sim \mathcal{N}(0, \mathbf{I}_d) \quad \text{and} \quad \beta \sim \mathcal{N}(r_1 \beta_1 + r_2 \beta_2, (1 - r_1^2 - r_2^2) \mathbf{I}_d).$$

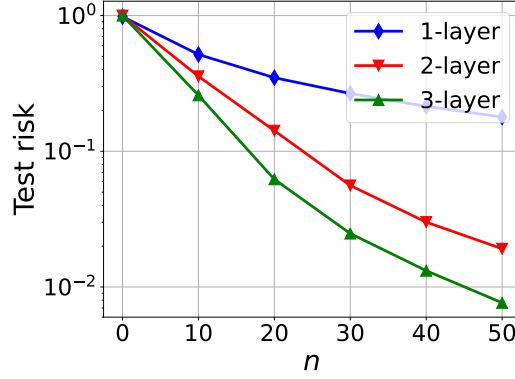
Input features are randomly sampled  $\mathbf{x}_i^{(k)} \sim \mathcal{N}(0, \mathbf{I}_d)$  and  $y_i^{(k)} = \beta^\top \mathbf{x}_i^{(k)}$  ( $\sigma = 0$ ),  $k \in \{1, 2\}$ . Additionally, delimiters  $\bar{\mathbf{d}}_0, \dots, \bar{\mathbf{d}}_K$  are randomly sampled from  $\mathcal{N}(0, \mathbf{I}_p)$ .

**Implementation setting.** We train 1-layer linear attention and GLA models for solving multi-prompt ICL problem as described in Section 5.1. For GLA model, we consider sigmoid-type gating function given by scalar gating:  $g(\mathbf{z}) = \phi(\mathbf{w}_g^\top \mathbf{z}) \mathbf{1} \mathbf{1}^\top$ , or vector gating:  $g(\mathbf{z}) = \phi(\mathbf{W}_g \mathbf{z}) \mathbf{1} \mathbf{1}^\top$  where  $\phi(\mathbf{z}) = (1 + e^{-\mathbf{z}})^{-1}$  is the activation function. Note that although the theoretical results are based on the model constructions (c.f. (8) and (21)), we do not restrict the attention weights in our implementation. We train each model for 10000 iterations with batch size 256 and Adam optimizer with learning rate  $10^{-3}$ . Similar to the previous work (Li et al., 2024), since our study focuses on the optimization landscape, ICL problems using linear attention/GLA models are non-convex, and experiments are implemented via gradient descent, we repeat 10 model trainings from different model initialization and data sampling (e.g., different choice of delimiters) and results are presented as the minimal test risk among those 10 trails. Results presented have been normalized by  $d$ .

**Experimental results.** Based on the experimental setting, we can obtain the correlation matrix and vector following Definition 1

$$\mathbf{R} = \begin{bmatrix} \mathbf{1}_n \mathbf{1}_n^\top & \mathbf{0} \\ \mathbf{0} & \mathbf{1}_n \mathbf{1}_n^\top \end{bmatrix} \quad \text{and} \quad \mathbf{r} = \begin{bmatrix} r_1 \mathbf{1}_n^\top & r_2 \mathbf{1}_n^\top \end{bmatrix}^\top.$$

Then dotted curves display our theoretical results derive using  $\Sigma = \mathbf{I}$  and  $\mathbf{R}, \mathbf{r}$  above. Specifically, in Figure 1, black dashed curves represent  $\mathcal{L}_{\text{WPGD}}^*$  following (18) and blues dashed curves represent  $\mathcal{L}_{\text{GLA}}^*$  following (25). We consider scenarios where  $(r_1, r_2) \in \{(0, 1), (0.2, 0.8), (0.5, 0.5), (0.8, 0.2)\}$  and results are presented in Figures (1a), (1b), (1c) and (1d), respectively.

Figure 2: Multi-layer GLA experiments with  $(r_1, r_2) = (0, 1)$ .

- GLA-wo achieves the worst performance among all the methods. We claim that it is due to the randomness of input tokens as discussed in Section 5.1. Thanks to the introduction of delimiters as described in (20), data and gating is decoupled and a task-dependent weighting is learnt. Hence, GLA is able to achieve comparable performance to the optimal one ( $\mathcal{L}_{\text{WPGD}}^*$ , red dashed). Note that GLA-wo performs even worse than LinAtt. It comes from the fact the weighting induced by GLA-wo varies over different input prompts and it can not implement all ones weight.
- The alignments between LinAtt (blue solid) and blue dashed curves validate our Corollary 3. In Figures 1a, 1b and 1c, the alignments between GLA (red solid) and  $\mathcal{L}_{\text{WPGD}}$  (black dashed) verify our Theorem 5, specifically, Equation 24. While in 1c and 1d, GLA achieves the same performance as LinAtt. It is due to the fact that GLA can not weight the history higher than its present. Then the equal-weighting, e.g.,  $\omega = 1$ , is the optimal weighting given such constraint. What’s more, the alignment between GLA-vector (cyan curves) and red dashed in Figure 1d validates our vector gating theorem in Theorem 6.

## A.1 MULTI-LAYER EXPERIMENTS

In this section, we present additional experiments on multi-layer GLA models. We adopt the same experimental setup as described in Figure 1a and Appendix A, with parameters set to  $(r_1, r_2) = (0, 1)$ . The results are displayed in Figure 2, where the blue, red, and green curves correspond to the performance of one-, two-, and three-layer GLA models, respectively, with the y-axis presented in log-scale. According to Theorem 2, an  $L$ -layer GLA performs  $L$  steps of WPGD, suggesting that deeper models should yield improved predictive performance. The experimental findings in Figure 2 align with the theoretical predictions of Theorem 2.

## B GLA $\Leftrightarrow$ WPGD

### B.1 PROOF OF THEROEM 1

Recap the problem settings from Section 2 where in-context samples are given by

$$\mathbf{Z} = [\mathbf{z}_1 \cdots \mathbf{z}_n \mathbf{z}_{n+1}]^\top = \begin{bmatrix} \mathbf{x}_1 & \cdots & \mathbf{x}_n & \mathbf{x}_{n+1} \\ \mathbf{y}_1 & \cdots & \mathbf{y}_n & 0 \end{bmatrix}^\top$$

and let the value, key and query embeddings at time  $i$  be

$$\mathbf{v}_i = \mathbf{W}_v \mathbf{z}_i, \quad \mathbf{k}_i = \mathbf{W}_k \mathbf{z}_i, \quad \text{and} \quad \mathbf{q}_i = \mathbf{W}_q \mathbf{z}_i.$$

Then we can rewrite the GLA output (c.f. (1)) as follows:

$$\begin{aligned} \mathbf{o}_i &= \mathbf{S}_i \mathbf{q}_i \quad \text{and} \quad \mathbf{S}_i = \mathbf{G}_i \odot \mathbf{S}_{i-1} + \mathbf{v}_i \mathbf{k}_i^\top \\ &= \sum_{j=1}^i \mathbf{G}_{j:i} \odot \mathbf{v}_j \mathbf{k}_j^\top \end{aligned}$$

where we define

$$\mathbf{G}_{j:i} = \mathbf{G}_{j+1} \odot \mathbf{G}_{j+2} \cdots \mathbf{G}_i, \quad j < i, \quad \text{and} \quad \mathbf{G}_{i:i} = \mathbf{1}\mathbf{1}^\top.$$

Consider the prediction based on the last token, then we obtain

$$\mathbf{o}_{n+1} = \mathbf{S}_{n+1} \mathbf{q}_{n+1} \quad \text{and} \quad \mathbf{S}_{n+1} = \sum_{j=1}^{n+1} \mathbf{G}_{j:n+1} \odot \mathbf{v}_j \mathbf{k}_j^\top.$$

**Construction 1:** Recall the model construction from (8) where

$$\mathbf{W}_k = \begin{bmatrix} \mathbf{P}_k^\top & \mathbf{0} \\ \mathbf{0}^\top & 0 \end{bmatrix}, \quad \mathbf{W}_q = \begin{bmatrix} \mathbf{P}_q^\top & \mathbf{0} \\ \mathbf{0}^\top & 0 \end{bmatrix} \quad \text{and} \quad \mathbf{W}_v = \begin{bmatrix} \mathbf{0}_{d \times d} & \mathbf{0} \\ \mathbf{0}^\top & 1 \end{bmatrix}. \quad (26)$$

Then, given each token  $\mathbf{z}_i = [\mathbf{x}_i^\top \ y_i]^\top$ ,  $i \in [n]$ , single-layer GLA returns

$$\mathbf{v}_i = \begin{bmatrix} \mathbf{0} \\ y_i \end{bmatrix}, \quad \mathbf{k}_i = \begin{bmatrix} \mathbf{P}_k^\top \mathbf{x}_i \\ 0 \end{bmatrix}, \quad \text{and} \quad \mathbf{q}_i = \begin{bmatrix} \mathbf{P}_q^\top \mathbf{x}_i \\ 0 \end{bmatrix},$$

and we obtain

$$\mathbf{v}_i \mathbf{k}_i^\top = \begin{bmatrix} \mathbf{0}_{d \times d} & \mathbf{0} \\ y_i \mathbf{x}_i^\top \mathbf{P}_k & 0 \end{bmatrix}, \quad i \leq n, \quad \text{and} \quad \mathbf{v}_{n+1} \mathbf{k}_{n+1}^\top = \mathbf{0}_{(d+1) \times (d+1)}.$$

Therefore, since only  $d$  entries in  $\mathbf{v}_i \mathbf{k}_i^\top$  matrix are nonzero, given  $\odot$  as the Hadamard product, only the corresponding  $d$  entries in all  $\mathbf{G}_i$  matrices are useful. Based on this observation, let

$$\mathbf{G}_i = \begin{bmatrix} * & * \\ \mathbf{g}_i^\top & * \end{bmatrix} \quad \text{and} \quad \mathbf{G}_{j:i} = \begin{bmatrix} * & * \\ \mathbf{g}_{j:i}^\top & * \end{bmatrix}$$

where  $\mathbf{g}_{j:i} = \mathbf{g}_{j+1} \odot \mathbf{g}_{j+2} \cdots \mathbf{g}_i \in \mathbb{R}^d$  for  $j < i$  and  $\mathbf{g}_{i:i} = \mathbf{1}_d$ .

