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

025 026

003

006

008

009 010

011

012 013

014

015

016

017

018

021

023

027 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 029 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 031 time approaches to sequence modeling. The initial approaches included linear attention and statespace 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), 033 034 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. 037

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

$$\mathbf{S}_i = \mathbf{G}_i \odot \mathbf{S}_{i-1} + \mathbf{v}_i \mathbf{k}_i^{\mathsf{T}}, \quad \text{and} \quad \mathbf{o}_i = \mathbf{S}_i \mathbf{q}_i. \tag{1}$$

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:

053

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 $(G_i)_{i=1}^n$. 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 G_i are identical, is described by the prediction:

$$\hat{\mathbf{y}} = \mathbf{x}^{\mathsf{T}} \mathbf{P} \mathbf{X}^{\mathsf{T}} (\mathbf{y} \odot \boldsymbol{\omega}). \tag{2}$$

Here, $X \in \mathbb{R}^{n \times d}$ is the input feature matrix; $y \in \mathbb{R}^n$ is the associated label vector; $x \in \mathbb{R}^d$ represents the test/query input to predict; $P \in \mathbb{R}^{d \times d}$ is the preconditioning matrix; and $\omega \in \mathbb{R}^n$ weights the individual samples. When ω is fixed, we drop "data-dependent" and simply refer to this algorithm as the WPGD algorithm. However, for GLA, $\omega := \omega(X, y)$ 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⇔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 $(X, y, x, y) \sim D$, with y being the label associated with the query, we investigate the population risk minimization:

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

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

- 080
081
081
082 \diamond 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 $\omega \in C$, where the restriction 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

060

065

067

068

069

070

071

073 074 075

076

077

078

079

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-090 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 092 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); 094 Jelassi et al. (2024) establish theoretical and empirical shortcomings of recurrent models for solving 095 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 096 097 closely related to implicit self-attention frameworks. For example, the work by Zimerman et al. 098 (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. 102

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,

108 Von Oswald et al. (2023) provide a construction of weights in linear attention-only transformers that 109 can replicate GD steps for a linear regression task on in-context examples. Notably, they observe 110 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. 111 (2024) focus on the loss landscape of ICL for linear attention models. For a single-layer model trained 112 on in-context prompts for random linear regression tasks, Mahankali et al. (2023); Ahn et al. (2024) 113 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., 114 2023) analyzes the challenges of causal masking in causal language models (causalLM), showing 115 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, 116 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 restricted to the trivial setting of equal 117 118 weights, relying on the standard prompt model with IID examples and shared tasks. In contrast, we 119 propose novel multitask datasets and prompt models where nontrivial weighting is provably optimal. This allows us to characterize the loss landscape of WPGD and explore advanced GLA models, 120 linking them to data-dependent WPGD algorithms. 121

122 123

124

133 134

139

140

149 150 151

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., *a* and *A*, represent vectors and matrices. The identity matrix of size *n* is I_n . 1 and 0 denote the all-one and all-zero vectors or matrices of proper size. $\mathcal{N}(\mu, \Sigma)$ is the Gaussian distribution with mean μ and covariance Σ . The symbol \odot denotes the Hadamard product and \oslash denotes Hadamard division. Given $a_{i+1}, \dots, a_j \in \mathbb{R}^d$, we use $a_{i:j}$ to denote $a_{i+1} \odot \dots \odot a_j$ for i < j, and $a_{i:i} = 1_d$ 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

$$\mathbf{Z} = \begin{bmatrix} z_1 \cdots z_n \ z_{n+1} \end{bmatrix}^\top = \begin{bmatrix} \mathbf{x}_1 \cdots \mathbf{x}_n & \mathbf{x}_{n+1} \\ y_1 & \cdots & y_n & 0 \end{bmatrix}^\top \in \mathbb{R}^{(n+1) \times (d+1)},$$
(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 $F \in \mathbb{R}^{(n+1)\times(d+1)} \to \mathbb{R}$ 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

$$\mathcal{L}(F) = \mathbb{E}_{\mathcal{D}}\left[(y - F(\mathbf{Z}))^2 \right].$$
(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(Z) takes the form $\hat{y} = z_{n+1}^{\top} W_q W_k^{\top} Z^{\top} Z W_v h$ where $W_k, W_q, W_v \in \mathbb{R}^{(d+1) \times (d+1)}$ are attention parameters, and $h \in \mathbb{R}^{d+1}$ is the linear prediction head. We assume that the in-context input-label pairs follow a *shared-task distribution*, where $\beta \sim \mathcal{N}(0, \Sigma_{\beta})$, x_i are i.i.d. with $x_i \sim \mathcal{N}(0, \Sigma_x)$, and $y_i \sim \mathcal{N}(\beta^{\top} x_i, \sigma^2)$, where $\sigma \ge 0$ 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 $\hat{\beta}$ coincides with the one-step optimal preconditioned gradient descent. In particular, we have $\hat{\beta} = P^* X^{\top} y$, where

$$\boldsymbol{P}^{\star} = \operatorname*{argmin}_{\boldsymbol{P} \in \mathbb{R}^{d \times d}} \mathbb{E}_{\mathcal{D}} \left[\left(\boldsymbol{y} - \boldsymbol{x}^{\top} \boldsymbol{P} \boldsymbol{X}^{\top} \boldsymbol{y} \right)^{2} \right] \quad \text{with} \quad \boldsymbol{X} := \begin{bmatrix} \boldsymbol{x}_{1} & \cdots & \boldsymbol{x}_{n} \end{bmatrix}^{\top} \quad \text{and} \quad \boldsymbol{y} := \begin{bmatrix} \boldsymbol{y}_{1}, \cdots, \boldsymbol{y}_{n} \end{bmatrix}^{\top}.$$
(6)

Linear attention and gating. Given the input prompt Z, let $Q = ZW_q$, $K = ZW_k$ and $V = ZW_y$ 152 be the corresponding query, key, and value embedding matrices, respectively. The output of causal 153 linear attention at time *i* can be computed in a recurrent form as $S_i = S_{i-1} + v_i k_i^{\top}$ and $o_i = S_i q_i$ where 154 $q_i, k_i, v_i \in \mathbb{R}^{d+1}$ are the query, key, value embeddings of z_i and $S_0 = 0$. This recurrent form implies 155 that linear attention has $O(d^2)$ cost, that is independent of N, to generate per-token. As presented in 156 (1), GLA follows the same structure as linear attention but with a gating mechanism, which equips 157 the model with the option to pass or supress 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. 158 159

160 We will show that GLA can weigh the context window through gating, thus, its capabilities are linked 161 to the WPGD algorithm described in (7). This will in turn facilitate GLA to effectively learn *multitask* prompt distributions described by $y_i \sim \mathcal{N}(\boldsymbol{\beta}_i^{\mathsf{T}} \boldsymbol{x}_i, \sigma^2)$ with $\boldsymbol{\beta}_i$'s not necessarily identical.

162 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.

166

167 3.1 GLA AS A DATA-DEPENDENT WPGD PREDICTOR

Data-Dependent WPGD. Given X and y as defined in (6), consider the weighted least squares objective $\mathcal{L}(\beta) = \sum_{i=1}^{n} \Omega_i \cdot (y_i - \beta^T x_i)^2$ with weights $\Omega \in \mathbb{R}^n$. To optimize this, we use gradient descent (GD) starting from zero initialization, $\beta_0 = 0$ with a step size of $\eta = 1/2$. One step of standard GD is given by

178

185

187

196 197

204 205

206

207 208 $\boldsymbol{\beta}_1 = \boldsymbol{\beta}_0 - \eta \nabla \mathcal{L}(\boldsymbol{\beta}_0) = \sum_{i=1}^n \Omega_i \cdot \boldsymbol{x}_i \boldsymbol{y}_i = \boldsymbol{X}^\top (\boldsymbol{\Omega} \odot \boldsymbol{y}).$

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

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

179 Above is the basic *scalar-weighted* WPGD predictor which weights individual datapoints. It turns out, 180 *vector-valued gating* can facilitate a more general estimator which weights individual coordinates. 181 To this aim, we introduce an extension as follows: Let $P_1, P_2 \in \mathbb{R}^{d \times d}$ denote the preconditioning 182 matrices, and let $\Omega \in \mathbb{R}^{n \times d}$ denote the *vector-valued weighting* matrix. Note that Ω is now a matrix 183 rather than vector to facilitate coordinate-wise weighting and will remain consistent throughout the 184 paper. We can similarly define

$$\boldsymbol{\beta}_{1}^{\mathrm{gu}}(\boldsymbol{P}_{1},\boldsymbol{P}_{2},\boldsymbol{\Omega}) := \boldsymbol{P}_{2}(\boldsymbol{X}\boldsymbol{P}_{1}\odot\boldsymbol{\Omega})^{\mathsf{T}}\boldsymbol{y}$$
(7a)

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

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

We note that by removing the preconditioning matrices P_1 , P_2 , and the weighting matrix Ω 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 $\Omega = \omega \mathbf{1}_d^{\mathsf{T}}$, 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:

$$\boldsymbol{W}_{k} = \begin{bmatrix} \boldsymbol{P}_{k}^{\top} & \boldsymbol{0} \\ \boldsymbol{0} & \boldsymbol{0} \end{bmatrix}, \quad \boldsymbol{W}_{q} = \begin{bmatrix} \boldsymbol{P}_{q}^{\top} & \boldsymbol{0} \\ \boldsymbol{0} & \boldsymbol{0} \end{bmatrix} \quad \text{and} \quad \boldsymbol{W}_{v} = \begin{bmatrix} \boldsymbol{0}_{d \times d} & \boldsymbol{0} \\ \boldsymbol{0} & \boldsymbol{1} \end{bmatrix}, \tag{8}$$

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 \mathbb{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

$$f_{GLA}(\mathbf{Z}) := \boldsymbol{o}_{n+1}^{\top} \boldsymbol{h} = \hat{\boldsymbol{\beta}}^{\top} \boldsymbol{x} \quad where \quad \hat{\boldsymbol{\beta}} = \boldsymbol{\beta}_{1}^{gd}(\boldsymbol{P}_{k}, \boldsymbol{P}_{q}, \boldsymbol{\Omega}).$$

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

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

212 Observe that, crucially, since g_i (or G_i) is associated with z_i , z_i influences the weighting of all history 213 $z_{j < i}$. We defer the proof of Theorem 1 to the Appendix B.1. It is noticeable that only d of the total 214 $(d + 1)^2$ entries in each gating matrix G_i are useful due to the model construction presented in (8). 215 However, if we relax the weight restriction, e.g., $W_v = [\mathbf{0}_{(d+1) \times d} \mathbf{1}_{d+1}]$, then the weighting matrix Ω in

Theorem 1 is associated with all rows of the G_i matrices. We defer the discussion to Appendix B.1.

216 3.2 CAPABILITIES OF MULTI-LAYER GLA

218 Ahn et al. (2024) demonstrated that, with appropriate construction, an *L*-layer linear attention model 219 performs *L*-step preconditioned gradient descent on the dataset $(\mathbf{x}_i, y_i)_{i=1}^n$ provided within the prompt. 220 In this work, we study multi-layer GLA and analyze the associated algorithm class it can emulate. 221 It is worth mentioning that Ahn et al. (2024) does not consider *causal masking* which is integral to 222 multilayer GLA due to its recurrent nature described in (1). Our analysis will capture the impact of 222 gating and causal mask through *n* separate gradient descent trajectories that are coupled.

Consider an *L*-layer GLA model. For $\ell \in [L]$, let Z_{ℓ} and O_{ℓ} denote the input and output of the ℓ 'th layer. In practice, residual connections are commonly applied. Hence, we define the updated output of the ℓ 'th layer (after applying the residual connection) as $\tilde{O}_{\ell} := Z_{\ell} + O_{\ell}$. Note that \tilde{O}_{ℓ} also serves as the input to the $(\ell + 1)$ 'th layer, i.e., $Z_{\ell+1} = \tilde{O}_{\ell}$. In the following, we focus on (d + 1)'th entries of each token's output at each layer, denoted by $\tilde{o}_{i,\ell} := (\tilde{O}_{\ell})_{i,d+1}$ 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 ℓ '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., $G_i = g(\mathbf{x}_i)$, and let Ω be defined as in Theorem 1. Additionally, denote the masking as $M_i = \begin{bmatrix} I_i & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \in \mathbb{R}^{n \times n}$, and let $\hat{\beta}_0, \beta_{i,0} = \mathbf{0}$ for $i \in [n]$.

