ACCELERATING BLOCK COORDINATE DESCENT FOR LLM FINETUNING VIA LANDSCAPE CORRECTION

Anonymous authors

Paper under double-blind review

ABSTRACT

Training and finetuning large language models (LLMs) are resource-intensive tasks, with memory limitations being a key bottleneck. A classic optimization method, block coordinate descent (BCD), offers solutions by segmenting the trainable parameters into multiple blocks and optimizing one active block at a time while freezing the others, thereby significantly reducing memory cost. However, we identify that blindly applying BCD to train LLMs can be inefficient for two reasons. First, optimizing only the active block requires backpropagating through multiple deeper yet inactive blocks, resulting in wasteful computations. Second, the frozen blocks, when they are not quite close to optimality, can narrow the optimization landscape, potentially misguiding the training of the active block. To address these issues simultaneously, we propose integrating BCD with *landscape correction*, which unfreezes the inactive blocks and updates them in a cost-efficient manner during the same backpropagation as the update to the active block. We show that our method empirically improves vanilla BCD with minimal additional computation and memory. Experiments on 8B and 70B models demonstrate that our proposed method surpasses memory efficient baselines and matches Adam's downstream performance while reducing memory cost by 80% compared to Adam.

027 028 029

030

004

005

010 011

012

013

014

015

016

017

018

019

021

023

025

026

1 INTRODUCTION

Large language models (LLMs) have gained significant popularity within the research community and industry. Training these models typically entails a pretraining phase on a vast dataset, followed by a series of finetuning adjustments to tailor the model for specific domain tasks. Both phases demand extensive computational resources, with memory being a primary constraint. For instance, optimizing a model containing N billion parameters using standard mixed-precision training requires at least 18N gigabytes of GPU memory (refer to Section 2 for additional details). Given the usual limitations of GPU memory, this constraint impedes researchers from experimenting with larger models.

To address this practical challenge, researchers have developed memory efficient algorithms for LLM training such as parameter efficient finetuning (PEFT), including Adapter Houlsby et al. (2019), LoRA Hu et al. (2021), prompt tuning Lester et al. (2021), prefix tuning Li & Liang (2021), etc. These techniques focus on training a small set of additional parameters while maintaining the original pretrained model unchanged. Other memory efficient methods for full parameter training have also been investigated. For example, Galore Zhao et al. (2024) applies a low-rank space projection to both the gradient and the optimizer's states to reduce memory consumption.

In addition to the existing approaches, a classic optimization paradigm, known as *block coordinate descent* (BCD), holds a strong potential for memory efficient LLM training and finetuning. Intuitively, BCD reduces the memory cost by partitioning the trainable parameters into several blocks and optimizing over only one active block at a time. For instance, segmenting an *N*-billion parameter model into *D* blocks decreases the memory consumption from 18N to $2N + \frac{16N}{D}$ GB, as only the gradient and optimizer states of the active block need to be stored. In fact, BCD has been the method of choice in many data science problems, with a wide array of variants developed for improving memory, performance, convergence, and efficiency; see, e.g., (Hsieh et al., 2008; Chang & Roth, 2011; Yu et al., 2012; Treister & Turek, 2014; Richtárik & Takáč, 2014).

061

062

063

064

065

066

067

078

079

081

082

083

085

090

091 092

093 094

095

096

In stark contrast to the previous memory efficient methods, and despite its intuitive memory benefits,
BCD has been overlooked and rarely explored in the context of LLMs. It was not until very recently
that Luo et al. (2024) proposed BAdam, which integrates BCD into an LLM finetuning framework
by training each active block with several Adam steps. Even though BAdam has shown preliminary
success in reducing memory cost during training and improving performance at test time, its direct
use of vanilla BCD leaves at least two fundamental aspects to be questioned:

- Computing the (stochastic) gradient for a *single* active block via backpropagation necessitates calculating the partial derivatives of the activations of *multiple* deeper yet inactive layers. This is wasteful of computation as these partial derivatives are not even used to update their corresponding weights.
- Given that the training objective is highly nonconvex, and since all blocks are frozen except for the active one, BCD tends to be misled by *its local view of the optimization landscape*, which potentially slows down its convergence speed.

Standing on the ground offered by BAdam, we push the research frontier by proposing a simple 068 remedy to the above two issues. Our method, termed BREAD, is a blend of two components: (1) 069 Similarly to BAdam, we update the active block using several Adam steps (the BCD component); (2) differently, we unfreeze the inactive blocks and update them using lightweight memory efficient 071 optimization techniques (the landscape correction component). Since landscape correction utilizes 072 the gradients of the activations that are already calculated, it adds minimal additional computation 073 and in fact addresses the wasteful computation issue. Furthermore, the landscape correction com-074 ponent provides BCD a better view of the optimization landscape for better updating the current 075 active block, thereby addressing the second point of concern. In combination, BREAD maintains 076 the memory efficient feature of BCD with improved learning capability and faster convergence. Our 077 main contributions are outlined as follows:

- Limitations of Standard BCD in LLMs: Our research identifies two fundamental limits of vanilla BCD when applied to neural networks (hence LLMs): The wasted computation of gradients during backpropagation and the suboptimal landscape caused by freezing inactive blocks. These limitations partly explain why the application of BCD in neural networks is uncommon.
- Blending BCD with Landscape Correction: We propose a new method termed BREAD, which combines BCD with a landscape correction technique to address these two limitations simultaneously. It unfreezes some of (or all) the inactive blocks and updates them using memory efficient optimization techniques. BREAD maintains the memory efficiency of BCD with improved optimization ability.
 - Excellent Performance: Our experiments on instruction tuning and preference optimization with the Llama 3.1-8B and Llama 3.1-70B models demonstrate that BREAD clearly outperforms state-of-the-art memory efficient training methods and achieves comparable downstream performance to that of Adam on five math benchmarks and MT-bench scores.