Combing all together, and letting  $\mathbf{X} = [\mathbf{x}_1 \ \mathbf{x}_2 \ \cdots \ \mathbf{x}_n]^\top$  and  $\mathbf{y} = [y_1 \ y_2 \ \cdots \ y_n]^\top$ , we obtain

$$\mathbf{o}_{n+1} = \mathbf{S}_{n+1} \mathbf{q}_{n+1} = \begin{bmatrix} \mathbf{0}_{d \times d} & \mathbf{0} \\ \sum_{j=1}^n y_j \mathbf{x}_j^\top \mathbf{P}_k \odot \mathbf{g}_{j:n+1}^\top & 0 \end{bmatrix} \begin{bmatrix} \mathbf{P}_q^\top \mathbf{x} \\ 0 \end{bmatrix} = \begin{bmatrix} \mathbf{0} \\ \mathbf{x}^\top \mathbf{P}_q (\mathbf{X} \mathbf{P}_k \odot \mathbf{\Omega})^\top \mathbf{y} \end{bmatrix}$$

where

$$\mathbf{\Omega} = [\mathbf{g}_{1:n+1} \ \mathbf{g}_{2:n+1} \ \cdots \ \mathbf{g}_{n:n+1}] \in \mathbb{R}^{n \times d}.$$

Then if taking the last entry of  $\mathbf{o}_{n+1}$  as final prediction, we get

$$\hat{y} := \mathbf{o}_{n+1,d+1} = \mathbf{x}^\top \hat{\boldsymbol{\beta}} \quad \text{where} \quad \hat{\boldsymbol{\beta}} = \mathbf{P}_q (\mathbf{X} \mathbf{P}_k \odot \mathbf{\Omega})^\top \mathbf{y}.$$

It completes the proof of Theorem 1.

**Construction 2:** Based on the construction given in (26), only  $d$  elements of  $\mathbf{G}_i$  matrices are useful. One might ask about the effect of other entries of  $\mathbf{G}_i$ . Therefore, in the following, we introduce an other model construction showing that different row of  $\mathbf{G}_i$  implements WPGD with different weighting. Similarly, let  $\mathbf{W}_k, \mathbf{W}_q$  be the same as (26) but with  $\mathbf{W}_v$  constructed by

$$\mathbf{W}_v = [\mathbf{0}_{(d+1) \times d} \ \mathbf{u}] \quad \text{where} \quad \mathbf{u} = [u_1 \ u_2 \ \cdots \ u_{d+1}]^\top \in \mathbb{R}^{d+1}.$$

Then the value embeddings have the form of  $\mathbf{v}_i = y_i \mathbf{u}$ , which gives

$$\mathbf{v}_i \mathbf{k}_i^\top = \mathbf{u} [y_i \mathbf{x}_i^\top \mathbf{P}_k \ \mathbf{0}].$$

Next, let

$$\mathbf{G}_i = \begin{bmatrix} (\mathbf{g}_i^1)^\top & * \\ (\mathbf{g}_i^2)^\top & * \\ \vdots & \vdots \\ (\mathbf{g}_i^{d+1})^\top & * \end{bmatrix} \quad \text{and} \quad \mathbf{G}_{j:i} = \begin{bmatrix} (\mathbf{g}_{j:i}^1)^\top & * \\ (\mathbf{g}_{j:i}^2)^\top & * \\ \vdots & \vdots \\ (\mathbf{g}_{j:i}^{d+1})^\top & * \end{bmatrix}$$

where  $\mathbf{g}_i' \in \mathbb{R}^d$  corresponds to the  $i'$ -th row of  $\mathbf{G}_i$  and  $\mathbf{g}_{j:i}' = \mathbf{g}_{j+1}' \odot \mathbf{g}_{j+1}' \cdots \mathbf{g}_i'$ . Then we get the output

$$\mathbf{o}_{n+1} = \begin{bmatrix} \sum_{j=1}^n u_1 y_j \mathbf{x}_j^\top \mathbf{P}_k \odot (\mathbf{g}_{j:n+1}^1)^\top & 0 \\ \sum_{j=1}^n u_2 y_j \mathbf{x}_j^\top \mathbf{P}_k \odot (\mathbf{g}_{j:n+1}^2)^\top & 0 \\ \vdots & \vdots \\ \sum_{j=1}^n u_{d+1} y_j \mathbf{x}_j^\top \mathbf{P}_k \odot (\mathbf{g}_{j:n+1}^{d+1})^\top & 0 \end{bmatrix} \begin{bmatrix} \mathbf{P}_q^\top \mathbf{x} \\ 0 \end{bmatrix} = \begin{bmatrix} \mathbf{x}^\top \mathbf{P}_q (X \mathbf{P}_k \odot \mathbf{\Omega}_1)^\top \mathbf{y} \\ \mathbf{x}^\top \mathbf{P}_q (X \mathbf{P}_k \odot \mathbf{\Omega}_2)^\top \mathbf{y} \\ \vdots \\ \mathbf{x}^\top \mathbf{P}_q (X \mathbf{P}_k \odot \mathbf{\Omega}_{d+1})^\top \mathbf{y} \end{bmatrix}$$

where

$$\mathbf{\Omega}_i = u_i \begin{bmatrix} \mathbf{g}_{1:n+1}^i & \mathbf{g}_{2:n+1}^i & \cdots & \mathbf{g}_{n:n+1}^i \end{bmatrix} \in \mathbb{R}^{n \times d}, \quad i \leq d+1.$$

Therefore, consider  $(d+1)$ -dimensional output  $\mathbf{o}_{n+1}$ . Each entry implements a 1-step WPGD with same preconditioners  $\mathbf{P}_k, \mathbf{P}_q$  and different weighting matrices  $\mathbf{\Omega}$ 's. The weighting matrix of  $i$ 'th entry is determined by the  $i$ 'th row of all gating matrices. Note that if consider the last entry of  $\mathbf{o}_{n+1}$  as prediction, it returns the same result as Construction 1 above, where only last rows of  $\mathbf{G}_i$ 's are useful.

Additionally, suppose that the final prediction  $\hat{\mathbf{y}}$  is given after a linear head  $\mathbf{h}$ , that is,  $\hat{\mathbf{y}} = \mathbf{h}^\top \mathbf{o}_{n+1}$ , and let  $\mathbf{h} = [h_1 \ h_2 \ \cdots \ h_{d+1}]^\top \in \mathbb{R}^{d+1}$ . Then

$$\hat{\mathbf{y}} = \mathbf{h}^\top \mathbf{o}_{n+1} = \mathbf{x}^\top \mathbf{P}_q (X \mathbf{P}_k \odot \bar{\mathbf{\Omega}})^\top \mathbf{y} \quad (27)$$

where

$$\bar{\mathbf{\Omega}} = \sum_{i=1}^{d+1} h_i \mathbf{\Omega}_i = \sum_{i=1}^{d+1} h_i u_i \begin{bmatrix} \mathbf{g}_{1:n+1}^i & \mathbf{g}_{2:n+1}^i & \cdots & \mathbf{g}_{n:n+1}^i \end{bmatrix} \in \mathbb{R}^{n \times d}. \quad (28)$$

Then, single-layer GLA still returns 1-step WPGD with updated weighting matrix.

## B.2 PROOF OF THEOREM 2

**Theorem 7** (Extended version of Theorem 2). *Consider an  $L$ -layer GLA with  $\ell$ 'th layer parameterized by  $\mathbf{P}_{k,\ell}, \mathbf{P}_{q,\ell} \in \mathbb{R}^{d \times d}$  as in (8) and with corresponding gating vectors  $\mathbf{g}_i^\ell, i \in [n+1], \ell \in [L]$ . Let  $\hat{\mathbf{y}}_{i,\ell}$  be the  $(d+1)$ 'th entry of the  $i$ 'th token of the  $\ell$ 'th layer input (or  $(\ell-1)$ 'th layer output after residual). Additionally, denote  $\mathbf{\Omega}_\ell = [\mathbf{g}_{1:n+1}^\ell \ \cdots \ \mathbf{g}_{n:n+1}^\ell]^\top$  and  $\bar{X}_\ell = X \mathbf{P}_{k,\ell} \odot \mathbf{\Omega}_\ell$ . Let  $\mathbf{B}_\ell = [\beta_{1,\ell} \ \cdots \ \beta_{n,\ell}]^\top$  where  $\beta_{i,0} = \mathbf{0}$  for  $i \in [n+1]$  and  $\mathbf{M}_i = \begin{bmatrix} \mathbf{I}_i & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \in \mathbb{R}^{n \times n}$ . Then it satisfies that for*

- $i \leq n, \hat{\mathbf{y}}_{i,\ell} = \mathbf{y}_i - \mathbf{x}_i^\top \beta_{i,\ell-1}$  where  $\beta_{i,\ell} = \beta_{i,\ell-1} + \mathbf{P}_{q,\ell} (\nabla_{i,\ell} \odot \mathbf{g}_{i:n+1}^\ell)$
- and  $\hat{\mathbf{y}}_{n+1,\ell} = \mathbf{x}^\top \beta_{\ell-1}$  where  $\beta_\ell = (1 + \alpha_\ell) \beta_{\ell-1} + \mathbf{P}_{q,\ell} (\nabla_{n,\ell} \odot \mathbf{g}_{n+1}^\ell)$  and  $\alpha_\ell = \mathbf{x}^\top \mathbf{P}_{q,\ell} \mathbf{P}_{k,\ell}^\top \mathbf{x}$ .

Here, we define  $\nabla_{i,\ell} = \bar{X}_\ell^\top \mathbf{M}_i ((X \odot \mathbf{B}_{\ell-1}) \mathbf{1} - \mathbf{y})$ .