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

233 234

236 237 238

239 240 241

- 235
- For $i \leq n$, $\tilde{o}_{i,\ell} = y_i \boldsymbol{x}_i^\top \boldsymbol{\beta}_{i,\ell}$ where $\boldsymbol{\beta}_{i,\ell} = \boldsymbol{\beta}_{i,\ell-1} + \boldsymbol{P}_{q,\ell} (\nabla_{i,\ell} \oslash \boldsymbol{g}_{i:n+1})$,

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

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

$$\nabla_{i,\ell} = \bar{X}_{\ell}^{\top} M_i \left(\hat{y}_{\ell} - y \right).$$

We defer the proof of Theorem 2 to the Appendix B.2. Theorem 2 states that an L-layer GLA 242 implements L steps of WPGD but with gradient in a recurrent form. To recap, given data (X, y)243 and prediction $\hat{\beta}$, the gradient with respect to the squared loss takes the form $X^{\top}(X\hat{\beta} - y)$, up to 244 some constant c. In comparison, $P_{q,\ell}(\nabla_{i,\ell} \otimes g_{i:n+1})$ similarly acts as a gradient but incorporates 245 layer-wise feature preconditioners $(P_{q,\ell}, P_{k,\ell})$, data weighting (Ω) , and causality $(g_{i:n+1}, M_i)$. Here, 246 M_i represents causal masking, ensuring that at time *i*, only inputs from $j \le i$ 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 247 248 causal masking limits convergence by introducing sequence biases, akin to online gradient descent 249 with non-decaying step sizes. 250

To simplify the theorem statement, we assume that the gating function depends only on the input feature, e.g., $G_i = g(x_i)$, 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 α_{ℓ} 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 ℓ 'th layer be parameterized by $P_{q,\ell}$, $P_{k,\ell}$ as in (8) and define $P_{\ell} := P_{q,\ell}P_{k,\ell}^{\top}$, $\ell \in [L]$. Let $\hat{\beta}_0, \beta_{i,0} = 0$ for $i \in [n]$. Then, the (d + 1)'th entry of the *i*'th token of the ℓ 'th layer outputs satisfies:

- For $i \leq n$, $\tilde{o}_{i,\ell} = y_i \mathbf{x}_i^\top \boldsymbol{\beta}_{i,\ell}$ where $\boldsymbol{\beta}_{i,\ell} = \boldsymbol{\beta}_{i,\ell-1} + \boldsymbol{P}_\ell \nabla_{i,\ell}$,
- 259 260 261 262
- $\tilde{o}_{n+1,\ell} = -\mathbf{x}^{\top} \hat{\boldsymbol{\beta}}_{\ell}$ where $\hat{\boldsymbol{\beta}}_{\ell} = (1 + \alpha_{\ell}) \hat{\boldsymbol{\beta}}_{\ell-1} + \boldsymbol{P}_{\ell} \nabla_{n,\ell}$ and $\alpha_{\ell} = \mathbf{x}^{\top} \boldsymbol{P}_{\ell} \mathbf{x}$.

Here, we define $\nabla_{i,\ell} = X^{\top} M_i (\hat{y}_{\ell} - y)$ with \hat{y}_{ℓ}, 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).}

270 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 $g(z_i) = G_i$ in GLA, as detailed in Theorem 1. Now let us consider the widely used types of gating functions, such as $G_i = \alpha_i \mathbf{1}_{d+1}^{\mathsf{T}}$ (Yang et al., 2023; Katsch, 2023; Qin et al., 2024; Peng et al., 2024) or $G_i = \gamma_i \mathbf{1}_{d+1} \mathbf{1}_{d+1}^{\mathsf{T}}$ (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 $\begin{bmatrix} * & * \\ g_i \mathbf{1}_d^{\mathsf{T}} & * \end{bmatrix}$, thus simplifying the predictor to a sample-weighted PGD, as given by

280 281

283 284 285

286

302 303

305

315 316

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

where $P = P_q P_k^{\top}$ 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 ygiven x via the 1-step WPGD prediction $\hat{y} = x^T \hat{\beta}$, with $\hat{\beta}$ as in (10). The algorithm learning problem is given by (3) which minimizes the WPGD risk $\mathbb{E}_{\mathcal{D}}[(y - x^T PX(\omega \odot y))^2]$.

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 ω consists of all ones and does not specifically explore the optimization landscape of ω 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 $r_{ij} = \mathbb{E}[\beta_i^{\mathsf{T}}\beta_j]/d$ for $i, j \in [n + 1]$, and the task and correlation matrices

$$\boldsymbol{\beta} := \boldsymbol{\beta}_{n+1}, \quad \boldsymbol{B} = [\boldsymbol{\beta}_1 \ \dots \ \boldsymbol{\beta}_n]^{\mathsf{T}}, \quad \boldsymbol{R} = \frac{1}{d} \mathbb{E}[\boldsymbol{B}\boldsymbol{B}^{\mathsf{T}}], \quad and \quad \boldsymbol{r} = \frac{1}{d} \mathbb{E}[\boldsymbol{B}\boldsymbol{\beta}]. \tag{11}$$

Additionally, for any $i, j \in [n + 1]$, $\beta_i - r_{ij}\beta_j$ is independent of β_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). $(\boldsymbol{\beta}_i)_{i=1}^{n+1}$ are drawn according to the correlated task model of Definition 1, $(\boldsymbol{x}_i)_{i=1}^{n+1} \in \mathbb{R}^d$ are IID following $\boldsymbol{x}_i \sim \mathcal{N}(0, \boldsymbol{\Sigma})$ and $y_i \sim \mathcal{N}(\boldsymbol{x}_i^{\mathsf{T}} \boldsymbol{\beta}_i, \sigma^2)$ for $i \in [n+1]$.

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

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

Assumption A. For the correlation vector \mathbf{r} from (11), we have $\mathbf{r} = \mathbf{E}\mathbf{a}$ for some $\mathbf{a} = [a_1, \dots, a_n]^\top \in \mathbb{R}^n$ 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, ω) of the WPGD objective in (3).

327 328

330

331 332

341

347

349

355

360

361

362

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

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

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

where $\{s_i\}_{i=1}^d$ and $\{\lambda_i\}_{i=1}^n$ are the eigenvalues of Σ 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 $\mathcal{L}(\mathbf{P}, \boldsymbol{\omega})$ in (3) has a stationary point $(\mathbf{P}^{\star}, \boldsymbol{\omega}^{\star})$, up to rescaling, defined as

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

where γ^{\star} is a fixed point of composite function $h(g(\gamma))$.

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

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

$$\Delta_{\Sigma} \cdot \Delta_{R} < M + s_{\min},\tag{15}$$

where Δ_{Σ} and Δ_{R} denote the effective spectral gaps of Σ and R, respectively, as given in (12); s_{\min} is the smallest eigenvalue of Σ ; 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 $\mathcal{L}(\mathbf{P}, \boldsymbol{\omega})$ has a unique (global) minima ($\mathbf{P}^{\star}, \boldsymbol{\omega}^{\star}$), up to re-scaling, given by (14).

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

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

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

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

Substituting the expression for $\omega(\gamma)$ into $\gamma = \frac{\omega^{\top} R \omega}{\|\omega\|^2}$, and applying Assumption A, we obtain the equation $\gamma = h(g(\gamma))$. We then show that whenever $\Delta_{\Sigma} \cdot \Delta_R < M + s_{\min}$, 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 γ^* into (16) implies that (P^*, ω^*), 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.

Theorem 4 establishes mild conditions under which a unique (global) minimum (P^* , ω^*) 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 (P, ω) optimization.

Remark 1 An interesting observation about the optimal gating parameter ω^* is its connection to the correlation matrix R, which captures the task correlations in a multitask learning setting. Specifically, the optimal gating given in (14) highlights how ω^* depends directly on both the task correlation matrix R and the vector 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 $\Delta_{\Sigma} \cdot \Delta_R < M + s_{\min}$, 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 $\Sigma = I$. Then, $\Delta_{\Sigma} = 0$, satisfying Condition (15), and we have $g(\gamma^*) = \frac{1}{d+\sigma^2+1}$, which yields

$$P^{\star} = I$$
, and $\omega^{\star} = \left(R + (d + \sigma^2 + 1)I\right)^{-1} r$. (17)

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

$$\mathcal{L}_{WPGD}^{\star} = d + \sigma^2 - d \cdot \boldsymbol{r}^{\top} \left(\boldsymbol{R} + (d + \sigma^2 + 1) \boldsymbol{I} \right)^{-1} \boldsymbol{r}.$$
 (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 ω is restricted and data-dependent, meaning that $\mathcal{L}_{\text{WPGD}}^{\star}$ 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.

400 5.1 Multi-task prompt model

384 385 386

387 388

389 390

391 392

399

409 410

419

420

424

425

426

427

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

$$\mathcal{L}_{\mathsf{GLA}}^{\star} := \min_{\boldsymbol{P}_{k,q},g} \mathcal{L}_{\mathsf{GLA}}(\boldsymbol{P}_{k},\boldsymbol{P}_{q},g) \quad \text{where} \quad \mathcal{L}_{\mathsf{GLA}}(\boldsymbol{P}_{k},\boldsymbol{P}_{q},g) = \mathbb{E}_{\mathcal{D}}\Big[(y - f_{\mathsf{GLA}}(\boldsymbol{Z}))^{2}\Big]. \tag{19}$$

411 Here, $g \in G$ represents the gating function.

Note that 1) the task vectors $(\boldsymbol{\beta}_k)_{k=1}^K$ are not explicitly shown in the prompt, 2) examples $(\boldsymbol{x}_i^{(k)}, y_i^{(k)})$ are randomly drawn, and 3) the gating function is applied to the tokens/input samples $(\boldsymbol{Z}_k)_{k=1}^K$. 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 $(d_k)_{k=1}^K$ be the delimiters that determine stop of the tasks. Specifically, the final prompt is given by

 $\mathbf{Z} = \begin{bmatrix} \mathbf{Z}_1^{\mathsf{T}} & d_1 & \cdots & \mathbf{Z}_K^{\mathsf{T}} & d_K & z \end{bmatrix}^{\mathsf{T}}.$ (20)

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

- For task prompts Z_k : Contextual features are set to a fixed vector $\bar{d}_0 \neq 0$.
- For delimiters d_k : Data features are set to zero (e.g., $x_i = 0$ and $y_i = 0$) so that $d_k = [\mathbf{0}_{d+1} \, \bar{d}_k]$ where \bar{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 $(d_k)_{k=1}^K$, 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



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.

WPGD loss $\mathcal{L}_{WPGD}^{\star}$ 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.

451 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.

464

465 466 467

470

471

472

473 474 475

476 477

478

479 480

481

450

452

$$\tilde{W}_{k} = \begin{bmatrix} W_{k} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}, \quad \tilde{W}_{q} = \begin{bmatrix} W_{q} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \text{ and } \tilde{W}_{v} = \begin{bmatrix} W_{v} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}.$$
(21)

Here, $\tilde{W}_{k,q,v} \in \mathbb{R}^{(d+p+1)\times(d+p+1)}$ and $W_{k,q,v} \in \mathbb{R}^{(d+1)\times(d+1)}$ 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.

Assumption C. The correlation between context tasks $(\boldsymbol{\beta}_k)_{k=1}^K$ and query task $\boldsymbol{\beta}$ satisfies $\mathbb{E}[\boldsymbol{\beta}_i^{\mathsf{T}}\boldsymbol{\beta}_j] = 0$ and $\mathbb{E}[\boldsymbol{\beta}_i^{\mathsf{T}}\boldsymbol{\beta}] \le \mathbb{E}[\boldsymbol{\beta}_j^{\mathsf{T}}\boldsymbol{\beta}]$ for $1 \le i \le j \le K$.

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

$$\boldsymbol{X} = \begin{bmatrix} \boldsymbol{X}_1^\top & \cdots & \boldsymbol{X}_K^\top \end{bmatrix}^\top \in \mathbb{R}^{n \times d} \quad \text{and} \quad \boldsymbol{y} = \begin{bmatrix} \boldsymbol{y}_1^\top & \cdots & \boldsymbol{y}_K^\top \end{bmatrix}^\top \in \mathbb{R}^n.$$
(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 $\mathcal{L}_{WPGD}(P, \omega)$ from (3) with dataset (X, y) defined in (22). Suppose Assumption *B* holds and consider GLA with scalar gating $g(z) = \phi(w_g^{\top} z) \mathbf{1} \mathbf{1}^{\top}$ where w_g is the trainable parameter. Consider input prompt **Z** defined in (20) and model constructions described in (21). Then the optimal risk $\mathcal{L}_{GLA}^{\star}$ defined in (19) obeys

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

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

$$\mathcal{L}_{GLA}^{\star} = \mathcal{L}_{WPGD}^{\star}.$$
 (24)

Assumption B ensures that any ω in W can be achieved by an appropriate choice of gating parameters. Furthermore, Assumption C guarantees that the optimal choice of ω under the WPGD objective lies within the search space 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

486 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 $\mathcal{L}_{GLA}^{\star}$ obtained from (23) is lower than the optimal WPGD loss $\mathcal{L}_{WPGD}^{\star}$. Further experimental details are deferred to Appendix A. 487 488 489

490 Loss landscape of vector gating. Till now, much of our discussion has focused on the scalar gating 491 setting. It is important to highlight that, even in the scalar-weighting context, analyzing the WPGD setting. It is important to ingring that, even in the scatar-weighting context, analyzing the WPGD problem remains non-trivial due to the joint optimization over (P, ω) . However, as demonstrated in Theorem 5, scalar gating can only express weightings within the set W. If Assumption C does not hold, $\mathcal{L}_{GLA}^{\star}$ cannot achieve the optimal WPGD loss (see the misalignment between red solid curve, presenting $\mathcal{L}_{GLA}^{\star}$, and black dashed curve, presenting $\mathcal{L}_{WPGD}^{\star}$ in Figure 1d). We argue that vector gating overcomes this limitation by applying distinct weighting mechanisms across different dimensions, facilitating stronger expressivity. 492 493 494 495 496

497 **Theorem 6** (Vector Gating). Recall input prompt Z from (20) and model constructions from (21) but with $W_v = [\mathbf{0}_{(d+1)\times d} \mathbf{u}]$. Suppose Assumption **B** holds and consider GLA with vector gating 498 $g(z) = \phi(W_g z) \mathbf{1}^{\top}$. Here, \mathbf{u} and W_g are trainable parameters. Consider Problem (19), where we employ a vector gating $g(z) = \phi(W_g z) \mathbf{1}^{\top}$. Let $\mathcal{L}_{\text{GLA-v}}^{\star}$ denote its optimal risk, and $\mathcal{L}_{\text{WPGD}}^{\star}$ be defined as 499 500 in (3). Then, the optimal risk obeys $\mathcal{L}_{GLA-v}^{\star} = \mathcal{L}_{WPGDv}^{\star}$ 501

502 In Theorem 5, the equivalence between $\mathcal{L}_{GLA}^{\star}$ and $\mathcal{L}_{WPGD}^{\star}$ 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 $\mathcal{L}_{GLA-\nu}^{\star} = \mathcal{L}_{WPGD}^{\star}$. Specifically, under the bounded activation model of Assumption **B**, scalar gating is unable to express non-monotonic weighting schemes. For instance, 504 505 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 506 507 508 509 curve, which outperform the scalar gating results (red solid) and align with the optimal WPGD loss (black dashed). 510

Loss landscape of 1-layer linear attention. Inspired by the fact that linear attention implements all ones gating, that is, $G_i \equiv 1$. Consider training a single-layer linear attention and let $f_{ATT}(Z) :=$ 511 512 $f_{\text{GLA}}(\mathbf{Z}, \mathbf{G}_i \equiv \mathbf{1})$ be its prediction. Let $\mathcal{L}_{\text{ATT}}^{\star}$ be the corresponding optimal risk following (19). 513

Corollary 3. Consider a single-layer linear attention following model construction in (8) and 514 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 515 516

517

518 519

525

527 528

529

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

(25)

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

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

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.

DISCUSSION 6

530 To summarize, this work offers a fresh theoretical perspective on gated linear attention models 531 through in-context learning by showing that they can emulate data-dependent weighted preconditioned 532 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 534 mathematical interest in its own right, developed the first characterization of its optimization landscape, 536 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 538 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 539 on input features, are promising directions for future research.

540 References

548

549

550

559

560

561 562

563

565

566

567

- 542 Kwangjun Ahn, Xiang Cheng, Hadi Daneshmand, and Suvrit Sra. Transformers learn to implement 543 preconditioned gradient descent for in-context learning. *Advances in Neural Information Processing* 544
- 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 Timalsina, 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.
- 571 Chen Dun, Mirian Hipolito Garcia, Guoqing Zheng, Ahmed Hassan Awadallah, Anastasios Kyrillidis, and Robert Sim. Sweeping heterogeneity with smart mops: Mixture of prompts for Ilm task adaptation. *arXiv preprint arXiv:2310.02842*, 2023.
- 574 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.
- 577 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.
- 583 Khashayar Gatmiry, Nikunj Saunshi, Sashank J Reddi, Stefanie Jegelka, and Sanjiv Kumar. Can
 584 looped transformers learn to implement multi-step gradient descent for in-context learning? In
 585 Forty-first International Conference on Machine Learning.
- Riccardo Grazzi, Julien Siems, Simon Schrodi, 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.
- 593 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.

594 595 596	Angelos Katharopoulos, Apoorv Vyas, Nikolaos Pappas, and François Fleuret. Transformers are rnns: Fast autoregressive transformers with linear attention. In <i>International conference on machine</i> <i>learning</i> , pp. 5156–5165. PMLR, 2020.					
597 598	Tobias Katsch. Gateloop: Fully data-controlled linear recurrence for sequence modeling. <i>arXiv</i> preprint arXiv:2311.01927, 2023.					
599 600 601	Yingcong Li, Muhammed Emrullah Ildiz, Dimitris Papailiopoulos, and Samet Oymak. Transformers as algorithms: Generalization and stability in in-context learning. In <i>International Conference on Machine Learning</i> , pp. 19565–19594. PMLR, 2023.					
602 603 604	Yingcong Li, Ankit Singh Rawat, and Samet Oymak. Fine-grained analysis of in-context linear estimation: Data, architecture, and beyond. <i>arXiv preprint arXiv:2407.10005</i> , 2024.					
605 606	Ziqian Lin and Kangwook Lee. Dual operating modes of in-context learning. arXiv preparXiv:2402.18819, 2024.					
607 608 609	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. <i>arXiv preprint arXiv:2307.03576</i> , 2023.					
610 611 612	Sewon Min, Xinxi Lyu, Ari Holtzman, Mikel Artetxe, Mike Lewis, Hannaneh Hajishirzi, and Luke Zettlemoyer. Rethinking the role of demonstrations: What makes in-context learning work? <i>arXiv</i> preprint arXiv:2202.12837, 2022.					
613 614 615	Antonio Orvieto, Samuel L Smith, Albert Gu, Anushan Fernando, Caglar Gulcehre, Razvan Pascanu, and Soham De. Resurrecting recurrent neural networks for long sequences. In <i>International</i> <i>Conference on Machine Learning</i> , pp. 26670–26698. PMLR, 2023.					
616 617 618	Jongho Park, Jaeseung Park, Zheyang Xiong, Nayoung Lee, Jaewoong Cho, Samet Oymak, Kang- wook Lee, and Dimitris Papailiopoulos. Can mamba learn how to learn? a comparative study on in-context learning tasks. <i>arXiv preprint arXiv:2402.04248</i> , 2024.					
619 620 621	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. <i>arXiv preprint arXiv:2305.13048</i> , 2023.					
622 623 624	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. <i>arXiv preprint arXiv:2404.05892</i> , 2024.					
625 626	Hao Peng, Nikolaos Pappas, Dani Yogatama, Roy Schwartz, Noah A Smith, and Lingpeng Kong. Random feature attention. <i>arXiv preprint arXiv:2103.02143</i> , 2021.					
627 628	Zhen Qin, Songlin Yang, Weixuan Sun, Xuyang Shen, Dong Li, Weigao Sun, and Yiran Zhong. Hgrn2: Gated linear rnns with state expansion. <i>arXiv preprint arXiv:2404.07904</i> , 2024.					
629 630 631	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. <i>arXiv preprint arXiv:2405.15731</i> , 2024.					
633 634	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. <i>arXiv preprint arXiv:2307.08621</i> , 2023.					
636 637	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. <i>arXiv preprint arXiv:2405.05254</i> , 2024.					
639	A Vaswani. Attention is all you need. Advances in Neural Information Processing Systems, 2017.					
640 641 642	Johannes Von Oswald, Eyvind Niklasson, Ettore Randazzo, João Sacramento, Alexander Mordvintsev, Andrey Zhmoginov, and Max Vladymyrov. Transformers learn in-context by gradient descent. In <i>International Conference on Machine Learning</i> , pp. 35151–35174. PMLR, 2023.					
643 644 645	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. <i>arXiv preprint arXiv:2309.05858</i> , 2023.					
646 647	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. <i>arXiv preprint arXiv:2407.00256</i> , 2024.					

648 649 650	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? <i>arXiv preprint</i> <i>arXiv:2310.08391</i> , 2023.					
651 652	Sang Michael Xie, Aditi Raghunathan, Percy Liang, and Tengyu Ma. An explanation of in-context learning as implicit bayesian inference. In <i>International Conference on Learning Representations</i>					
653						
654	Songlin Yang, Bailin Wang, Yikang Shen, Rameswar Panda, and Yoon Kim. Gated linear attention transformers with hardware-efficient training. <i>arXiv preprint arXiv:2312.06635</i> , 2023.					
655	Duici Zhang, Spanger Ersi, and Dater I. Dortlatt. Trained transformers lager linear models in contact					
656 657	arXiv preprint arXiv:2306.09927, 2023.					
658 659	Ruiqi Zhang, Spencer Frei, and Peter L Bartlett. Trained transformers learn linear models in-context. <i>Journal of Machine Learning Research</i> , 25(49):1–55, 2024.					
660 661	Itamar Zimerman, Ameen Ali, and Lior Wolf. A unified implicit attention formulation for gated-linear recurrent sequence models. <i>arXiv preprint arXiv:2405.16504</i> , 2024.					
662						
663	Chang Zong, Jian Shao, Weiming Lu, and Yueting Zhuang. Stock movement prediction with multimodal stable fusion via gated cross-attention mechanism. <i>arXiv preprint arXiv:2406.06594</i> , 2024					
004	2024.					
665						
666						
667						
668						
669						
670						
671						
672						
673						
674						
675						
676						
670						
670						
600						
691						
682						
683						
684						
685						
686						
687						
688						
689						
690						
691						
692						
693						
694						
695						
696						
697						
698						
699						
700						
701						

7	0	2
7	0	3
7	0	4
7	0	5
7	0	6
7	0	7
7	0	8
7	0	9
7	1	0
7	1	1
7	1	2
7	1	3
7	1	4
7	1	5
7	1	6
7	1	7
7	1	8
7	1	9
7	2	0
7	2	1
7	2	2
7	2	3
7	2	4
7	2	5
7	2	6
7	2	7
7	2	8
7	2	9
7	3	0
7	3	1
7	3	2
7	3	3
7	3	4

CONTENTS A Implementation Detail **B GLA** \Leftrightarrow **WPGD** C Optimization Landscape of WPGD D Loss Landscape of 1-layer GLA

IMPLEMENTATION DETAIL А

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 $(r_1, r_2) := \left(\mathbb{E}[\beta_1^\top \beta]/d, \mathbb{E}[\beta_2^\top \beta]/d\right)$ denote the correlations between in-context tasks β_1, β_2 and query task β . We generate task vectors as follows:

$$\boldsymbol{\beta}_1, \boldsymbol{\beta}_2 \sim \mathcal{N}(0, \boldsymbol{I}_d) \text{ and } \boldsymbol{\beta} \sim \mathcal{N}(r_1\boldsymbol{\beta}_1 + r_2\boldsymbol{\beta}_2, (1 - r_1^2 - r_2^2)\boldsymbol{I}_d).$$

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

Implementation setting. We train 1-layer linear attention and GLA models for solving multiprompt ICL problem as described in Section 5.1. For GLA model, we consider sigmoid-type gating function given by scalar gating: $g(z) = \phi(W_g^{\top} z) \mathbf{1} \mathbf{1}^{\top}$, or vector gating: $g(z) = \phi(W_g z) \mathbf{1}^{\top}$ where $\phi(z) = (1 + e^{-z})^{-1}$ 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

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

Then dotted curves display our theoretical results derive using $\Sigma = I$ and R, r above. Specifically, in Figure 1, black dashed curves represent $\mathcal{L}_{WPGD}^{\star}$ following (18) and blues dashed curves represent $\mathcal{L}_{GLA}^{\star}$ 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.