*Proof.* Recapping the model construction from (8) and following the same analysis in Appendix B.1, for  $i \leq n$ , we obtain

$$\mathbf{S}_i = \begin{bmatrix} \mathbf{0}_{d \times d} & \mathbf{0} \\ \sum_{j=1}^i y_j \mathbf{x}_j^\top \mathbf{P}_k \odot \mathbf{g}_{j:i}^\top & 0 \end{bmatrix}.$$

Additionally, recap that we have

$$\mathbf{M}_i = \begin{bmatrix} \mathbf{I}_i & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \quad \text{and} \quad \mathbf{\Omega} = [\mathbf{g}_{i:n+1} \ \cdots \ \mathbf{g}_{n:n+1}].$$

Let  $\odot$  denote Hadamard division. Then

$$\begin{aligned} \sum_{j=1}^i y_j \mathbf{P}_k^\top \mathbf{x}_j \odot \mathbf{g}_{j:i} &= \left( \sum_{j=1}^i y_j \mathbf{P}_k^\top \mathbf{x}_j \odot \mathbf{g}_{j:n+1} \right) \odot \mathbf{g}_{i:n+1} \\ &= (X \mathbf{P}_k \odot \mathbf{\Omega})^\top \mathbf{M}_i \mathbf{y} \odot \mathbf{g}_{i:n+1}, \end{aligned}$$



Therefore,

$$\mathbf{o}_i = \mathbf{S}_i \mathbf{q}_i = \begin{bmatrix} \mathbf{0}_{d \times d} & \mathbf{0} \\ ((X\mathbf{P}_k \odot \mathbf{\Omega})^\top \mathbf{M}_i \mathbf{y} \odot \mathbf{g}_{i:n+1})^\top & 0 \end{bmatrix} \begin{bmatrix} \mathbf{P}_q^\top \mathbf{x}_i \\ 0 \end{bmatrix} = \begin{bmatrix} \mathbf{x}_i^\top \mathbf{P}_q (\bar{\mathbf{X}}^\top \mathbf{M}_i \mathbf{y} \odot \mathbf{g}_{i:n+1}) \\ 0 \end{bmatrix}. \quad (29)$$

where we define  $\bar{\mathbf{X}} := X\mathbf{P}_k \odot \mathbf{\Omega}$ . Similarly, we can get the last token output

$$\mathbf{o}_{n+1} = \begin{bmatrix} \mathbf{0} \\ \mathbf{x}^\top \mathbf{P}_q (\bar{\mathbf{X}}^\top \mathbf{y} \odot \mathbf{g}_{n+1}) \end{bmatrix}. \quad (30)$$

Next, we consider the multi-layer GLA model. To begin with, let us define the input and output of  $\ell$ 'th layer as

$$\begin{aligned} \mathbf{Z}_\ell &= [\mathbf{z}_{1,\ell} \quad \cdots \quad \mathbf{z}_{n,\ell} \quad \mathbf{z}_{n+1,\ell}]^\top \in \mathbb{R}^{(n+1) \times (d+1)}, \\ \mathbf{O}_\ell &= [\mathbf{o}_{1,\ell} \quad \cdots \quad \mathbf{o}_{n,\ell} \quad \mathbf{o}_{n+1,\ell}]^\top \in \mathbb{R}^{(n+1) \times (d+1)}, \end{aligned}$$

where  $\mathbf{Z}_1 = \mathbf{Z}$ . Then, given the residual connection of each layer, the input of  $(\ell + 1)$ 'th layer is given by

$$\mathbf{Z}_{\ell+1} = \mathbf{Z}_\ell + \mathbf{O}_\ell. \quad (31)$$

Note that  $\mathbf{Z}_{\ell+1}$  is also the output of  $\ell$ 'th layer after residual. Recall (29) which implies that the first  $d$  dimension of the output  $\mathbf{o}_i$  for all tokens  $i \in [n + 1]$  is zero. Therefore, the first  $d$  dimension of  $\mathbf{z}_{i,\ell}$  keeps the same as  $\mathbf{x}_i$  and let us write the input of  $\ell$ 'th layer (also the output of  $(\ell - 1)$ 'th layer after residual) as

$$\mathbf{Z}_\ell = \begin{bmatrix} \mathbf{x}_1 & \cdots & \mathbf{x}_n & \mathbf{x} \\ \hat{\mathbf{y}}_{1,\ell} & \cdots & \hat{\mathbf{y}}_{n,\ell} & \hat{\mathbf{y}}_{n+1,\ell} \end{bmatrix}^\top, \quad (32)$$

and  $\hat{\mathbf{y}}_{i,1} = \mathbf{y}_i$  for  $i \in [n]$  and  $\hat{\mathbf{y}}_{n+1,1} = 0$ . Suppose that the  $\ell$ 'th layer is parameterized by  $(\mathbf{P}_{q,\ell}, \mathbf{P}_{k,\ell})$  and let  $\mathbf{Z}_\ell$  be its input. Additionally, suppose the gating matrices for  $\ell$ 'th layer,  $i$ 'th token is

$$\mathbf{G}_i^\ell = \begin{bmatrix} * & * \\ (\mathbf{g}_i^\ell)^\top & * \end{bmatrix}.$$

- We first study  $\hat{\mathbf{y}}_{i,\ell}$  for  $i \leq n$ . Following (29), we obtain the output at time  $i$

$$\mathbf{o}_{i,\ell} = \begin{bmatrix} \mathbf{0} \\ \mathbf{x}_i^\top \mathbf{P}_{q,\ell} (\bar{\mathbf{X}}_\ell^\top \mathbf{M}_i \hat{\mathbf{y}}_\ell \odot \mathbf{g}_{i:n+1}^\ell) \end{bmatrix}$$

where  $\bar{\mathbf{X}}_\ell := X\mathbf{P}_{k,\ell} \odot \mathbf{\Omega}_\ell$  and

$$\begin{aligned} \hat{\mathbf{y}}_\ell &= [\hat{\mathbf{y}}_{1,\ell} \quad \cdots \quad \hat{\mathbf{y}}_{n,\ell}]^\top \in \mathbb{R}^n, \\ \mathbf{\Omega}_\ell &= [\mathbf{g}_{1:n+1}^\ell \quad \mathbf{g}_{2:n+1}^\ell \quad \cdots \quad \mathbf{g}_{n:n+1}^\ell]^\top \in \mathbb{R}^{n \times d}. \end{aligned}$$

Following the residual connection as in (31), we have  $\mathbf{z}_{i,\ell+1} = \mathbf{z}_{i,\ell} + \mathbf{o}_{i,\ell}$  and hence

$$\hat{\mathbf{y}}_{i,\ell+1} = \hat{\mathbf{y}}_{i,\ell} + \mathbf{x}_i^\top \mathbf{P}_{q,\ell} (\bar{\mathbf{X}}_\ell^\top \mathbf{M}_i \hat{\mathbf{y}}_\ell \odot \mathbf{g}_{i:n+1}^\ell). \quad (33)$$

Now consider the algorithm given in the theorem statement where  $\hat{\mathbf{y}}_{i,\ell} = \mathbf{y}_i - \mathbf{x}_i^\top \boldsymbol{\beta}_{i,\ell-1}$  and  $\boldsymbol{\beta}_{i,\ell} = \boldsymbol{\beta}_{i,\ell-1} + \mathbf{P}_{q,\ell} (\bar{\mathbf{X}}_\ell^\top \mathbf{M}_i ((X \odot \mathbf{B}_{\ell-1}) \mathbf{1} - \mathbf{y}) \odot \mathbf{g}_{i:n+1}^\ell)$ , which gives

$$\begin{aligned} \hat{\mathbf{y}}_{i,\ell+1} &= \mathbf{y}_i - \mathbf{x}_i^\top \boldsymbol{\beta}_{i,\ell} \\ \hat{\mathbf{y}}_{i,\ell} &= \mathbf{y}_i - \mathbf{x}_i^\top \boldsymbol{\beta}_{i,\ell-1}. \end{aligned} \quad (34)$$

Then

$$\begin{aligned} \hat{\mathbf{y}}_{i,\ell+1} - \hat{\mathbf{y}}_{i,\ell} &= -\mathbf{x}_i^\top (\boldsymbol{\beta}_{i,\ell} - \boldsymbol{\beta}_{i,\ell-1}) \\ &= -\mathbf{x}_i^\top \mathbf{P}_{q,\ell} (\bar{\mathbf{X}}_\ell^\top \mathbf{M}_i ((X \odot \mathbf{B}_{\ell-1}) \mathbf{1} - \mathbf{y}) \odot \mathbf{g}_{i:n+1}^\ell) \\ &= \mathbf{x}_i^\top \mathbf{P}_{q,\ell} (\bar{\mathbf{X}}_\ell^\top \mathbf{M}_i \hat{\mathbf{y}}_\ell \odot \mathbf{g}_{i:n+1}^\ell). \end{aligned} \quad (35)$$

The last equation uses (34), that

$$(X \odot B_{\ell-1})\mathbf{1} = [\mathbf{x}_1^\top \beta_{1,\ell} \cdots \mathbf{x}_n^\top \beta_{n,\ell}]^\top \implies (X \odot B_{\ell-1})\mathbf{1} - \mathbf{y} = -\hat{\mathbf{y}}_\ell. \quad (36)$$

The equality between (33) and (35) completes the proof for  $i \in [n]$ .

• Next, we consider the last token output, that is  $i = n + 1$ . In the following, we remove the subscript  $n + 1$  from some notations for simplification.

Similarly, we get the  $(n + 1)$ 'th output of  $\ell$ 'th layer

$$\mathbf{o}_{n+1,\ell} = \begin{bmatrix} \mathbf{0} \\ \mathbf{x}^\top \mathbf{P}_{q,\ell} (\tilde{\mathbf{X}}_\ell^\top \hat{\mathbf{y}}_\ell \odot \mathbf{g}_{n+1}^\ell) \end{bmatrix} + \begin{bmatrix} \mathbf{0} \\ \mathbf{x}^\top \mathbf{P}_{q,\ell} \mathbf{P}_{k,\ell}^\top \mathbf{x} \cdot \hat{\mathbf{y}}_{n+1,\ell} \end{bmatrix}$$

where the second term comes from the fact that  $\hat{\mathbf{y}}_{n+1,\ell} \neq 0$  for  $\ell \neq 0$ .