• 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 ($\mathcal{L}_{wpGD}^{\star}$, 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}_{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

⁷⁸² In this section, we present additional experiments on multi-layer GLA models. We adopt the same ⁷⁸³ experimental setup as described in Figure 1a and Appendix A, with parameters set to $(r_1, r_2) = (0, 1)$. ⁷⁸⁴ The results are displayed in Figure 2, where the blue, red, and green curves correspond to the ⁷⁸⁵ performance of one-, two-, and three-layer GLA models, respectively, with the *y*-axis presented in ⁷⁸⁶ log-scale. According to Theorem 2, an *L*-layer GLA performs *L* steps of WPGD, suggesting that ⁷⁸⁷ 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

$$\mathbf{Z} = [\mathbf{z}_1 \cdots \mathbf{z}_n \ \mathbf{z}_{n+1}]^{\top} = \begin{bmatrix} \mathbf{x}_1 \cdots \mathbf{x}_n & \mathbf{x}_{n+1} \\ y_1 & \cdots & y_n & 0 \end{bmatrix}^{\top}$$

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

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

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

$$\boldsymbol{o}_i = \boldsymbol{S}_i \boldsymbol{q}_i$$
 and $\boldsymbol{S}_i = \boldsymbol{G}_i \odot \boldsymbol{S}_{i-1} + \boldsymbol{v}_i \boldsymbol{k}_i^{\mathsf{T}}$

$$=\sum_{j=1}^{l}\boldsymbol{G}_{j:i}\odot\boldsymbol{v}_{j}\boldsymbol{k}_{j}^{\mathsf{T}}$$

where we define

$$G_{j:i} = G_{j+1} \odot G_{j+2} \cdots G_i, \quad j < i, \text{ and } G_{i:i} = 11^{+}.$$

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

$$\boldsymbol{o}_{n+1} = \boldsymbol{S}_{n+1}\boldsymbol{q}_{n+1}$$
 and $\boldsymbol{S}_{n+1} = \sum_{j=1}^{n+1} \boldsymbol{G}_{j:n+1} \odot \boldsymbol{v}_j \boldsymbol{k}_j^{\mathsf{T}}.$

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

$$\boldsymbol{W}_{k} = \begin{bmatrix} \boldsymbol{P}_{k}^{\top} & \boldsymbol{0} \\ \boldsymbol{0}^{\top} & \boldsymbol{0} \end{bmatrix}, \quad \boldsymbol{W}_{q} = \begin{bmatrix} \boldsymbol{P}_{q}^{\top} & \boldsymbol{0} \\ \boldsymbol{0}^{\top} & \boldsymbol{0} \end{bmatrix} \quad \text{and} \quad \boldsymbol{W}_{v} = \begin{bmatrix} \boldsymbol{0}_{d \times d} & \boldsymbol{0} \\ \boldsymbol{0}^{\top} & \boldsymbol{1} \end{bmatrix}.$$
(26)

Then, given each token $z_i = [x_i^{\top} y_i]^{\top}, i \in [n]$, single-layer GLA returns

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

and we obtain

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

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

$$\boldsymbol{G}_{i} = \begin{bmatrix} * & * \\ \boldsymbol{g}_{i}^{\top} & * \end{bmatrix}$$
 and $\boldsymbol{G}_{j:i} = \begin{bmatrix} * & * \\ \boldsymbol{g}_{j:i}^{\top} & * \end{bmatrix}$

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

Combing all together, and letting $X = [x_1 \ x_2 \ \cdots \ x_n]^\top$ and $y = [y_1 \ y_2 \ \cdots \ y_n]^\top$, we obtain

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

832 where

$$\mathbf{\Omega} = \begin{bmatrix} \boldsymbol{g}_{1:n+1} & \boldsymbol{g}_{2:n+1} & \cdots & \boldsymbol{g}_{n:n+1} \end{bmatrix} \in \mathbb{R}^{n \times d}.$$

Then if taking the last entry of o_{n+1} as final prediction, we get

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

836837 It completes the proof of Theorem 1.

Construction 2: Based on the construction given in (26), only *d* elements of G_i matrices are useful. One might ask about the effect of other entries of G_i . Therefore, in the following, we introduce an other model construction showing that different row of G_i implements WPGD with different weighting. Similarly, let W_k , W_q be the same as (26) but with W_v constructed by

$$\boldsymbol{W}_{v} = \begin{bmatrix} \boldsymbol{0}_{(d+1) \times d} & \boldsymbol{u} \end{bmatrix}$$
 where $\boldsymbol{u} = \begin{bmatrix} u_{1} & u_{2} & \cdots & u_{d+1} \end{bmatrix}^{\top} \in \mathbb{R}^{d+1}$

Then the value embeddings have the form of $v_i = y_i u$, which gives

$$\boldsymbol{v}_i \boldsymbol{k}_i^{\mathsf{T}} = \boldsymbol{u} \begin{bmatrix} y_i \boldsymbol{x}_i^{\mathsf{T}} \boldsymbol{P}_k & \boldsymbol{0} \end{bmatrix}.$$

Next, let

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

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

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

where

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

Therefore, consider (d + 1)-dimensional output o_{n+1} . Each entry implements a 1-step WPGD with same preconditioners P_k , P_q and different weighting matrices Ω '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 o_{n+1} as prediction, it returns the same result as Construction 1 above, where only last rows of G_i 's are useful. Additionally, suppose that the final prediction \hat{y} is given after a linear head \boldsymbol{h} , that is, $\hat{y} = \boldsymbol{h}^{\top} \boldsymbol{o}_{n+1}$, and let $\boldsymbol{h} = [h_1 \ h_2 \ \cdots \ h_{d+1}]^{\top} \in \mathbb{R}^{d+1}$. Then

$$\hat{\mathbf{y}} = \boldsymbol{h}^{\mathsf{T}} \boldsymbol{o}_{n+1} = \boldsymbol{x}^{\mathsf{T}} \boldsymbol{P}_q \left(\boldsymbol{X} \boldsymbol{P}_k \odot \bar{\boldsymbol{\Omega}} \right)^{\mathsf{T}} \boldsymbol{y}$$
(27)

where

$$\bar{\boldsymbol{\Omega}} = \sum_{i=1}^{d+1} h_i \boldsymbol{\Omega}_i = \sum_{i=1}^{d+1} h_i u_i \begin{bmatrix} \boldsymbol{g}_{1:n+1}^i & \boldsymbol{g}_{2:n+1}^i & \cdots & \boldsymbol{g}_{n:n+1}^i \end{bmatrix} \in \mathbb{R}^{n \times d}.$$
(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 ℓ 'th layer parameterized by $P_{k,\ell}, P_{q,\ell} \in \mathbb{R}^{d \times d}$ as in (8) and with corresponding gating vectors $g_{i,i}^{\ell}$ $i \in [n+1], \ell \in [L]$. Let $\hat{y}_{i,\ell}$ be the (d+1)'th entry of the i'th token of the ℓ 'th layer input (or $(\ell-1)$ 'th layer output after residual). Additionally, denote $\Omega_{\ell} = [g_{1:n+1}^{\ell} \cdots g_{n:n+1}^{\ell}]^{\top}$ and $\bar{X}_{\ell} = XP_{k,\ell} \odot \Omega_{\ell}$. Let $B_{\ell} = [\beta_{1,\ell} \cdots \beta_{n,\ell}]^{\top}$ where $\beta_{i,0} = 0$ for $i \in [n+1]$ and $M_i = \begin{bmatrix} I_i & 0\\ 0 & 0 \end{bmatrix} \in \mathbb{R}^{n \times n}$. Then it satisfies that for

•
$$i \leq n, \ \hat{y}_{i,\ell} = y_i - \boldsymbol{x}_i^{\mathsf{T}} \boldsymbol{\beta}_{i,\ell-1} \text{ where } \boldsymbol{\beta}_{i,\ell} = \boldsymbol{\beta}_{i,\ell-1} + \boldsymbol{P}_{q,\ell} \left(\nabla_{i,\ell} \oslash \boldsymbol{g}_{i:n+1}^{\ell} \right)$$

• and
$$\hat{y}_{n+1,\ell} = \mathbf{x}^{\top} \boldsymbol{\beta}_{\ell-1}$$
 where $\boldsymbol{\beta}_{\ell} = (1 + \alpha_{\ell}) \boldsymbol{\beta}_{\ell-1} + \boldsymbol{P}_{q,\ell} \left(\nabla_{n,\ell} \oslash \boldsymbol{g}_{n+1}^{\ell} \right)$ and $\alpha_{\ell} = \mathbf{x}^{\top} \boldsymbol{P}_{q,\ell} \boldsymbol{P}_{k,\ell}^{\top} \mathbf{x}$.

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

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

$$\boldsymbol{S}_{i} = \begin{bmatrix} \boldsymbol{0}_{d \times d} & \boldsymbol{0} \\ \sum_{j=1}^{i} y_{j} \boldsymbol{x}_{j}^{\mathsf{T}} \boldsymbol{P}_{k} \odot \boldsymbol{g}_{j:i}^{\mathsf{T}} & \boldsymbol{0} \end{bmatrix}$$

Additionally, recap that we have

$$M_i = \begin{bmatrix} I_i & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}$$
 and $\Omega = \begin{bmatrix} g_{i:n+1} & \cdots & g_{n:n+1} \end{bmatrix}$.

Let ⊘ denote Hadamard division. Then

$$\sum_{j=1}^{i} y_j \boldsymbol{P}_k^{\mathsf{T}} \boldsymbol{x}_j \odot \boldsymbol{g}_{j:i} = \left(\sum_{j=1}^{i} y_j \boldsymbol{P}_k^{\mathsf{T}} \boldsymbol{x}_j \odot \boldsymbol{g}_{j:n+1} \right) \oslash \boldsymbol{g}_{i:n+1}$$
$$= (\boldsymbol{X} \boldsymbol{P}_k \odot \boldsymbol{\Omega})^{\mathsf{T}} \boldsymbol{M}_i \boldsymbol{y} \oslash \boldsymbol{g}_{i:n+1},$$

Therefore,

$$\boldsymbol{o}_{i} = \boldsymbol{S}_{i}\boldsymbol{q}_{i} = \begin{bmatrix} \boldsymbol{0}_{d \times d} & \boldsymbol{0} \\ ((\boldsymbol{X}\boldsymbol{P}_{k} \odot \boldsymbol{\Omega})^{\top}\boldsymbol{M}_{i}\boldsymbol{y} \oslash \boldsymbol{g}_{i:n+1})^{\top} & \boldsymbol{0} \end{bmatrix} \begin{bmatrix} \boldsymbol{P}_{q}^{\top}\boldsymbol{x}_{i} \\ \boldsymbol{0} \end{bmatrix} = \begin{bmatrix} \boldsymbol{0} \\ \boldsymbol{x}_{i}^{\top}\boldsymbol{P}_{q}\left(\bar{\boldsymbol{X}}^{\top}\boldsymbol{M}_{i}\boldsymbol{y} \oslash \boldsymbol{g}_{i:n+1}\right) \end{bmatrix}.$$
(29)

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

$$\boldsymbol{o}_{n+1} = \begin{bmatrix} \boldsymbol{0} \\ \boldsymbol{x}^{\top} \boldsymbol{P}_q \left(\bar{\boldsymbol{X}}^{\top} \boldsymbol{y} \oslash \boldsymbol{g}_{n+1} \right) \end{bmatrix}.$$
(30)

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

912
913
914

$$\mathbf{Z}_{\ell} = \begin{bmatrix} z_{1,\ell} & \cdots & z_{n,\ell} & z_{n+1,\ell} \end{bmatrix}^{\top} \in \mathbb{R}^{(n+1)\times(d+1)},$$

 $\mathbf{O}_{\ell} = \begin{bmatrix} \boldsymbol{o}_{1,\ell} & \cdots & \boldsymbol{o}_{n,\ell} & \boldsymbol{o}_{n+1,\ell} \end{bmatrix}^{\top} \in \mathbb{R}^{(n+1)\times(d+1)},$

where $Z_1 = Z$. Then, given the residual connection of each layer, the input of $(\ell + 1)$ 'th layer is given by 917 (21)

$$\mathbf{Z}_{\ell+1} = \mathbf{Z}_{\ell} + \mathbf{O}_{\ell}.\tag{31}$$

918 Note that $Z_{\ell+1}$ is also the output of ℓ 'th layer after residual. Recall (29) which implies that the first *d* 919 dimension of the output o_i for all tokens $i \in [n + 1]$ is zero. Therefore, the first *d* dimension of $z_{i,\ell}$ 920 keeps the same as x_i and let us write the input of ℓ 'th layer (also the output of $(\ell - 1)$ 'th layer after 921 residual) as

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

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

$$\boldsymbol{G}_{i}^{\ell} = \begin{bmatrix} \ast & \ast \\ (\boldsymbol{g}_{i}^{\ell})^{\mathsf{T}} & \ast \end{bmatrix}$$

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

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

where $\bar{X}_{\ell} := X P_{k,\ell} \odot \Omega_{\ell}$ and

$$\hat{\boldsymbol{y}}_{\ell} = \begin{bmatrix} \hat{y}_{1,\ell} & \cdots & \hat{y}_{n,\ell} \end{bmatrix}^{\mathsf{T}} \in \mathbb{R}^{n}, \\ \boldsymbol{\Omega}_{\ell} = \begin{bmatrix} \boldsymbol{g}_{1:n+1}^{\ell} & \boldsymbol{g}_{2:n+1}^{\ell} & \cdots & \boldsymbol{g}_{n:n+1}^{\ell} \end{bmatrix}^{\mathsf{T}} \in \mathbb{R}^{n \times d}.$$

Following the residual connection as in (31), we have $z_{i\ell+1} = z_{i\ell} + o_{i\ell}$ and hence

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

Now consider the algorithm given in the theorem statement where $\hat{y}_{i,\ell} = y_i - x_i^\top \beta_{i,\ell-1}$ and $\beta_{i,\ell} = \beta_{i,\ell-1} + P_{q,\ell} \left(\bar{X}_{\ell}^\top M_i((X \odot B_{\ell-1})\mathbf{1} - \mathbf{y}) \oslash g_{i;n+1}^\ell \right)$, which gives

$$\hat{y}_{i,\ell+1} = y_i - \boldsymbol{x}_i^{\top} \boldsymbol{\beta}_{i,\ell}$$

$$\hat{y}_{i,\ell} = y_i - \boldsymbol{x}_i^{\top} \boldsymbol{\beta}_{i,\ell-1}.$$
(34)

Then

$$\hat{y}_{i,\ell+1} - \hat{y}_{i,\ell} = -\boldsymbol{x}_i^{\top} (\boldsymbol{\beta}_{i,\ell} - \boldsymbol{\beta}_{i,\ell-1})
= -\boldsymbol{x}_i^{\top} \boldsymbol{P}_{q,\ell} \left(\bar{\boldsymbol{X}}_{\ell}^{\top} \boldsymbol{M}_i ((\boldsymbol{X} \odot \boldsymbol{B}_{\ell-1}) \boldsymbol{1} - \boldsymbol{y}) \oslash \boldsymbol{g}_{i:n+1}^{\ell} \right)
= \boldsymbol{x}_i^{\top} \boldsymbol{P}_{q,\ell} \left(\bar{\boldsymbol{X}}_{\ell}^{\top} \boldsymbol{M}_i \hat{\boldsymbol{y}}_{\ell} \oslash \boldsymbol{g}_{i:n+1}^{\ell} \right).$$
(35)

The last equation uses (34), that

$$(\boldsymbol{X} \odot \boldsymbol{B}_{\ell-1})\boldsymbol{1} = [\boldsymbol{x}_1^{\mathsf{T}} \boldsymbol{\beta}_{1,\ell} \cdots \boldsymbol{x}_n^{\mathsf{T}} \boldsymbol{\beta}_{n,\ell}]^{\mathsf{T}} \implies (\boldsymbol{X} \odot \boldsymbol{B}_{\ell-1})\boldsymbol{1} - \boldsymbol{y} = -\hat{\boldsymbol{y}}_{\ell}.$$
 (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 ℓ 'th layer

$$\boldsymbol{o}_{n+1,\ell} = \begin{bmatrix} \boldsymbol{0} \\ \boldsymbol{x}^{\mathsf{T}} \boldsymbol{P}_{q,\ell} \left(\bar{\boldsymbol{X}}_{\ell}^{\mathsf{T}} \hat{\boldsymbol{y}}_{\ell} \oslash \boldsymbol{g}_{n+1}^{\ell} \right) \end{bmatrix} + \begin{bmatrix} \boldsymbol{0} \\ \boldsymbol{x}^{\mathsf{T}} \boldsymbol{P}_{q,\ell} \boldsymbol{P}_{k,\ell}^{\mathsf{T}} \boldsymbol{x} \cdot \hat{\boldsymbol{y}}_{n+1,\ell} \end{bmatrix}$$

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

963 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

$$\hat{y}_{n+1,\ell+1} = \hat{y}_{n+1,\ell} + \boldsymbol{x}^{\top} \boldsymbol{P}_{q,\ell} \left(\bar{\boldsymbol{X}}_{\ell}^{\top} \hat{\boldsymbol{y}}_{\ell} \otimes \boldsymbol{g}_{n+1}^{\ell} \right) + \alpha_{\ell} \cdot \hat{y}_{n+1,\ell} = (1 + \alpha_{\ell}) \hat{y}_{n+1,\ell} + \boldsymbol{x}^{\top} \boldsymbol{P}_{q,\ell} \left(\bar{\boldsymbol{X}}_{\ell}^{\top} \hat{\boldsymbol{y}}_{\ell} \otimes \boldsymbol{g}_{n+1}^{\ell} \right).$$
(37)

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

970
$$\hat{y}_{n+1,\ell+1} = -\boldsymbol{x}^{\mathsf{T}} \boldsymbol{\beta}_{n+1,\ell}$$
971
$$\hat{y}_{n+1,\ell+1} = -\boldsymbol{x}^{\mathsf{T}} \boldsymbol{\beta}_{n+1,\ell}$$

$$\hat{y}_{n+1,\ell} = -\boldsymbol{x}^{\mathsf{T}} \boldsymbol{\beta}_{n+1,\ell-1}.$$

Then $\hat{y}_{n+1\,\ell+1} - (1+\alpha_{\ell})\hat{y}_{n+1\,\ell} = -\boldsymbol{x}^{\top} \left(\boldsymbol{\beta}_{n+1,\ell} - (1+\alpha_{\ell})\boldsymbol{\beta}_{n+1,\ell-1}\right)$ $= -\mathbf{x}^{\top} \mathbf{P}_{q,\ell} \left(\bar{\mathbf{X}}_{\ell}^{\top} ((\mathbf{X} \odot \mathbf{B}_{\ell-1}) \mathbf{1} - \mathbf{y}) \oslash \mathbf{g}_{n+1}^{\ell} \right)$ $= \boldsymbol{x}^{\top} \boldsymbol{P}_{q,\ell} \left(\bar{\boldsymbol{X}}_{\ell}^{\top} \hat{\boldsymbol{y}}_{\ell} \otimes \boldsymbol{g}_{n+1}^{\ell} \right)$ which is the same as (37) by using the fact from (36). С **OPTIMIZATION LANDSCAPE OF WPGD** C.1 **PROOF OF THEOREM 3** *Proof.* Recapping the objective from (3) and following Definition 2, we have $\mathcal{L}(\boldsymbol{P},\boldsymbol{\omega}) = \mathbb{E}\left| \left(\boldsymbol{y} - \boldsymbol{x}^{\top} \boldsymbol{P} \boldsymbol{X}(\boldsymbol{\omega} \odot \boldsymbol{y}) \right)^2 \right|$ $= \mathbb{E}\left[y^{2}\right] - 2\mathbb{E}\left[y\boldsymbol{x}^{\top}\boldsymbol{P}\boldsymbol{X}(\boldsymbol{\omega}\odot\boldsymbol{y})\right] + \mathbb{E}\left[\left(\boldsymbol{x}^{\top}\boldsymbol{P}\boldsymbol{X}(\boldsymbol{\omega}\odot\boldsymbol{y})\right)^{2}\right].$ Let $y = \mathbf{x}^{\mathsf{T}} \boldsymbol{\beta} + \boldsymbol{\xi}$ and $y_i = \mathbf{x}_i^{\mathsf{T}} \boldsymbol{\beta}_i + \boldsymbol{\xi}_i$, for $i \in [n]$, where $\boldsymbol{\xi}, \boldsymbol{\xi}_i \sim \mathcal{N}(0, \sigma^2)$ are i.i.d. Then, $\mathbb{E}[v^2] = \mathbb{E}[(\boldsymbol{x}^{\mathsf{T}}\boldsymbol{\beta} + \boldsymbol{\xi})^2] = \operatorname{tr}(\boldsymbol{\Sigma}) + \sigma^2.$ and $\mathbb{E}\left[y\mathbf{x}^{\top} \mathbf{P} \mathbf{X}(\boldsymbol{\omega} \odot \mathbf{y})\right] = \mathbb{E}\left[(\boldsymbol{\beta}^{\top} \mathbf{x} + \boldsymbol{\xi})\mathbf{x}^{\top} \mathbf{P} \sum_{i=1}^{n} \omega_{i} \mathbf{x}_{i}(\mathbf{x}_{i}^{\top} \boldsymbol{\beta}_{i} + \boldsymbol{\xi}_{i})\right]$ $= \mathbb{E} \left| \boldsymbol{\beta}^{\mathsf{T}} \boldsymbol{x} \boldsymbol{x}^{\mathsf{T}} \boldsymbol{P} \sum_{i=1}^{n} \omega_{i} \boldsymbol{x}_{i} \boldsymbol{x}_{i}^{\mathsf{T}} \boldsymbol{\beta}_{i} \right|$ $= \operatorname{tr}\left(\boldsymbol{\Sigma}\boldsymbol{P}\boldsymbol{\Sigma}\sum_{i=1}^{n}\omega_{i}\mathbb{E}\left[\boldsymbol{\beta}_{i}\boldsymbol{\beta}^{\mathsf{T}}\right]\right)$ $= \operatorname{tr} \left(\Sigma^2 \boldsymbol{P} \right) \boldsymbol{\omega}^{\mathsf{T}} \boldsymbol{r}.$ Here, the last equality comes from the fact that since $\beta_i - r_{ij}\beta_j$ is independent of β_j for $i, j \in [n+1]$ following Definition 1, we have $\mathbb{E}[\beta_i\beta^{\top}] = r_{i,n+1}I_d$ and $\sum_{i=1}^{n} \omega_i \mathbb{E}[\beta_i\beta^{\top}]$ returns $\omega^{\dagger} \mathbf{r} \cdot \mathbf{I}_d$. Hence, $\mathbb{E}\left[\left(\boldsymbol{x}^{\top}\boldsymbol{P}\boldsymbol{X}(\boldsymbol{\omega}\odot\boldsymbol{y})\right)^{2}\right] = \mathbb{E}\left[\boldsymbol{x}^{\top}\boldsymbol{P}\left(\sum_{i=1}^{n}\omega_{i}(\boldsymbol{x}_{i}^{\top}\boldsymbol{\beta}_{i}+\boldsymbol{\xi}_{i})\boldsymbol{x}_{i}\right)\left(\sum_{i=1}^{n}\omega_{i}\boldsymbol{x}_{i}^{\top}(\boldsymbol{x}_{i}^{\top}\boldsymbol{\beta}_{i}+\boldsymbol{\xi}_{i})\right)\boldsymbol{P}^{\top}\boldsymbol{x}\right]$ $= \operatorname{tr}\left(\boldsymbol{P}^{\mathsf{T}}\boldsymbol{\Sigma}\boldsymbol{P}\mathbb{E}\left|\sum_{i=1}^{n}\omega_{i}^{2}(\boldsymbol{x}_{i}^{\mathsf{T}}\boldsymbol{\beta}_{i}+\xi_{i})^{2}\boldsymbol{x}_{i}\boldsymbol{x}_{i}^{\mathsf{T}}\right|\right)$ + tr $\left(\boldsymbol{P}^{\mathsf{T}} \boldsymbol{\Sigma} \boldsymbol{P} \mathbb{E} \left| \sum_{i=1}^{n} \omega_{i} \omega_{j} (\boldsymbol{x}_{i}^{\mathsf{T}} \boldsymbol{\beta}_{i} + \xi_{i}) \boldsymbol{x}_{i} \boldsymbol{x}_{j}^{\mathsf{T}} (\boldsymbol{x}_{j}^{\mathsf{T}} \boldsymbol{\beta}_{j} + \xi_{j}) \right| \right)$ where $\operatorname{tr}\left(\boldsymbol{P}^{\mathsf{T}}\boldsymbol{\Sigma}\boldsymbol{P}\mathbb{E}\left[\sum_{i=1}^{n}\omega_{i}^{2}(\boldsymbol{x}_{i}^{\mathsf{T}}\boldsymbol{\beta}_{i}+\boldsymbol{\xi}_{i})^{2}\boldsymbol{x}_{i}\boldsymbol{x}_{i}^{\mathsf{T}}\right]\right)=\operatorname{tr}\left(\boldsymbol{P}^{\mathsf{T}}\boldsymbol{\Sigma}\boldsymbol{P}\mathbb{E}\left[\sum_{i=1}^{n}\omega_{i}^{2}(\boldsymbol{x}_{i}^{\mathsf{T}}\boldsymbol{\beta}_{i}\boldsymbol{\beta}_{i}^{\mathsf{T}}\boldsymbol{x}_{i}+\sigma^{2})\boldsymbol{x}_{i}\boldsymbol{x}_{i}^{\mathsf{T}}\right]\right)$ $= \|\boldsymbol{\omega}\|_{\ell_{2}}^{2} \operatorname{tr} \left(\boldsymbol{P}^{\mathsf{T}} \boldsymbol{\Sigma} \boldsymbol{P} \left(\mathbb{E} \left[\boldsymbol{x} \boldsymbol{x}^{\mathsf{T}} \boldsymbol{x} \boldsymbol{x}^{\mathsf{T}} \right] + \sigma^{2} \boldsymbol{\Sigma} \right) \right)$ $= \|\omega\|_{\ell_2}^2 \left(\operatorname{tr} \left(\Sigma P^{\mathsf{T}} \Sigma P \right) \left(\operatorname{tr} \left(\Sigma \right) + \sigma^2 \right) + \operatorname{tr} \left(\Sigma^2 P^{\mathsf{T}} \Sigma P \right) \right),$ and $\operatorname{tr}\left(\boldsymbol{P}^{\mathsf{T}}\boldsymbol{\Sigma}\boldsymbol{P}\mathbb{E}\left[\sum_{i=1}^{n}\omega_{i}\omega_{j}(\boldsymbol{x}_{i}^{\mathsf{T}}\boldsymbol{\beta}_{i}+\xi_{i})\boldsymbol{x}_{i}\boldsymbol{x}_{j}^{\mathsf{T}}(\boldsymbol{x}_{j}^{\mathsf{T}}\boldsymbol{\beta}_{j}+\xi_{j})\right]\right)=\operatorname{tr}\left(\boldsymbol{P}^{\mathsf{T}}\boldsymbol{\Sigma}\boldsymbol{P}\mathbb{E}\left[\sum_{i=1}^{n}\omega_{i}\omega_{j}\boldsymbol{x}_{i}\boldsymbol{x}_{i}^{\mathsf{T}}\boldsymbol{\beta}_{i}\boldsymbol{\beta}_{j}^{\mathsf{T}}\boldsymbol{x}_{j}\boldsymbol{x}_{j}^{\mathsf{T}}\right]\right)$ $= \operatorname{tr} \left(\Sigma^2 \boldsymbol{P}^{\mathsf{T}} \Sigma \boldsymbol{P} \right) \boldsymbol{\omega}^{\mathsf{T}} \boldsymbol{R} \boldsymbol{\omega}.$