Let  $\alpha_\ell := \mathbf{x}^\top \mathbf{P}_{q,\ell} \mathbf{P}_{k,\ell}^\top \mathbf{x}$ . Given  $\mathbf{Z}_{\ell+1} = \mathbf{Z}_\ell + \mathbf{O}_\ell$ , we obtain

$$\begin{aligned} \hat{\mathbf{y}}_{n+1,\ell+1} &= \hat{\mathbf{y}}_{n+1,\ell} + \mathbf{x}^\top \mathbf{P}_{q,\ell} (\tilde{\mathbf{X}}_\ell^\top \hat{\mathbf{y}}_\ell \odot \mathbf{g}_{n+1}^\ell) + \alpha_\ell \cdot \hat{\mathbf{y}}_{n+1,\ell} \\ &= (1 + \alpha_\ell) \hat{\mathbf{y}}_{n+1,\ell} + \mathbf{x}^\top \mathbf{P}_{q,\ell} (\tilde{\mathbf{X}}_\ell^\top \hat{\mathbf{y}}_\ell \odot \mathbf{g}_{n+1}^\ell). \end{aligned} \quad (37)$$

Now, consider the algorithm given in the theorem statement where  $\hat{\mathbf{y}}_{n+1,\ell} = -\mathbf{x}^\top \beta_{n+1,\ell-1}$  and  $\beta_{n+1,\ell} = (1 + \alpha_\ell) \beta_{n+1,\ell-1} + \mathbf{P}_{q,\ell} (\tilde{\mathbf{X}}_\ell^\top ((X \odot B_{\ell-1})\mathbf{1} - \mathbf{y}) \odot \mathbf{g}_{n+1}^\ell)$ , which gives

$$\begin{aligned} \hat{\mathbf{y}}_{n+1,\ell+1} &= -\mathbf{x}^\top \beta_{n+1,\ell} \\ \hat{\mathbf{y}}_{n+1,\ell} &= -\mathbf{x}^\top \beta_{n+1,\ell-1}. \end{aligned}$$

Then

$$\begin{aligned} \hat{\mathbf{y}}_{n+1,\ell+1} - (1 + \alpha_\ell) \hat{\mathbf{y}}_{n+1,\ell} &= -\mathbf{x}^\top (\beta_{n+1,\ell} - (1 + \alpha_\ell) \beta_{n+1,\ell-1}) \\ &= -\mathbf{x}^\top \mathbf{P}_{q,\ell} (\tilde{\mathbf{X}}_\ell^\top ((X \odot B_{\ell-1})\mathbf{1} - \mathbf{y}) \odot \mathbf{g}_{n+1}^\ell) \\ &= \mathbf{x}^\top \mathbf{P}_{q,\ell} (\tilde{\mathbf{X}}_\ell^\top \hat{\mathbf{y}}_\ell \odot \mathbf{g}_{n+1}^\ell) \end{aligned}$$

which is the same as (37) by using the fact from (36).  $\square$

## C OPTIMIZATION LANDSCAPE OF WPGD

### C.1 PROOF OF THEOREM 3

*Proof.* Recapping the objective from (3) and following Definition 2, we have

$$\begin{aligned} \mathcal{L}(\mathbf{P}, \omega) &= \mathbb{E} \left[ \left( \mathbf{y} - \mathbf{x}^\top \mathbf{P} \mathbf{X}(\omega \odot \mathbf{y}) \right)^2 \right] \\ &= \mathbb{E} [\mathbf{y}^2] - 2\mathbb{E} [\mathbf{y} \mathbf{x}^\top \mathbf{P} \mathbf{X}(\omega \odot \mathbf{y})] + \mathbb{E} \left[ \left( \mathbf{x}^\top \mathbf{P} \mathbf{X}(\omega \odot \mathbf{y}) \right)^2 \right]. \end{aligned}$$

Let  $\mathbf{y} = \mathbf{x}^\top \beta + \xi$  and  $y_i = \mathbf{x}_i^\top \beta_i + \xi_i$ , for  $i \in [n]$ , where  $\xi, \xi_i \sim \mathcal{N}(0, \sigma^2)$  are i.i.d. Then,

$$\mathbb{E}[\mathbf{y}^2] = \mathbb{E}[(\mathbf{x}^\top \beta + \xi)^2] = \text{tr}(\Sigma) + \sigma^2,$$

and

$$\begin{aligned} \mathbb{E} [\mathbf{y} \mathbf{x}^\top \mathbf{P} \mathbf{X}(\omega \odot \mathbf{y})] &= \mathbb{E} \left[ (\beta^\top \mathbf{x} + \xi) \mathbf{x}^\top \mathbf{P} \sum_{i=1}^n \omega_i \mathbf{x}_i (\mathbf{x}_i^\top \beta_i + \xi_i) \right] \\ &= \mathbb{E} \left[ \beta^\top \mathbf{x} \mathbf{x}^\top \mathbf{P} \sum_{i=1}^n \omega_i \mathbf{x}_i \mathbf{x}_i^\top \beta_i \right] \\ &= \text{tr} \left( \Sigma \mathbf{P} \Sigma \sum_{i=1}^n \omega_i \mathbb{E} [\beta_i \beta_i^\top] \right) \\ &= \text{tr} (\Sigma^2 \mathbf{P}) \omega^\top \mathbf{r}. \end{aligned}$$

Here, the last equality comes from the fact that since  $\beta_i - r_{ij}\beta_j$  is independent of  $\beta_j$  for  $i, j \in [n+1]$  following Definition 1, we have  $\mathbb{E}[\beta_i\beta_j^\top] = r_{i,n+1}\mathbf{I}_d$  and  $\sum_{i=1}^n \omega_i \mathbb{E}[\beta_i\beta_j^\top]$  returns  $\omega^\top \mathbf{r} \cdot \mathbf{I}_d$ .

Hence,

$$\begin{aligned} \mathbb{E}\left[(\mathbf{x}^\top \mathbf{P} \mathbf{X}(\omega \odot \mathbf{y}))^2\right] &= \mathbb{E}\left[\mathbf{x}^\top \mathbf{P} \left(\sum_{i=1}^n \omega_i (\mathbf{x}_i^\top \beta_i + \xi_i) \mathbf{x}_i\right) \left(\sum_{i=1}^n \omega_i \mathbf{x}_i^\top (\mathbf{x}_i^\top \beta_i + \xi_i)\right) \mathbf{P}^\top \mathbf{x}\right] \\ &= \text{tr}\left(\mathbf{P}^\top \Sigma \mathbf{P} \mathbb{E}\left[\sum_{i=1}^n \omega_i^2 (\mathbf{x}_i^\top \beta_i + \xi_i)^2 \mathbf{x}_i \mathbf{x}_i^\top\right]\right) \\ &\quad + \text{tr}\left(\mathbf{P}^\top \Sigma \mathbf{P} \mathbb{E}\left[\sum_{i \neq j} \omega_i \omega_j (\mathbf{x}_i^\top \beta_i + \xi_i) \mathbf{x}_i \mathbf{x}_j^\top (\mathbf{x}_j^\top \beta_j + \xi_j)\right]\right), \end{aligned}$$

where

$$\begin{aligned} \text{tr}\left(\mathbf{P}^\top \Sigma \mathbf{P} \mathbb{E}\left[\sum_{i=1}^n \omega_i^2 (\mathbf{x}_i^\top \beta_i + \xi_i)^2 \mathbf{x}_i \mathbf{x}_i^\top\right]\right) &= \text{tr}\left(\mathbf{P}^\top \Sigma \mathbf{P} \mathbb{E}\left[\sum_{i=1}^n \omega_i^2 (\mathbf{x}_i^\top \beta_i \beta_i^\top \mathbf{x}_i + \sigma^2) \mathbf{x}_i \mathbf{x}_i^\top\right]\right) \\ &= \|\omega\|_{\ell_2}^2 \text{tr}\left(\mathbf{P}^\top \Sigma \mathbf{P} (\mathbb{E}[\mathbf{x} \mathbf{x}^\top] + \sigma^2 \mathbf{I})\right) \\ &= \|\omega\|_{\ell_2}^2 (\text{tr}(\Sigma \mathbf{P}^\top \Sigma \mathbf{P}) (\text{tr}(\Sigma) + \sigma^2) + \text{tr}(\Sigma^2 \mathbf{P}^\top \Sigma \mathbf{P})), \end{aligned}$$

and

$$\begin{aligned} \text{tr}\left(\mathbf{P}^\top \Sigma \mathbf{P} \mathbb{E}\left[\sum_{i \neq j} \omega_i \omega_j (\mathbf{x}_i^\top \beta_i + \xi_i) \mathbf{x}_i \mathbf{x}_j^\top (\mathbf{x}_j^\top \beta_j + \xi_j)\right]\right) &= \text{tr}\left(\mathbf{P}^\top \Sigma \mathbf{P} \mathbb{E}\left[\sum_{i \neq j} \omega_i \omega_j \mathbf{x}_i \mathbf{x}_i^\top \beta_i \beta_j^\top \mathbf{x}_j \mathbf{x}_j^\top\right]\right) \\ &= \text{tr}(\Sigma^2 \mathbf{P}^\top \Sigma \mathbf{P}) \omega^\top \mathbf{R} \omega. \end{aligned}$$

Combining all together and letting  $M := \text{tr}(\Sigma) + \sigma^2$ , we obtain

$$\begin{aligned} \mathcal{L}(\mathbf{P}, \omega) &= M - 2\text{tr}(\Sigma^2 \mathbf{P}) \omega^\top \mathbf{r} \\ &\quad + M \|\omega\|_{\ell_2}^2 \text{tr}(\Sigma \mathbf{P}^\top \Sigma \mathbf{P}) + (\|\omega\|_{\ell_2}^2 + \omega^\top \mathbf{R} \omega) \text{tr}(\Sigma^2 \mathbf{P}^\top \Sigma \mathbf{P}). \end{aligned} \quad (38)$$