1026 Combining all together and letting $M := tr(\Sigma) + \sigma^2$, we obtain

$$\mathcal{L}(\boldsymbol{P},\boldsymbol{\omega}) = \boldsymbol{M} - 2\operatorname{tr}\left(\boldsymbol{\Sigma}^{2}\boldsymbol{P}\right)\boldsymbol{\omega}^{\mathsf{T}}\boldsymbol{r} + \boldsymbol{M} \|\boldsymbol{\omega}\|_{\ell_{2}}^{2} \operatorname{tr}\left(\boldsymbol{\Sigma}\boldsymbol{P}^{\mathsf{T}}\boldsymbol{\Sigma}\boldsymbol{P}\right) + (\|\boldsymbol{\omega}\|_{\ell_{2}}^{2} + \boldsymbol{\omega}^{\mathsf{T}}\boldsymbol{R}\boldsymbol{\omega})\operatorname{tr}\left(\boldsymbol{\Sigma}^{2}\boldsymbol{P}^{\mathsf{T}}\boldsymbol{\Sigma}\boldsymbol{P}\right).$$
(38)

For simplicity, and without loss of generality, let

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

Then, we obtain

$$\mathcal{L}(\tilde{\boldsymbol{P}},\omega) = M - 2\operatorname{tr}\left(\boldsymbol{\Sigma}\tilde{\boldsymbol{P}}\right)\omega^{\mathsf{T}}\boldsymbol{r} + M \|\boldsymbol{\omega}\|_{\ell_{2}}^{2}\operatorname{tr}\left(\tilde{\boldsymbol{P}}^{\mathsf{T}}\tilde{\boldsymbol{P}}\right) + (\|\boldsymbol{\omega}\|_{\ell_{2}}^{2} + \boldsymbol{\omega}^{\mathsf{T}}\boldsymbol{R}\omega)\operatorname{tr}\left(\boldsymbol{\Sigma}\tilde{\boldsymbol{P}}^{\mathsf{T}}\tilde{\boldsymbol{P}}\right).$$

$$\tag{40}$$

Further, the gradients can be written as

$$\nabla_{\tilde{\boldsymbol{P}}} \mathcal{L}(\tilde{\boldsymbol{P}}, \boldsymbol{\omega}) = -2\boldsymbol{\omega}^{\mathsf{T}} \boldsymbol{r} \boldsymbol{\Sigma} + 2M \|\boldsymbol{\omega}\|_{\ell_2}^2 \, \tilde{\boldsymbol{P}} + 2(\|\boldsymbol{\omega}\|_{\ell_2}^2 + \boldsymbol{\omega}^{\mathsf{T}} \boldsymbol{R} \boldsymbol{\omega}) \boldsymbol{\Sigma} \tilde{\boldsymbol{P}},\tag{41}$$

$$\nabla_{\omega} \mathcal{L}(\tilde{\boldsymbol{P}}, \boldsymbol{\omega}) = -2 \operatorname{tr}\left(\boldsymbol{\Sigma} \tilde{\boldsymbol{P}}\right) \boldsymbol{r} + 2M \operatorname{tr}\left(\tilde{\boldsymbol{P}}^{\top} \tilde{\boldsymbol{P}}\right) \boldsymbol{\omega} + 2 \operatorname{tr}\left(\boldsymbol{\Sigma} \tilde{\boldsymbol{P}}^{\top} \tilde{\boldsymbol{P}}\right) (\boldsymbol{I}_{n} + \boldsymbol{R}) \boldsymbol{\omega}.$$
(42)

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

$$\tilde{\boldsymbol{P}} = \left(M \|\boldsymbol{\omega}\|_{\ell_2}^2 \boldsymbol{I} + (\|\boldsymbol{\omega}\|_{\ell_2}^2 + \boldsymbol{\omega}^{\mathsf{T}} \boldsymbol{R} \boldsymbol{\omega}) \boldsymbol{\Sigma} \right)^{-1} \boldsymbol{\Sigma} \boldsymbol{\omega}^{\mathsf{T}} \boldsymbol{r}$$

$$= \frac{\boldsymbol{\omega}^{\mathsf{T}} \boldsymbol{r}}{M \|\boldsymbol{\omega}\|_{\ell_2}^2} \left(\frac{\|\boldsymbol{\omega}\|_{\ell_2}^2 + \boldsymbol{\omega}^{\mathsf{T}} \boldsymbol{R} \boldsymbol{\omega}}{M \|\boldsymbol{\omega}\|_{\ell_2}^2} \boldsymbol{I} + \boldsymbol{\Sigma}^{-1} \right)^{-1}$$

$$= \frac{\boldsymbol{\omega}^{\mathsf{T}} \boldsymbol{r}}{M \|\boldsymbol{\omega}\|_{\ell_2}^2} \left(\frac{\boldsymbol{\gamma} + 1}{M} \cdot \boldsymbol{I} + \boldsymbol{\Sigma}^{-1} \right)^{-1},$$
(43a)
$$= \frac{\boldsymbol{\omega}^{\mathsf{T}} \boldsymbol{r}}{M \|\boldsymbol{\omega}\|_{\ell_2}^2} \left(\frac{\boldsymbol{\gamma} + 1}{M} \cdot \boldsymbol{I} + \boldsymbol{\Sigma}^{-1} \right)^{-1},$$

1054 where $\gamma = \omega^{\top} \mathbf{R} \omega / \|\omega\|_{\ell_2}^2$.

1055 Further, 1056

$$\omega = \left(\left(M \operatorname{tr} \left(\tilde{\boldsymbol{P}}^{\mathsf{T}} \tilde{\boldsymbol{P}} \right) + \operatorname{tr} \left(\Sigma \tilde{\boldsymbol{P}}^{\mathsf{T}} \tilde{\boldsymbol{P}} \right) \right) \boldsymbol{I} + \operatorname{tr} \left(\Sigma \tilde{\boldsymbol{P}}^{\mathsf{T}} \tilde{\boldsymbol{P}} \right) \boldsymbol{R} \right)^{-1} \operatorname{tr} \left(\Sigma \tilde{\boldsymbol{P}} \right) \boldsymbol{r}$$

$$= \frac{\operatorname{tr} \left(\Sigma \tilde{\boldsymbol{P}} \right)}{M \operatorname{tr} \left(\tilde{\boldsymbol{P}}^{\mathsf{T}} \tilde{\boldsymbol{P}} \right) + \operatorname{tr} \left(\Sigma \tilde{\boldsymbol{P}}^{\mathsf{T}} \tilde{\boldsymbol{P}} \right)} \left(\boldsymbol{I} + \frac{\operatorname{tr} \left(\Sigma \tilde{\boldsymbol{P}}^{\mathsf{T}} \tilde{\boldsymbol{P}} \right)}{M \operatorname{tr} \left(\tilde{\boldsymbol{P}}^{\mathsf{T}} \tilde{\boldsymbol{P}} \right) + \operatorname{tr} \left(\Sigma \tilde{\boldsymbol{P}}^{\mathsf{T}} \tilde{\boldsymbol{P}} \right)} \boldsymbol{R} \right)^{-1} \boldsymbol{r}.$$
(43b)

1062 Let

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

1066 Then, we get

$$\begin{aligned} \frac{\operatorname{tr}\left(\boldsymbol{\Sigma}\tilde{\boldsymbol{P}}^{\mathsf{T}}\tilde{\boldsymbol{P}}\right)}{M\operatorname{tr}\left(\tilde{\boldsymbol{P}}^{\mathsf{T}}\tilde{\boldsymbol{P}}\right) + \operatorname{tr}\left(\boldsymbol{\Sigma}\tilde{\boldsymbol{P}}^{\mathsf{T}}\tilde{\boldsymbol{P}}\right)} &= \left(1 + M\frac{\operatorname{tr}\left(\tilde{\boldsymbol{P}}^{\mathsf{T}}\tilde{\boldsymbol{P}}\right)}{\operatorname{tr}\left(\boldsymbol{\Sigma}\tilde{\boldsymbol{P}}^{\mathsf{T}}\tilde{\boldsymbol{P}}\right)}\right)^{-1} \\ &= \left(1 + M\frac{\operatorname{tr}\left(\boldsymbol{\Sigma}_{\gamma}^{-2}\right)}{\operatorname{tr}\left(\boldsymbol{\Sigma}\boldsymbol{\Sigma}_{\gamma}^{-2}\right)}\right)^{-1} \\ &= \left(1 + M\sum_{i=1}^{d}\frac{s_{i}^{2}}{(M + (\gamma + 1)s_{i})^{2}}\left(\sum_{i=1}^{d}\frac{s_{i}^{3}}{(M + (\gamma + 1)s_{i})^{2}}\right)^{-1}\right)^{-1} \\ &=: g(\gamma). \end{aligned}$$

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

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

1084

1087 1088 1089

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

(46)

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

$$\gamma = \frac{\mathbf{r}^{\top}(g(\gamma)\mathbf{R} + \mathbf{I})^{-1}\mathbf{R}(g(\gamma)\mathbf{R} + \mathbf{I})^{-1}\mathbf{r}}{\mathbf{r}^{\top}(g(\gamma)\mathbf{R} + \mathbf{I})^{-2}\mathbf{r}}$$
$$= \sum_{i=1}^{n} \frac{\lambda_{i}a_{i}^{2}}{(1 + g(\gamma)\lambda_{i})^{2}} \left(\sum_{i=1}^{n} \frac{a_{i}^{2}}{(1 + g(\gamma)\lambda_{i})^{2}}\right)^{-1}$$
(45)

1093

1094 where the second equality follows from Assumption A and the fact that $\mathbf{R} = \mathbf{E} \operatorname{diag}(\lambda) \mathbf{E}^{\mathsf{T}}$ denotes the eigen decomposition of \mathbf{R} , with $\lambda = [\lambda_1, \dots, \lambda_n]^{\mathsf{T}} \in \mathbb{R}^n_+$.

 $\tilde{\boldsymbol{P}} = C(\boldsymbol{r}, \boldsymbol{\omega}, \boldsymbol{\Sigma}) \cdot \left(\frac{\boldsymbol{\gamma}^{\star} + 1}{M} \cdot \boldsymbol{I} + \boldsymbol{\Sigma}^{-1}\right)^{-1}, \text{ and }$

 $\boldsymbol{P}(\boldsymbol{\gamma}) = C(\boldsymbol{r}, \boldsymbol{\omega}, \boldsymbol{\Sigma}) \cdot \boldsymbol{\Sigma}^{-\frac{1}{2}} \left(\frac{\boldsymbol{\gamma}^{\star} + 1}{\sigma^{2} + \operatorname{tr}(\boldsymbol{\Sigma})} \cdot \boldsymbol{\Sigma} + \boldsymbol{I} \right)^{-1} \boldsymbol{\Sigma}^{-\frac{1}{2}}, \text{ and}$

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

 $\boldsymbol{\omega} = c(\boldsymbol{r}, \boldsymbol{\omega}, \boldsymbol{\Sigma}) \cdot \left(g(\boldsymbol{\gamma}^{\star}) \cdot \boldsymbol{R} + \boldsymbol{I} \right)^{-1} \boldsymbol{r}.$

 $=: h(g(\gamma)),$

1099 1100

1102 1103 for some $C(\mathbf{r}, \omega, \Sigma) = \frac{\omega^{\top} \mathbf{r}}{M \|\omega\|_{\ell_2}^2}$ and $c(\mathbf{r}, \omega, \Sigma) = \frac{\operatorname{tr}(\Sigma \tilde{P})}{M \operatorname{tr}(\tilde{P}^{\top} \tilde{P}) + \operatorname{tr}(\Sigma \tilde{P}^{\top} \tilde{P})}$. 1104

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

1106

1109

1110 1111 This completes the proof.

- 1112

 1113
 C.2
 Proof of Theorem 4
- 1114 We first provide the following Lemma.

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

 $\omega(\gamma) = c(\mathbf{r}, \omega, \Sigma) \cdot \left(g(\gamma^{\star}) \cdot \mathbf{R} + \mathbf{I}\right)^{-1} \mathbf{r}.$

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

1119 1120 1121

1122

1118

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

1123
1124 where
$$M = \sigma^2 + \sum_{i=1}^d s_i$$
.

1125 Suppose $\Delta_{\Sigma} \cdot \Delta_{R} < M + s_{\min}$, where Δ_{Σ} and Δ_{R} denote the effective spectral gaps of Σ and R, 1126 respectively, as given in (12); and s_{\min} is the smallest eigenvalue of Σ . We have that

1127
1128
1129
$$\left|\frac{\partial g}{\partial \gamma} \cdot \frac{\partial h}{\partial g}\right| \le \frac{\Delta_{\Sigma}^2 \cdot \Delta_R^2}{(M + s_{\min})^2} < 1$$

1130 1131 *Proof.* Let

1132
1133
$$B(\gamma) = \sum_{i=1}^{d} \frac{s_i^3}{(M + (\gamma + 1)s_i)^2}, \quad C(\gamma) = \sum_{i=1}^{d} \frac{s_i^2}{(M + (\gamma + 1)s_i)^2}, \quad A(\gamma) = 1 + M \frac{C(\gamma)}{B(\gamma)}.$$

The derivatives of
$$B(y)$$
 and $C(y)$ are

$$B'(y) = -2 \sum_{i=1}^{d} \frac{s_{i}^{4}}{(M + (y + 1)s_{i})^{3}}, \quad C'(y) = -2 \sum_{i=1}^{d} \frac{s_{i}^{3}}{(M + (y + 1)s_{i})^{3}},$$
The gradient of $g(y)$ is

$$\frac{\partial g}{\partial y} = -M \left(\frac{1}{A(y)B(y)}\right)^{2} (C'(y)B(y) - C(y)B'(y)).$$
(49)
If can be seen that

$$\left(\frac{1}{A(y)}\right)^{2} \leq M^{-2} \left(\sum_{i=1}^{d} \frac{s_{i}^{3}}{(M + (y + 1)s_{i})^{2}}\right)^{2} \left(\sum_{i=1}^{d} \frac{s_{i}^{2}}{(M + (y + 1)s_{i})^{2}}\right)^{-2}.$$
(50a)
Further, we have

$$C'(y)B(y) - C(y)B'(y) = -2 \sum_{i=1}^{d} \frac{s_{i}^{3}}{(M + (y + 1)s_{i})^{2}} \sum_{i=1}^{d} \frac{s_{i}^{3}}{(M + (y + 1)s_{i})^{2}}\right)^{-2}.$$
(50b)

$$= \sum_{i=1}^{d} \sum_{j=1}^{d} \frac{s_{j}^{2}}{(M + (y + 1)s_{i})^{2}} \sum_{i=1}^{d} \frac{s_{i}^{4}}{(M + (y + 1)s_{i})^{3}},$$
where

$$T_{ij} = s_{i}^{2}(M + (y + 1)s_{i})s_{i}^{4} + s_{i}^{4}s_{i}^{2}(M + (y + 1)s_{i})^{3},$$
where

$$T_{ij} = s_{i}^{2}(M + (y + 1)s_{i})s_{j}^{4} + s_{i}^{4}s_{i}^{2}(M + (y + 1)s_{i})s_{j}^{3}},$$
(50c)

$$= s_{i}^{2}s_{i}^{2}(M + (y + 1)s_{i})s_{i}^{4} + s_{i}^{4}s_{i}^{2}(M + (y + 1)s_{i})s_{i}^{3},$$
Thus, substituting (50a) and (50b) into (49), we obtain

$$\left|\frac{\partial g}{\partial y}\right| \leq M \cdot M^{-1}\left(\sum_{i=1}^{d} \frac{s_{i}^{2}}{(M + (y + 1)s_{i})^{2}}\right)^{-2}\sum_{i,j=1}^{d} \frac{s_{i}^{2}s_{j}^{2}(s_{i} - s_{j})^{2}}{(M + (y + 1)s_{i})^{3}(M + (y + 1)s_{i})^{3}}.$$
(51)
Next, we derive $\frac{\partial g}{\partial y}$. Let

$$D(\hat{y}) = \sum_{i=1}^{n} \frac{\lambda_{i}a_{i}^{2}}{(1 + \gamma_{i}\lambda_{i})^{2}}, \quad E(\hat{y}) = \sum_{i=1}^{n} \frac{\lambda_{i}a_{i}^{2}}{(1 + \gamma_{i}\lambda_{i})^{3}}.$$
The derivative of h with respect to \hat{y} is yet by

The derivative of *h* with respect to $\bar{\gamma}$ is given by $\frac{\partial h}{\partial \bar{\gamma}} = -\left(\frac{1}{E(\bar{\gamma})}\right)^2 \left(E(\bar{\gamma})D'(\bar{\gamma}) - D(\bar{\gamma})E'(\bar{\gamma})\right).$

(52)

¹¹⁸⁸ Substituting into (52), we get

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

1193 and

 $E(\bar{\gamma})D'(\bar{\gamma}) - D(\bar{\gamma})E'(\bar{\gamma}) = 2\sum_{i=1}^{n} \frac{\lambda_{i}^{2}a_{i}^{2}}{(1 + \bar{\gamma}\lambda_{i})^{3}} \sum_{i=1}^{n} \frac{a_{i}^{2}}{(1 + \bar{\gamma}\lambda_{i})^{2}}$ $- 2\sum_{i=1}^{n} \frac{\lambda_{i}a_{i}^{2}}{(1 + \bar{\gamma}\lambda_{i})^{2}} \sum_{i=1}^{n} \frac{a_{i}^{2}\lambda_{i}}{(1 + \bar{\gamma}\lambda_{i})^{3}}$ $= \sum_{i=1}^{n} \sum_{j=1}^{n} \frac{\bar{T}_{ij}}{(1 + \bar{\gamma}\lambda_{i})^{3}(1 + \bar{\gamma}\lambda_{j})^{3}}$ $= \sum_{i=1}^{n} \sum_{j=1}^{n} \frac{a_{i}^{2}a_{j}^{2}\left(\lambda_{i}^{2} + \lambda_{j}^{2} - 2\lambda_{i}\lambda_{j}\right)}{(1 + \bar{\gamma}\lambda_{i})^{3}(1 + \bar{\gamma}\lambda_{j})^{3}}.$ (53b)

1206 Here,

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

1212 Hence, substituting (53a) and (53b) into (52) gives

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

¹²¹⁷ Now, for the combined derivative, we have

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

1225 Note that $M + (\gamma + 1)s_i$ and $1 + \bar{\gamma}\lambda_j$ are nonnegative for all *i*, *j*. Hence,

$$\left|\frac{\partial g}{\partial \gamma} \cdot \frac{\partial h}{\partial \bar{\gamma}}\right| \leq \left(\sum_{i=1}^{d} \frac{s_i^2 \left(M + (\gamma + 1)s_i\right)}{\left(M + (\gamma + 1)s_i\right)^3}\right)^{-2}$$

1230
1231
1232
$$\cdot \sum_{i,j=1}^{d} \frac{s_i^2 s_j^2 \cdot \Delta_1 \cdot \left(M + (\gamma + 1) s_j\right) (M + (\gamma + 1) s_i)}{(M + (\gamma + 1) s_i)^3 \left(M + (\gamma + 1) s_j\right)^3}$$

1232
1233
1234

$$i,j=1$$
 $(M + (\gamma + 1))$
 $\left(\sum_{i}^{n} \frac{a_{i}^{2}(1 + \bar{\gamma}\lambda_{i})}{(1 - \bar{\gamma}\lambda_{i}^{2})}\right)^{-2}$

1235
1236
$$\left(\sum_{i=1}^{2} (1 + \bar{\gamma}\lambda_i)^3\right)$$

$$\Rightarrow a_i^2 a_i^2 \cdot \Delta_2 \cdot (1 + \bar{\gamma}\lambda_i) (1 + \bar{\gamma}\lambda_i)$$

1237
1238
1239

$$\cdot \sum_{i,j \in S} \frac{\alpha_i \alpha_j - \alpha_2 - (1 + \bar{\gamma}\lambda_j) (1 + \bar{\gamma}\lambda_j)}{(1 + \bar{\gamma}\lambda_j)^3 (1 + \bar{\gamma}\lambda_j)^3}$$

1240 where

$$\Delta_{1} := \max_{i,j} \frac{(s_{i} - s_{j})^{2}}{\left(M + (\gamma + 1)s_{j}\right)(M + (\gamma + 1)s_{i})}, \quad \Delta_{2} := \max_{i,j \in \mathcal{S}} \frac{(\lambda_{i} - \lambda_{j})^{2}}{(1 + \bar{\gamma}\lambda_{i})(1 + \bar{\gamma}\lambda_{j})}.$$
(55)

 1242 Here, $S = \{i \in [n] | \lambda_i \neq 0\} \subset [n].$

1245

1246

1260 1261

1266

1269

1279

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

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

1247 where Δ_{Σ} and Δ_{R} are the spectral gaps of Σ and R; and s_{\min} is the smallest eigenvalue of Σ ; and 1248 $M = \sigma^{2} + \sum_{i=1}^{d} s_{i}$.

1250 *Proof of Theorem 4.* Lemma 1 shows that $|\partial h(g(\gamma))/\partial \gamma| < 1$, and as a result, the mapping $h(g(\gamma))$ 1251 on \mathbb{R}_+ is a contraction mapping. Therefore, by the Banach fixed-point theorem, this guarantees the 1252 existence of a unique root, denoted as $\gamma = \gamma^*$. This completes the proof of **T1**. In the following, we 1253 provide the proof of **T2**. Substituting the unique γ^* into (16) and using the fact that $\tilde{P} = \sqrt{\Sigma}P\sqrt{\Sigma}$, 1254 we obtain (P^*, ω^*) , as given in (14), as a global minima of (3).