For simplicity, and without loss of generality, let

$$\tilde{\mathbf{P}} = \sqrt{\Sigma} \mathbf{P} \sqrt{\Sigma}. \quad (39)$$

Then, we obtain

$$\begin{aligned} \mathcal{L}(\tilde{\mathbf{P}}, \omega) &= M - 2\text{tr}(\Sigma \tilde{\mathbf{P}}) \omega^\top \mathbf{r} \\ &\quad + M \|\omega\|_{\ell_2}^2 \text{tr}(\tilde{\mathbf{P}}^\top \tilde{\mathbf{P}}) + (\|\omega\|_{\ell_2}^2 + \omega^\top \mathbf{R} \omega) \text{tr}(\Sigma \tilde{\mathbf{P}}^\top \tilde{\mathbf{P}}). \end{aligned} \quad (40)$$

Further, the gradients can be written as

$$\nabla_{\tilde{\mathbf{P}}} \mathcal{L}(\tilde{\mathbf{P}}, \omega) = -2\omega^\top \mathbf{r} \Sigma + 2M \|\omega\|_{\ell_2}^2 \tilde{\mathbf{P}} + 2(\|\omega\|_{\ell_2}^2 + \omega^\top \mathbf{R} \omega) \Sigma \tilde{\mathbf{P}}, \quad (41)$$

$$\nabla_{\omega} \mathcal{L}(\tilde{\mathbf{P}}, \omega) = -2\text{tr}(\Sigma \tilde{\mathbf{P}}) \mathbf{r} + 2M \text{tr}(\tilde{\mathbf{P}}^\top \tilde{\mathbf{P}}) \omega + 2\text{tr}(\Sigma \tilde{\mathbf{P}}^\top \tilde{\mathbf{P}}) (\mathbf{I}_n + \mathbf{R}) \omega. \quad (42)$$

Using the first-order optimality condition, and setting  $\nabla_{\tilde{\mathbf{P}}} \mathcal{L}(\tilde{\mathbf{P}}, \omega) = 0$  and  $\nabla_{\omega} \mathcal{L}(\tilde{\mathbf{P}}, \omega) = 0$ , we obtain

$$\begin{aligned} \tilde{\mathbf{P}} &= \left(M \|\omega\|_{\ell_2}^2 \mathbf{I} + (\|\omega\|_{\ell_2}^2 + \omega^\top \mathbf{R} \omega) \Sigma\right)^{-1} \Sigma \omega^\top \mathbf{r} \\ &= \frac{\omega^\top \mathbf{r}}{M \|\omega\|_{\ell_2}^2} \left(\frac{\|\omega\|_{\ell_2}^2 + \omega^\top \mathbf{R} \omega}{M \|\omega\|_{\ell_2}^2} \mathbf{I} + \Sigma^{-1}\right)^{-1} \\ &= \frac{\omega^\top \mathbf{r}}{M \|\omega\|_{\ell_2}^2} \left(\frac{\gamma + 1}{M} \cdot \mathbf{I} + \Sigma^{-1}\right)^{-1}, \end{aligned} \quad (43a)$$

where  $\gamma = \omega^\top \mathbf{R} \omega / \|\omega\|_{\ell_2}^2$ .

Further,

$$\begin{aligned}\omega &= \left( (M \text{tr}(\tilde{P}^\top \tilde{P}) + \text{tr}(\Sigma \tilde{P}^\top \tilde{P})) I + \text{tr}(\Sigma \tilde{P}^\top \tilde{P}) R \right)^{-1} \text{tr}(\Sigma \tilde{P}) r \\ &= \frac{\text{tr}(\Sigma \tilde{P})}{M \text{tr}(\tilde{P}^\top \tilde{P}) + \text{tr}(\Sigma \tilde{P}^\top \tilde{P})} \left( I + \frac{\text{tr}(\Sigma \tilde{P}^\top \tilde{P})}{M \text{tr}(\tilde{P}^\top \tilde{P}) + \text{tr}(\Sigma \tilde{P}^\top \tilde{P})} R \right)^{-1} r.\end{aligned}\quad (43b)$$

Let

$$\Sigma_\gamma := \frac{\gamma + 1}{M} \cdot I + \Sigma^{-1}.$$

Then, we get

$$\begin{aligned}\frac{\text{tr}(\Sigma \tilde{P}^\top \tilde{P})}{M \text{tr}(\tilde{P}^\top \tilde{P}) + \text{tr}(\Sigma \tilde{P}^\top \tilde{P})} &= \left( 1 + M \frac{\text{tr}(\tilde{P}^\top \tilde{P})}{\text{tr}(\Sigma \tilde{P}^\top \tilde{P})} \right)^{-1} \\ &= \left( 1 + M \frac{\text{tr}(\Sigma_\gamma^{-2})}{\text{tr}(\Sigma \Sigma_\gamma^{-2})} \right)^{-1} \\ &= \left( 1 + M \sum_{i=1}^d \frac{s_i^2}{(M + (\gamma + 1)s_i)^2} \left( \sum_{i=1}^d \frac{s_i^3}{(M + (\gamma + 1)s_i)^2} \right)^{-1} \right)^{-1} \\ &=: g(\gamma).\end{aligned}$$

Here, the last equality follows from eigen decomposition  $\Sigma = U \text{diag}(s) U^\top$  with  $s = [s_1, \dots, s_d]^\top \in \mathbb{R}_{++}^d$ .

Now, plugging  $\tilde{P}$  defined in (43a) within  $\omega$  given in (43b), we obtain

$$\omega = \frac{\text{tr}(\Sigma \tilde{P})}{M \text{tr}(\tilde{P}^\top \tilde{P}) + \text{tr}(\Sigma \tilde{P}^\top \tilde{P})} \cdot (g(\gamma) \cdot R + I)^{-1} r. \quad (44)$$

Using the above formulae for  $\omega$ , we rewrite  $\gamma = \omega^\top R \omega / \|\omega\|_{\ell_2}^2$  as

$$\begin{aligned}\gamma &= \frac{r^\top (g(\gamma) R + I)^{-1} R (g(\gamma) R + I)^{-1} r}{r^\top (g(\gamma) R + I)^{-2} r} \\ &= \sum_{i=1}^n \frac{\lambda_i a_i^2}{(1 + g(\gamma) \lambda_i)^2} \left( \sum_{i=1}^n \frac{a_i^2}{(1 + g(\gamma) \lambda_i)^2} \right)^{-1} \\ &=: h(g(\gamma)),\end{aligned}\quad (45)$$

where the second equality follows from Assumption A and the fact that  $R = E \text{diag}(\lambda) E^\top$  denotes the eigen decomposition of  $R$ , with  $\lambda = [\lambda_1, \dots, \lambda_n]^\top \in \mathbb{R}_+^n$ .

Now, let  $\gamma^*$  denote a fixed point of composite function  $h(g(\gamma))$ . From (43a) and (44), we obtain

$$\begin{aligned}\tilde{P} &= C(r, \omega, \Sigma) \cdot \left( \frac{\gamma^* + 1}{M} \cdot I + \Sigma^{-1} \right)^{-1}, \quad \text{and} \\ \omega &= c(r, \omega, \Sigma) \cdot (g(\gamma^*) \cdot R + I)^{-1} r.\end{aligned}\quad (46)$$

for some  $C(r, \omega, \Sigma) = \frac{\omega^\top r}{M \|\omega\|_{\ell_2}^2}$  and  $c(r, \omega, \Sigma) = \frac{\text{tr}(\Sigma \tilde{P})}{M \text{tr}(\tilde{P}^\top \tilde{P}) + \text{tr}(\Sigma \tilde{P}^\top \tilde{P})}$ .

Now, using the our definition  $\tilde{P} = \sqrt{\Sigma} P \sqrt{\Sigma}$ , we obtain

$$\begin{aligned}P(\gamma) &= C(r, \omega, \Sigma) \cdot \Sigma^{-\frac{1}{2}} \left( \frac{\gamma^* + 1}{\sigma^2 + \text{tr}(\Sigma)} \cdot \Sigma + I \right)^{-1} \Sigma^{-\frac{1}{2}}, \quad \text{and} \\ \omega(\gamma) &= c(r, \omega, \Sigma) \cdot (g(\gamma^*) \cdot R + I)^{-1} r.\end{aligned}$$

This completes the proof.  $\square$

## C.2 PROOF OF THEOREM 4

We first provide the following Lemma.

**Lemma 1.** *Let the functions  $h : \mathbb{R}_+ \rightarrow \mathbb{R}_+$  and  $g : \mathbb{R}_+ \rightarrow \mathbb{R}_+$  be defined as*

$$h(\bar{\gamma}) = \sum_{i=1}^n \frac{\lambda_i a_i^2}{(1 + \bar{\gamma} \lambda_i)^2} \left( \sum_{i=1}^n \frac{a_i^2}{(1 + \bar{\gamma} \lambda_i)^2} \right)^{-1}, \quad (47)$$

$$g(\gamma) = \left( 1 + M \sum_{i=1}^d \frac{s_i^2}{(M + (\gamma + 1)s_i)^2} \left( \sum_{i=1}^d \frac{s_i^3}{(M + (\gamma + 1)s_i)^2} \right)^{-1} \right)^{-1}, \quad (48)$$

where  $M = \sigma^2 + \sum_{i=1}^d s_i$ .

Suppose  $\Delta_{\Sigma} \cdot \Delta_{\mathbf{R}} < M + s_{\min}$ , where  $\Delta_{\Sigma}$  and  $\Delta_{\mathbf{R}}$  denote the effective spectral gaps of  $\Sigma$  and  $\mathbf{R}$ , respectively, as given in (12); and  $s_{\min}$  is the smallest eigenvalue of  $\Sigma$ . We have that

$$\left| \frac{\partial g}{\partial \gamma} \cdot \frac{\partial h}{\partial g} \right| \leq \frac{\Delta_{\Sigma}^2 \cdot \Delta_{\mathbf{R}}^2}{(M + s_{\min})^2} < 1.$$

*Proof.* Let

$$B(\gamma) = \sum_{i=1}^d \frac{s_i^3}{(M + (\gamma + 1)s_i)^2}, \quad C(\gamma) = \sum_{i=1}^d \frac{s_i^2}{(M + (\gamma + 1)s_i)^2}, \quad A(\gamma) = 1 + M \frac{C(\gamma)}{B(\gamma)}.$$

The derivatives of  $B(\gamma)$  and  $C(\gamma)$  are

$$B'(\gamma) = -2 \sum_{i=1}^d \frac{s_i^4}{(M + (\gamma + 1)s_i)^3}, \quad C'(\gamma) = -2 \sum_{i=1}^d \frac{s_i^3}{(M + (\gamma + 1)s_i)^3}.$$

The gradient of  $g(\gamma)$  is

$$\frac{\partial g}{\partial \gamma} = -M \left( \frac{1}{A(\gamma)B(\gamma)} \right)^2 (C'(\gamma)B(\gamma) - C(\gamma)B'(\gamma)). \quad (49)$$

It can be seen that

$$\left( \frac{1}{A(\gamma)} \right)^2 \leq M^{-2} \left( \sum_{i=1}^d \frac{s_i^3}{(M + (\gamma + 1)s_i)^2} \right)^2 \left( \sum_{i=1}^d \frac{s_i^2}{(M + (\gamma + 1)s_i)^2} \right)^{-2},$$

which implies that

$$\left( \frac{1}{A(\gamma)B(\gamma)} \right)^2 \leq M^{-2} \left( \sum_{i=1}^d \frac{s_i^2}{(M + (\gamma + 1)s_i)^2} \right)^{-2}. \quad (50a)$$

Further, we have

$$\begin{aligned} C'(\gamma)B(\gamma) - C(\gamma)B'(\gamma) &= -2 \sum_{i=1}^d \frac{s_i^3}{(M + (\gamma + 1)s_i)^3} \sum_{i=1}^d \frac{s_i^3}{(M + (\gamma + 1)s_i)^2} \\ &\quad + 2 \sum_{i=1}^d \frac{s_i^2}{(M + (\gamma + 1)s_i)^2} \sum_{i=1}^d \frac{s_i^4}{(M + (\gamma + 1)s_i)^3} \\ &= \sum_{i=1}^d \sum_{j=1}^d \frac{T_{ij}}{(M + (\gamma + 1)s_i)^3 (M + (\gamma + 1)s_j)^3} \\ &= M \cdot \sum_{i=1}^d \sum_{j=1}^d \frac{s_i^2 s_j^2 (s_i - s_j)^2}{(M + (\gamma + 1)s_i)^3 (M + (\gamma + 1)s_j)^3}, \end{aligned} \quad (50b)$$



where

$$\begin{aligned} T_{ij} &= s_i^2(M + (\gamma + 1)s_i)s_j^4 + s_i^4s_j^2(M + (\gamma + 1)s_j) \\ &\quad - s_i^3s_j^3(M + (\gamma + 1)s_j) - s_i^3(M + (\gamma + 1)s_i)s_j^3 \\ &= s_i^2s_j^2\left(M \cdot (s_j^2 + s_i^2 - 2s_is_j) + (\gamma + 1)(s_is_j^2 + s_i^2s_j - s_is_j^2 - s_i^2s_j)\right). \end{aligned} \quad (50c)$$

Thus, substituting (50a) and (50b) into (49), we obtain

$$\left| \frac{\partial g}{\partial \gamma} \right| \leq M \cdot M^{-1} \left( \sum_{i=1}^d \frac{s_i^2}{(M + (\gamma + 1)s_i)} \right)^{-2} \sum_{i,j=1}^d \frac{s_i^2s_j^2(s_i - s_j)^2}{(M + (\gamma + 1)s_i)^3(M + (\gamma + 1)s_j)^3}. \quad (51)$$

Next, we derive  $\frac{\partial h}{\partial \bar{\gamma}}$ . Let

$$D(\bar{\gamma}) = \sum_{i=1}^n \frac{\lambda_i a_i^2}{(1 + \bar{\gamma}\lambda_i)^2}, \quad E(\bar{\gamma}) = \sum_{i=1}^n \frac{a_i^2}{(1 + \bar{\gamma}\lambda_i)^2}.$$

We have

$$D'(\bar{\gamma}) = -2 \sum_{i=1}^n \frac{\lambda_i^2 a_i^2}{(1 + \bar{\gamma}\lambda_i)^3}, \quad E'(\bar{\gamma}) = -2 \sum_{i=1}^n \frac{\lambda_i a_i^2}{(1 + \bar{\gamma}\lambda_i)^3}.$$

The derivative of  $h$  with respect to  $\bar{\gamma}$  is given by

$$\frac{\partial h}{\partial \bar{\gamma}} = - \left( \frac{1}{E(\bar{\gamma})} \right)^2 (E(\bar{\gamma})D'(\bar{\gamma}) - D(\bar{\gamma})E'(\bar{\gamma})). \quad (52)$$

Substituting into (52), we get

$$\left( \frac{1}{E(\bar{\gamma})} \right)^2 = \left( \sum_{i=1}^n \frac{a_i^2}{(1 + \bar{\gamma}\lambda_i)^2} \right)^{-2}, \quad (53a)$$

and

$$\begin{aligned} E(\bar{\gamma})D'(\bar{\gamma}) - D(\bar{\gamma})E'(\bar{\gamma}) &= 2 \sum_{i=1}^n \frac{\lambda_i^2 a_i^2}{(1 + \bar{\gamma}\lambda_i)^3} \sum_{i=1}^n \frac{a_i^2}{(1 + \bar{\gamma}\lambda_i)^2} \\ &\quad - 2 \sum_{i=1}^n \frac{\lambda_i a_i^2}{(1 + \bar{\gamma}\lambda_i)^2} \sum_{i=1}^n \frac{a_i^2 \lambda_i}{(1 + \bar{\gamma}\lambda_i)^3} \\ &= \sum_{i=1}^n \sum_{j=1}^n \frac{\bar{T}_{ij}}{(1 + \bar{\gamma}\lambda_i)^3(1 + \bar{\gamma}\lambda_j)^3} \\ &= \sum_{i=1}^n \sum_{j=1}^n \frac{a_i^2 a_j^2 (\lambda_i^2 + \lambda_j^2 - 2\lambda_i \lambda_j)}{(1 + \bar{\gamma}\lambda_i)^3(1 + \bar{\gamma}\lambda_j)^3}. \end{aligned} \quad (53b)$$

Here,

$$\begin{aligned} \bar{T}_{ij} &= \lambda_i^2 a_i^2 a_j^2 (1 + \bar{\gamma}\lambda_j) + a_i^2 (1 + \bar{\gamma}\lambda_i) \lambda_j^2 a_j^2 \\ &\quad - \lambda_i a_i^2 (1 + \bar{\gamma}\lambda_i) a_j^2 \lambda_j - a_i^2 \lambda_i \lambda_j a_j^2 (1 + \bar{\gamma}\lambda_j) \\ &= a_i^2 a_j^2 (\lambda_i^2 (1 + \bar{\gamma}\lambda_j) + (1 + \bar{\gamma}\lambda_i) \lambda_j^2 - \lambda_i (1 + \bar{\gamma}\lambda_i) \lambda_j - \lambda_i \lambda_j (1 + \bar{\gamma}\lambda_j)). \end{aligned} \quad (53c)$$

Hence, substituting (53a) and (53b) into (52) gives

$$\frac{\partial h}{\partial \bar{\gamma}} = - \left( \sum_{i=1}^n \frac{a_i^2}{(1 + \bar{\gamma}\lambda_i)^2} \right)^{-2} \sum_{i,j=1}^n \frac{a_i^2 a_j^2 (\lambda_i - \lambda_j)^2}{(1 + \bar{\gamma}\lambda_i)^3(1 + \bar{\gamma}\lambda_j)^3}. \quad (54)$$

Now, for the combined derivative, we have

$$\begin{aligned} \left| \frac{\partial g}{\partial \gamma} \cdot \frac{\partial h}{\partial \bar{\gamma}} \right| &\leq \left( \sum_{i=1}^d \frac{s_i^2}{(M + (\gamma + 1)s_i)^2} \right)^{-2} \sum_{i,j=1}^d \frac{s_i^2 s_j^2 (s_i - s_j)^2}{(M + (\gamma + 1)s_i)^3 (M + (\gamma + 1)s_j)^3} \\ &\quad \cdot \left( \sum_{i=1}^n \frac{a_i^2}{(1 + \bar{\gamma}\lambda_i)^2} \right)^{-2} \sum_{i,j=1}^n \frac{a_i^2 a_j^2 (\lambda_i - \lambda_j)^2}{(1 + \bar{\gamma}\lambda_i)^3 (1 + \bar{\gamma}\lambda_j)^3}. \end{aligned}$$

Note that  $M + (\gamma + 1)s_i$  and  $1 + \bar{\gamma}\lambda_j$  are nonnegative for all  $i, j$ . Hence,

$$\begin{aligned} \left| \frac{\partial g}{\partial \gamma} \cdot \frac{\partial h}{\partial \bar{\gamma}} \right| &\leq \left( \sum_{i=1}^d \frac{s_i^2 (M + (\gamma + 1)s_i)}{(M + (\gamma + 1)s_i)^3} \right)^{-2} \\ &\quad \cdot \sum_{i,j=1}^d \frac{s_i^2 s_j^2 \cdot \Delta_1 \cdot (M + (\gamma + 1)s_j) (M + (\gamma + 1)s_i)}{(M + (\gamma + 1)s_i)^3 (M + (\gamma + 1)s_j)^3} \\ &\quad \cdot \left( \sum_{i=1}^n \frac{a_i^2 (1 + \bar{\gamma}\lambda_i)}{(1 + \bar{\gamma}\lambda_i)^3} \right)^{-2} \\ &\quad \cdot \sum_{i,j \in \mathcal{S}} \frac{a_i^2 a_j^2 \cdot \Delta_2 \cdot (1 + \bar{\gamma}\lambda_i) (1 + \bar{\gamma}\lambda_j)}{(1 + \bar{\gamma}\lambda_i)^3 (1 + \bar{\gamma}\lambda_j)^3}, \end{aligned}$$

where

$$\Delta_1 := \max_{i,j} \frac{(s_i - s_j)^2}{(M + (\gamma + 1)s_j) (M + (\gamma + 1)s_i)}, \quad \Delta_2 := \max_{i,j \in \mathcal{S}} \frac{(\lambda_i - \lambda_j)^2}{(1 + \bar{\gamma}\lambda_i) (1 + \bar{\gamma}\lambda_j)}. \quad (55)$$

Here,  $\mathcal{S} = \{i \in [n] | \lambda_i \neq 0\} \subset [n]$ .