1255 Next, we claim that $(\mathbf{P}^*, \boldsymbol{\omega}^*)$ is the unique global minimizer of $\mathcal{L}(\mathbf{P}, \boldsymbol{\omega})$ up to rescaling, i.e., any 1256 other minimizer $(\hat{\mathbf{P}}, \hat{\boldsymbol{\omega}})$ must be related to $(\mathbf{P}^*, \boldsymbol{\omega}^*)$ by scaling factors α and β , such that $\hat{\mathbf{P}} = \alpha \mathbf{P}^*$ 1257 and $\hat{\boldsymbol{\omega}} = \beta \boldsymbol{\omega}^*$, for some $\alpha, \beta > 0$.

¹²⁵⁸ The loss function is given by

$$\mathcal{L}(\boldsymbol{P},\boldsymbol{\omega}) = \boldsymbol{M} - 2\operatorname{tr}\left(\boldsymbol{\Sigma}^{2}\boldsymbol{P}\right)\boldsymbol{\omega}^{\mathsf{T}}\boldsymbol{r} + \boldsymbol{M}\|\boldsymbol{\omega}\|^{2}\operatorname{tr}\left(\boldsymbol{\Sigma}\boldsymbol{P}^{\mathsf{T}}\boldsymbol{\Sigma}\boldsymbol{P}\right) + (\|\boldsymbol{\omega}\|^{2} + \boldsymbol{\omega}^{\mathsf{T}}\boldsymbol{R}\boldsymbol{\omega})\operatorname{tr}\left(\boldsymbol{\Sigma}^{2}\boldsymbol{P}^{\mathsf{T}}\boldsymbol{\Sigma}\boldsymbol{P}\right)$$

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

1264
$$\mathcal{L}(\alpha \boldsymbol{P}, \beta \boldsymbol{\omega}) = M - 2\alpha\beta \operatorname{tr}(\boldsymbol{\Sigma}^2 \boldsymbol{P})\boldsymbol{\omega}^\top \boldsymbol{r} + M\alpha^2 \beta^2 ||\boldsymbol{\omega}||^2 \operatorname{tr}(\boldsymbol{\Sigma} \boldsymbol{P}^\top \boldsymbol{\Sigma} \boldsymbol{P}) + \alpha^2 \beta^2 (||\boldsymbol{\omega}||^2 + \boldsymbol{\omega}^\top \boldsymbol{R} \boldsymbol{\omega}) \operatorname{tr}(\boldsymbol{\Sigma}^2 \boldsymbol{P}^\top \boldsymbol{\Sigma} \boldsymbol{P}).$$
1265 Define

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

1267 Thus, the rescaled loss function becomes 1268

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

¹²⁷⁰ For (P^*, ω^*) to be a minimizer, the partial derivatives of the loss function with respect to P and ω ¹²⁷¹ must vanish at (P^*, ω^*) . However, we consider the effect of the rescaling in terms of α and β . To ¹²⁷² find the stationary points of $\mathcal{L}(\alpha P, \beta \omega)$, we differentiate with respect to α and β :

1273
1274
1275
1276

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

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

¹²⁷⁷ Setting these to zero, we obtain the system

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

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

$$\hat{P} = \alpha P^{\star}$$
 and $\hat{\omega} = \beta \omega^{2}$

1284 for some scalars $\alpha, \beta > 0$ satisfying (56).

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

1288 1289 C.3 Proof of Corollary 2

1290 *Proof.* Since by assumption $\Sigma = I$, it follows from (13b) that

1

1292
1293
1293
1294

$$g(\gamma^{\star}) = \left(1 + (d + \sigma^2) \sum_{i=1}^{d} \frac{1}{(d + \sigma^2 + \gamma^{\star} + 1)^2} \left(\sum_{i=1}^{d} \frac{1}{(d + \sigma^2 + \gamma^{\star} + 1)^2}\right)\right)$$

d

$$5 = \frac{1}{d}$$

1 1

Substituting this into (14) gives

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

Now, recall that the objective function is given by

$$\mathcal{L}(\omega) = M - 2\operatorname{tr}(\Sigma^2 P)\omega^{\mathsf{T}} r + M \|\omega\|_{\ell_2} \operatorname{tr}(\Sigma P^{\mathsf{T}} \Sigma P) + (\|\omega\|^2 + \omega^{\mathsf{T}} R\omega)\operatorname{tr}(\Sigma^2 P^{\mathsf{T}} \Sigma P)$$

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

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

$$\mathcal{L}(\boldsymbol{\omega}^{\star}) = (\sigma^{2} + d) - 2 \cdot d \cdot \boldsymbol{r}^{\mathsf{T}} \boldsymbol{\omega}^{\star} + (\sigma^{2} + d) \cdot \left\|\boldsymbol{\omega}^{\star}\right\|_{\ell_{2}}^{2} d + d\left(\|\boldsymbol{\omega}^{\star}\|^{2} + \boldsymbol{\omega}^{\mathsf{T}} \boldsymbol{R} \boldsymbol{\omega}^{\star}\right).$$

The expression simplifies as

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

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

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

and

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

Thus,

$$\begin{aligned} & |1319 \\ & |1320 \\ & |1320 \\ & |1321 \\ & |1322 \end{aligned} \qquad (d + \sigma^2 + 1) ||\omega^{\star}||^2 + \omega^{\star \top} R \omega^{\star} = r^{\top} \left(R + (d + \sigma^2 + 1) I \right)^{-1} \left((d + \sigma^2 + 1) I + R \right) \left(R + (d + \sigma^2 + 1) I \right)^{-1} r \\ & = r^{\top} \left(R + (d + \sigma^2 + 1) I \right)^{-1} r. \end{aligned}$$

Substituting this result back into the objective function gives

$$\mathcal{L}(\boldsymbol{\omega}^{\star}) = (\sigma^2 + d) - d \cdot \boldsymbol{r}^{\mathsf{T}} \left(\boldsymbol{R} + (d + \sigma^2 + 1) \boldsymbol{I} \right)^{-1} \boldsymbol{r}.$$

D LOSS LANDSCAPE OF 1-LAYER GLA

D.1 PROOF OF THEOREM 5

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

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

Define a set $\bar{\mathcal{W}} := \left\{ [\omega_1 \cdots \omega_n]^\top \in \mathbb{R}^n \mid 0 \le \omega_i \le \omega_j \le 1, \forall 1 \le i \le j \le n \right\}$ and we have $\mathcal{W} \in \bar{\mathcal{W}}$. Given scalar gating $G_i = \begin{bmatrix} * & * \\ g_i \mathbf{1}^T & * \end{bmatrix}$, following (10), the weighting vector returns

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

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 \overline{W} .

Next, we will show that

$$\omega^{\star} \in \mathcal{W}$$
 where $\omega^{\star} = \arg \min_{P, \omega \in \tilde{\mathcal{W}}} \mathcal{L}_{WPGD}(P, \omega).$

Define the weighting vector $\boldsymbol{\omega} = [\boldsymbol{\omega}_1^\top \cdots \boldsymbol{\omega}_K^\top]^\top \in \mathbb{R}^n$ where we have $\boldsymbol{\omega}_k = [\boldsymbol{\omega}_1^{(k)} \cdots \boldsymbol{\omega}_{n_k}^{(k)}]^\top \in \mathbb{R}^{n_k}$. For any $\omega \notin W$, there exist (i, j, k), i = j - 1 such that $\omega_i^{(k)} < \omega_i^{(k)}$. Given gradient in (42), we have that 1350 1351 $\nabla_{\omega_{i}^{(k)}} \mathcal{L} = c_{1} \cdot \omega_{i}^{(k)} + c_{2} \text{ and } \nabla_{\omega_{j}^{(k)}} \mathcal{L} = c_{1} \cdot \omega_{j}^{(k)} + c_{2} \text{ with for some } c_{1}, c_{2} \text{ with } c_{1} > 0. \quad \nabla_{\omega_{i}^{(k)}} \mathcal{L} < \nabla_{\omega_{j}^{(k)}} \mathcal{L}$ 1352
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. 1353
This results in showing that $\omega^{\star} \in \mathcal{W}$.

Finally, we will show that any $\omega \in W$ 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 W and assume that $\omega_K = \alpha < 1$ without loss of generality. Then such sample weighting can be achieved via the gating

$$\begin{bmatrix} \mathbf{1}_{n_1}^{\top} & \frac{\omega_1}{\omega_{1:K}} & \cdots & \mathbf{1}_{n_k}^{\top} & \frac{\omega_k}{\omega_{k:K}} & \cdots & \mathbf{1}_{n_K}^{\top} & \frac{\omega_K}{\omega_{K:K}} \end{bmatrix}^{\top}.$$

Let $\omega'_k := \frac{\omega_k}{\omega_{k\cdot \kappa}}$ and let w_g be in the form of

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

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

$$\phi(\tilde{\boldsymbol{w}}_{g}^{\top}\bar{\boldsymbol{d}}_{k}) \begin{cases} = 1, & k = 0 \\ = \omega_{k}', & k \in [K] \end{cases}$$

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

$$\boldsymbol{\omega}^{\star} = \left(g(\boldsymbol{\gamma}^{\star}) \cdot \boldsymbol{R} + \boldsymbol{I}\right)^{-1}$$

Since Assumption C holds and $n_1 = n_2 = \cdots = n_K := \bar{n}, \omega^*$ takes the form of $\omega^* = cr$ 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.

1377 1378

1384 1385

1358 1359

1361 1362

1363

1367

1372 1373

1379 D.2 Proof of Theorem 6 1380

1381 1382 *Proof.* Following the similar proof of Theorem 5, and letting $\tilde{W} := \left\{ \left[\omega_1 \mathbf{1}_{n_1}^\top \cdots \omega_K \mathbf{1}_{n_K}^\top \right]^\top \in \mathbb{R}^n \right\}$, we obtain

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

1386 Therefore, it remains to show that any $\omega \in \tilde{W}$ can be implemented via some gating function. Let 1387 $\omega = [\omega_1 \mathbf{1}_{n_1}^\top \cdots \omega_K \mathbf{1}_{n_K}^\top]$ be arbitrary weighting in \tilde{W} . Theorem 5 has shown that if $\omega_1 \le \omega_2 \le \cdots \le \omega_K$, 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:

- $\boldsymbol{\omega}_1 = \boldsymbol{\omega}_1 [\mathbf{1}_{n_1}^\top \cdots \mathbf{1}_{n_K}^\top]$
- 1393 $\boldsymbol{\omega}_2 = (\boldsymbol{\omega}_2 \boldsymbol{\omega}_1)[\mathbf{0}_{n_1}^{\mathsf{T}} \ \mathbf{1}_{n_2}^{\mathsf{T}} \ \cdots \ \mathbf{1}_{n_K}^{\mathsf{T}}]$ 1394
- 1394 1395 $\boldsymbol{\omega}_3 = (\boldsymbol{\omega}_3 - \boldsymbol{\omega}_2) [\mathbf{0}_{n_1}^{\mathsf{T}} \ \mathbf{0}_{n_2}^{\mathsf{T}} \ \mathbf{1}_{n_3}^{\mathsf{T}} \ \cdots \ \mathbf{1}_{n_K}^{\mathsf{T}}]$
 - $\boldsymbol{\omega}_{K} = (\omega_{3} \omega_{2}) [\mathbf{0}_{n_{1}}^{\mathsf{T}} \ \mathbf{0}_{n_{2}}^{\mathsf{T}} \ \mathbf{0}_{n_{3}}^{\mathsf{T}} \ \cdots \ \mathbf{0}_{n_{K-1}}^{\mathsf{T}} \mathbf{1}_{n_{K}}^{\mathsf{T}}]$

1308

1396

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

1401 From (27) and (28), let *i*'th dimension implements the weighting ω_i for $i \in [K]$. Specifically, let 1402 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 1403 $u_i = \omega_i - \omega_j$ with $\omega_0 = 0$. Then the composed weighting following (28) returns ω , which completes the proof. 1404 D.3 Proof of Corollary 3

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

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

$$\mathcal{L}(\boldsymbol{P}^{\star},\boldsymbol{\omega}=1) = d + \sigma^2 - 2cd\mathbf{1}^{\mathsf{T}}\boldsymbol{r} + (d + \sigma^2)c^2nd + (n + \mathbf{1}^{\mathsf{T}}\boldsymbol{R}\mathbf{1})c^2d$$
$$= d + \sigma^2 - cd\mathbf{1}^{\mathsf{T}}\boldsymbol{r}.$$

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

16 It completes the proof.