Finally, setting  $\bar{\gamma} = g(\gamma)$ , we obtain

$$|h'(g(\gamma)) \cdot g'(\gamma)| = \left| \frac{\partial g}{\partial \gamma} \cdot \frac{\partial h}{\partial \bar{\gamma}} \right| \leq \Delta_1 \cdot \Delta_2 \leq \frac{\Delta_\Sigma^2 \cdot \Delta_R^2}{(M + s_{\min})^2} < 1.$$

where  $\Delta_\Sigma$  and  $\Delta_R$  are the spectral gaps of  $\Sigma$  and  $R$ ; and  $s_{\min}$  is the smallest eigenvalue of  $\Sigma$ ; and  $M = \sigma^2 + \sum_{i=1}^d s_i$ .  $\square$

*Proof of Theorem 4.* Lemma 1 shows that  $|\partial h(g(\gamma))/\partial \gamma| < 1$ , and as a result, the mapping  $h(g(\gamma))$  on  $\mathbb{R}_+$  is a contraction mapping. Therefore, by the Banach fixed-point theorem, this guarantees the existence of a unique root, denoted as  $\gamma = \gamma^*$ . This completes the proof of **T1**. In the following, we provide the proof of **T2**. Substituting the unique  $\gamma^*$  into (16) and using the fact that  $\tilde{P} = \sqrt{\Sigma} P \sqrt{\Sigma}$ , we obtain  $(P^*, \omega^*)$ , as given in (14), as a global minima of (3).

Next, we claim that  $(P^*, \omega^*)$  is the unique global minimizer of  $\mathcal{L}(P, \omega)$  up to rescaling, i.e., any other minimizer  $(\hat{P}, \hat{\omega})$  must be related to  $(P^*, \omega^*)$  by scaling factors  $\alpha$  and  $\beta$ , such that  $\hat{P} = \alpha P^*$  and  $\hat{\omega} = \beta \omega^*$ , for some  $\alpha, \beta > 0$ .

The loss function is given by

$$\mathcal{L}(P, \omega) = M - 2\text{tr}(\Sigma^2 P) \omega^\top r + M \|\omega\|^2 \text{tr}(\Sigma P^\top \Sigma P) + (\|\omega\|^2 + \omega^\top R \omega) \text{tr}(\Sigma^2 P^\top \Sigma P)$$

Now, consider the effect of rescaling the variables  $P$  and  $\omega$  by introducing scalars  $\alpha$  and  $\beta$ , i.e., we substitute  $\alpha P$  and  $\beta \omega$  into the loss function

$$\mathcal{L}(\alpha P, \beta \omega) = M - 2\alpha \beta \text{tr}(\Sigma^2 P) \omega^\top r + M \alpha^2 \beta^2 \|\omega\|^2 \text{tr}(\Sigma P^\top \Sigma P) + \alpha^2 \beta^2 (\|\omega\|^2 + \omega^\top R \omega) \text{tr}(\Sigma^2 P^\top \Sigma P).$$

Define

$$A := \text{tr}(\Sigma^2 P) \omega^\top r, \quad B := \text{tr}(\Sigma P^\top \Sigma P), \quad C := \|\omega\|^2, \quad D := \omega^\top R \omega, \quad E := \text{tr}(\Sigma^2 P^\top \Sigma P).$$

Thus, the rescaled loss function becomes

$$\mathcal{L}(\alpha P, \beta \omega) = M - 2\alpha \beta A + M \alpha^2 \beta^2 B C + \alpha^2 \beta^2 (C + D) E.$$

For  $(\mathbf{P}^*, \omega^*)$  to be a minimizer, the partial derivatives of the loss function with respect to  $\mathbf{P}$  and  $\omega$  must vanish at  $(\mathbf{P}^*, \omega^*)$ . However, we consider the effect of the rescaling in terms of  $\alpha$  and  $\beta$ . To find the stationary points of  $\mathcal{L}(\alpha\mathbf{P}, \beta\omega)$ , we differentiate with respect to  $\alpha$  and  $\beta$ :

$$\begin{aligned}\frac{\partial \mathcal{L}}{\partial \alpha} &= -2\beta A + 2M\alpha\beta^2 BC + 2\alpha\beta^2(C + D)E, \\ \frac{\partial \mathcal{L}}{\partial \beta} &= -2\alpha A + 2M\alpha^2\beta BC + 2\alpha^2\beta(C + D)E.\end{aligned}$$

Setting these to zero, we obtain the system

$$\alpha\beta(MBC + (C + D)E) = A. \quad (56)$$

This condition must hold for any minimizer. Now, suppose there exists another minimizer  $(\hat{\mathbf{P}}, \hat{\omega})$  that also minimizes the loss function. By the first-order optimality conditions,  $\alpha\beta$  must remain constant. This implies that any other minimizer  $(\hat{\mathbf{P}}, \hat{\omega})$  must be proportional to the original minimizer  $(\mathbf{P}^*, \omega^*)$ , meaning

$$\hat{\mathbf{P}} = \alpha\mathbf{P}^* \quad \text{and} \quad \hat{\omega} = \beta\omega^*$$

for some scalars  $\alpha, \beta > 0$  satisfying (56).

Thus, any global minimizer  $(\hat{\mathbf{P}}, \hat{\omega})$  is a *scaled* version of  $(\mathbf{P}^*, \omega^*)$ , and no other distinct minimizer exists. This proves uniqueness up to rescaling.  $\square$

### C.3 PROOF OF COROLLARY 2

*Proof.* Since by assumption  $\Sigma = \mathbf{I}$ , it follows from (13b) that

$$\begin{aligned}g(\gamma^*) &= \left( 1 + (d + \sigma^2) \sum_{i=1}^d \frac{1}{(d + \sigma^2 + \gamma^* + 1)^2} \left( \sum_{i=1}^d \frac{1}{(d + \sigma^2 + \gamma^* + 1)^2} \right)^{-1} \right)^{-1} \\ &= \frac{1}{d + \sigma^2 + 1}.\end{aligned}$$

Substituting this into (14) gives

$$\mathbf{P}^* = \mathbf{I}, \quad \text{and} \quad \omega^* = \left( \mathbf{R} + (d + \sigma^2 + 1)\mathbf{I} \right)^{-1} \mathbf{r}.$$

Now, recall that the objective function is given by

$$\mathcal{L}(\omega) = M - 2\text{tr}(\Sigma^2 \mathbf{P})\omega^\top \mathbf{r} + M \|\omega\|_{\ell_2} \text{tr}(\Sigma \mathbf{P}^\top \Sigma \mathbf{P}) + (\|\omega\|^2 + \omega^\top \mathbf{R}\omega)\text{tr}(\Sigma^2 \mathbf{P}^\top \Sigma \mathbf{P}),$$

and, by assumption,  $M = \sigma^2 + d$ .

Substituting  $\mathbf{P}^* = \mathbf{I}$  and  $\omega^* = \left( \mathbf{R} + (d + \sigma^2 + 1)\mathbf{I} \right)^{-1} \mathbf{r}$  into the objective (38), and using  $\Sigma = \mathbf{I}$ , we get:

$$\mathcal{L}(\omega^*) = (\sigma^2 + d) - 2 \cdot d \cdot \mathbf{r}^\top \omega^* + (\sigma^2 + d) \cdot \|\omega^*\|_{\ell_2}^2 d + d (\|\omega^*\|^2 + \omega^{*\top} \mathbf{R}\omega^*).$$

The expression simplifies as

$$\mathcal{L}(\omega^*) = (\sigma^2 + d) - 2d \cdot \mathbf{r}^\top \left( \mathbf{R} + (d + \sigma^2 + 1)\mathbf{I} \right)^{-1} \mathbf{r} + (\sigma^2 + d) d \|\omega^*\|_{\ell_2}^2 + d (\|\omega^*\|^2 + \omega^{*\top} \mathbf{R}\omega^*).$$

Next, we compute  $\|\omega^*\|^2$  and  $\omega^{*\top} \mathbf{R}\omega^*$ . By definition, we have

$$\|\omega^*\|^2 = \mathbf{r}^\top \left( \mathbf{R} + (d + \sigma^2 + 1)\mathbf{I} \right)^{-2} \mathbf{r},$$

and

$$\omega^{*\top} \mathbf{R}\omega^* = \mathbf{r}^\top \left( \mathbf{R} + (d + \sigma^2 + 1)\mathbf{I} \right)^{-1} \mathbf{R} \left( \mathbf{R} + (d + \sigma^2 + 1)\mathbf{I} \right)^{-1} \mathbf{r}.$$

Thus,

$$\begin{aligned}(d + \sigma^2 + 1)\|\omega^*\|^2 + \omega^{*\top} \mathbf{R}\omega^* &= \mathbf{r}^\top \left( \mathbf{R} + (d + \sigma^2 + 1)\mathbf{I} \right)^{-1} \left( (d + \sigma^2 + 1)\mathbf{I} + \mathbf{R} \right) \left( \mathbf{R} + (d + \sigma^2 + 1)\mathbf{I} \right)^{-1} \mathbf{r} \\ &= \mathbf{r}^\top \left( \mathbf{R} + (d + \sigma^2 + 1)\mathbf{I} \right)^{-1} \mathbf{r}.\end{aligned}$$

Substituting this result back into the objective function gives

$$\mathcal{L}(\omega^*) = (\sigma^2 + d) - d \cdot \mathbf{r}^\top \left( \mathbf{R} + (d + \sigma^2 + 1)\mathbf{I} \right)^{-1} \mathbf{r}.$$

$\square$

## D LOSS LANDSCAPE OF 1-LAYER GLA

### D.1 PROOF OF THEOREM 5

*Proof.* We first prove that under Assumption B,  $\mathcal{L}_{\text{GLA}}^* = \min_{P \in \mathbb{R}^{d \times d}, \omega \in \mathcal{W}} \mathcal{L}_{\text{WPGD}}(P, \omega)$  where  $\mathcal{W}$  is the search space of weighting vector  $\omega \in \mathbb{R}^n$  defined as

$$\mathcal{W} := \left\{ [\omega_1 \mathbf{1}_{n_1}^\top \cdots \omega_K \mathbf{1}_{n_K}^\top]^\top \in \mathbb{R}^n \mid 0 \leq \omega_i \leq \omega_j \leq 1, \forall 1 \leq i \leq j \leq K \right\}.$$

Define a set  $\bar{\mathcal{W}} := \left\{ [\omega_1 \cdots \omega_n]^\top \in \mathbb{R}^n \mid 0 \leq \omega_i \leq \omega_j \leq 1, \forall 1 \leq i \leq j \leq n \right\}$  and we have  $\mathcal{W} \in \bar{\mathcal{W}}$ .

Given scalar gating  $G_i = \begin{bmatrix} * & * \\ g_i \mathbf{1}^\top & * \end{bmatrix}$ , following (10), the weighting vector returns

$$\omega := [g_{1:n+1} \cdots g_{n:n+1}]^\top.$$

Since GLA with scalar gating valued in  $[0, 1]$  following Assumption B, that is,  $g_i \in [0, 1]$ . Therefore, we have  $g_{i:n+1} \leq g_{j:n+1}$  for  $1 \leq i \leq j \leq n$ . Therefore, any weighting vector implemented by GLA gating should be inside  $\bar{\mathcal{W}}$ .

Next, we will show that

$$\omega^* \in \mathcal{W} \quad \text{where} \quad \omega^* = \arg \min_{P, \omega \in \bar{\mathcal{W}}} \mathcal{L}_{\text{WPGD}}(P, \omega).$$

Define the weighting vector  $\omega = [\omega_1^\top \cdots \omega_K^\top]^\top \in \mathbb{R}^n$  where we have  $\omega_k = [\omega_1^{(k)} \cdots \omega_{n_k}^{(k)}]^\top \in \mathbb{R}^{n_k}$ . For any  $\omega \notin \mathcal{W}$ , there exist  $(i, j, k)$ ,  $i = j - 1$  such that  $\omega_i^{(k)} < \omega_j^{(k)}$ . Given gradient in (42), we have that  $\nabla_{\omega_i^{(k)}} \mathcal{L} = c_1 \cdot \omega_i^{(k)} + c_2$  and  $\nabla_{\omega_j^{(k)}} \mathcal{L} = c_1 \cdot \omega_j^{(k)} + c_2$  with for some  $c_1, c_2$  with  $c_1 > 0$ .  $\nabla_{\omega_i^{(k)}} \mathcal{L} < \nabla_{\omega_j^{(k)}} \mathcal{L}$ . Therefore either increasing  $\omega_j^{(k)}$  (if  $\nabla_{\omega_i^{(k)}} \mathcal{L} < 0$ ) or decreasing  $\omega_j^{(k)}$  (if  $\nabla_{\omega_j^{(k)}} \mathcal{L} > 0$ ) will reduce the loss. This results in showing that  $\omega^* \in \mathcal{W}$ .

Finally, we will show that any  $\omega \in \mathcal{W}$  can be obtained via the GLA gating. Let  $\omega = [\omega_1 \mathbf{1}_{n_1}^\top \cdots \omega_K \mathbf{1}_{n_K}^\top]^\top$  be any vector in  $\mathcal{W}$  and assume that  $\omega_K = \alpha < 1$  without loss of generality. Then such sample weighting can be achieved via the gating

$$\left[ \mathbf{1}_{n_1}^\top \quad \frac{\omega_1}{\omega_{1:K}} \quad \cdots \quad \mathbf{1}_{n_k}^\top \quad \frac{\omega_k}{\omega_{k:K}} \quad \cdots \quad \mathbf{1}_{n_K}^\top \quad \frac{\omega_K}{\omega_{K:K}} \right]^\top.$$

Let  $\omega'_k := \frac{\omega_k}{\omega_{k:K}}$  and let  $\mathbf{w}_g$  be in the form of

$$\mathbf{w}_g = \begin{bmatrix} \mathbf{0}_{d+1} \\ \tilde{\mathbf{w}}_g \end{bmatrix} \in \mathbb{R}^{d+p+1}.$$

Then it remains to show that there exists  $\tilde{\mathbf{w}}_g$  satisfying:

$$\phi(\tilde{\mathbf{w}}_g^\top \tilde{\mathbf{d}}_k) \begin{cases} = 1, & k = 0 \\ = \omega'_k, & k \in [K] \end{cases}$$

Assumption B implies the feasible of the problem, which completes the proof of (23).

Recap the optimal weighting from (14) which takes the form of

$$\omega^* = (g(\gamma^*) \cdot \mathbf{R} + \mathbf{I})^{-1} \mathbf{r}.$$

Since Assumption C holds and  $n_1 = n_2 = \cdots = n_K := \bar{n}$ ,  $\omega^*$  takes the form of  $\omega^* = c\mathbf{r}$  for some positive constant  $c$ . Therefore, the optimal weighting (up to a scalar) is inside the set  $\mathcal{W}$ . Combining it with (23) completes the proof.  $\square$

## D.2 PROOF OF THEOREM 6

*Proof.* Following the similar proof of Theorem 5, and letting  $\tilde{\mathcal{W}} := \left\{ \left[ \omega_1 \mathbf{1}_{n_1}^\top \cdots \omega_K \mathbf{1}_{n_K}^\top \right]^\top \in \mathbb{R}^n \right\}$ , we obtain

$$\min_{P, \omega \in \tilde{\mathcal{W}}} \mathcal{L}_{\text{WPGD}}(P, \omega) = \min_{P, \omega} \mathcal{L}_{\text{WPGD}}(P, \omega).$$

Therefore, it remains to show that any  $\omega \in \tilde{\mathcal{W}}$  can be implemented via some gating function. Let  $\omega = [\omega_1 \mathbf{1}_{n_1}^\top \cdots \omega_K \mathbf{1}_{n_K}^\top]$  be arbitrary weighting in  $\tilde{\mathcal{W}}$ . Theorem 5 has shown that if  $\omega_1 \leq \omega_2 \leq \cdots \leq \omega_K$ , GLA with scalar function can implement such increasing weighting.

Now, inspired from Appendix B that all dimensions in the output implement individual WPGD. We can decouple the weighting into  $K$  separate weighting:

$$\begin{aligned} \omega_1 &= \omega_1 [\mathbf{1}_{n_1}^\top \cdots \mathbf{1}_{n_K}^\top] \\ \omega_2 &= (\omega_2 - \omega_1) [\mathbf{0}_{n_1}^\top \mathbf{1}_{n_2}^\top \cdots \mathbf{1}_{n_K}^\top] \\ \omega_3 &= (\omega_3 - \omega_2) [\mathbf{0}_{n_1}^\top \mathbf{0}_{n_2}^\top \mathbf{1}_{n_3}^\top \cdots \mathbf{1}_{n_K}^\top] \\ &\vdots \\ \omega_K &= (\omega_K - \omega_{K-1}) [\mathbf{0}_{n_1}^\top \mathbf{0}_{n_2}^\top \mathbf{0}_{n_3}^\top \cdots \mathbf{0}_{n_{K-1}}^\top \mathbf{1}_{n_K}^\top] \end{aligned}$$

and we have  $\omega = \sum_{k=1}^K \omega_k$ . Recap from Appendix B and consider the construction  $\mathbf{W}_v = [\mathbf{0}_{(d+1) \times d} \mathbf{u}]$ . Assumption B implies that  $K \leq p < d + p + 1$ .

From (27) and (28), let  $i$ 'th dimension implements the weighting  $\omega_i$  for  $i \in [K]$ . Specifically, let  $g^i$  implement weighting  $[\mathbf{0}_{n_1}^\top \cdots \mathbf{0}_{n_{i-1}}^\top \mathbf{1}_{n_i}^\top \cdots \mathbf{1}_{n_K}^\top]$  (which is feasible due to Theorem 5) and set  $u_i = \omega_i - \omega_j$  with  $\omega_0 = 0$ . Then the composed weighting following (28) returns  $\omega$ , which completes the proof.  $\square$

## D.3 PROOF OF COROLLARY 3

*Proof.* Recap from (43a) that given  $\Sigma = \mathbf{I}$  and  $\omega = \mathbf{1}$ ,

$$\begin{aligned} P^* &= \left( (d + \sigma^2) \mathbf{n} \mathbf{I} + (\mathbf{n} + \mathbf{1}^\top \mathbf{R} \mathbf{I}) \mathbf{I} \right)^{-1} \mathbf{1}^\top \mathbf{r} \\ &= \frac{\mathbf{1}^\top \mathbf{r}}{n(d + \sigma^2 + 1) + \mathbf{1}^\top \mathbf{R} \mathbf{I}} \mathbf{I} := c \mathbf{I}. \end{aligned}$$

Then taking it back to the loss function (c.f. (38)) obtains

$$\begin{aligned} \mathcal{L}(P^*, \omega = \mathbf{1}) &= d + \sigma^2 - 2cd \mathbf{1}^\top \mathbf{r} + (d + \sigma^2) c^2 nd + (\mathbf{n} + \mathbf{1}^\top \mathbf{R} \mathbf{I}) c^2 d \\ &= d + \sigma^2 - cd \mathbf{1}^\top \mathbf{r}. \end{aligned}$$

It completes the proof.  $\square$