2 PRELIMINARIES ON BLOCK COORDINATE DESCENT FOR LLM TRAINING

Our main focus lies in improving the efficiency of BCD when utilized to finetune LLMs. Therefore, we first review some preliminary concepts of LLM and BCD in this section.

Objective of training LLMs. Consider minimizing a general objective function $\min_{W} H(W) = \frac{1}{n} \sum_{j=1}^{n} h_j(W)$, where $W \in \mathbb{R}^d$, *n* is the number of data samples. In the context of training/finetuning LLMs, $h_j(W)$ represents the negative log-likelihood of the autoregressive probability $\mathbb{P}_{W}[y_j|x_j] = \prod_{s=1}^{m} \mathbb{P}_{W}[y_{j,s}|y_{j,1:s-1}, x_j]$ for the *j*-th prompt x_j and its corresponding *j*-th output y_j . In most LLM models, this autoregressive probability is modeled by a transformer architecture Vaswani et al. (2017), and thus $W \in \mathbb{R}^d$ encompasses all trainable parameters of the transformer, including the query, key, value, and output attention matrices, as well as the gate, up, and down projection matrices of each transformer layer.

BCD for LLM training. Instead of minimizing the objective function over the entire set of trainable parameters W, the key idea of BCD is to break down this high dimensional optimization problem into a series of lower dimensional ones, thereby significantly reducing the memory requirement of GPU RAM. Specifically, BCD first splits the model parameters into D block, i.e.,

$$\boldsymbol{W}_{\ell}^{t+1} \in \operatorname*{argmin}_{\boldsymbol{W}_{\ell} \in \mathbb{R}^{d_{\ell}}} \frac{1}{|\mathcal{N}|} \sum_{j \in \mathcal{N}} h_{j}(\boldsymbol{W}_{1}^{t+1}, \cdots, \boldsymbol{W}_{\ell-1}^{t+1}, \boldsymbol{W}_{\ell}, \boldsymbol{W}_{\ell+1}^{t}, \cdots, \boldsymbol{W}_{D}^{t})$$
(1)

where $\mathcal{N} \subseteq \{1, \dots, n\}$ is a batch of the training dataset. Updating from block $\ell = 1$ to block $\ell = D$ is counted as one block-epoch. Since it is intractable to solve (1) exactly, one can instead approximate the solution by implementing K Adam steps, as utilized in BAdam Luo et al. (2024).

120 BCD is a memory efficient full parameter optimization method for LLM training. It is evident 121 to see that BCD is a full parameter optimization method, as all the trainable parameters W will be 122 updated after one block-epoch. More importantly, BCD is also memory efficient.

123 Let us first analyze the memory consumption of the Adam optimizer under the mixed precision 124 training setting Micikevicius et al. (2017). The memory cost is attributed to the storage of the 125 model parameters, gradients, and optimizer states. We consider an LLM with N billion parameters 126 and express GPU memory consumption in gigabytes (GB). Initially, one must store the FP16 model 127 parameters for the backpropagation (BP) process, requiring 2N memory. Additionally, the optimizer maintains a copy of the model in FP32 precision, consuming another 4N memory. The gradients, 128 momentum, and second moment vectors are all stored in FP32 precision with each requiring 4N129 memory. Consequently, the total memory required is at least 18N. For example, in order to train a 130 Llama 3-8B or a Llama 3-70B model, Adam requires at least 144 GB or 1260 GB of GPU RAM, 131 respectively, which can be prohibitive in limited memory scenarios. 132

In sharp contrast to Adam, BCD only requires storing the FP32 model parameters, gradients, and optimizer states for the *active block* W_{ℓ} , which is only 1/D of the memory consumption needed for all the parameters. Thus, in addition to maintaining an FP16 model that requires 2N memory, BCD needs a total of only $2N + \frac{16N}{D}$ memory. Therefore, for training a Llama 3-8B or a Llama 3-70B model and when D = 32 or D = 80 (partition each transformer layer as a block), BCD only needs roughly 20 GB or 154 GB of GPU RAM, respectively, which is significantly cheaper compared to the costs of Adam. For a more detailed analysis on memory cost, we refer to Luo et al. (2024).

140 141

142

148 149 150

3 LIMITATIONS OF BCD FOR NEURAL NETWORKS

Although BCD is proven to be memory efficient for training and finetuning LLMs, we will illustrate
 two major limitations of the BCD optimization scheme when it is used for training models induced
 by a neural network structure in this section. The results apply to the setting of training and finetun ing LLMs as well. Motivated by these limitations, we will develop a more efficient training method
 based on BCD, which achieves superior performance compared to that of the vanilla BCD.

To ease our analysis, let us consider a L-layer feedforward neural network model:

$$\mathbf{z}_{\ell+1} = f_{\mathbf{W}_{\ell}}^{\ell}(\mathbf{z}_{\ell}), \ \forall 1 \le \ell \le L, \quad \text{with} \quad \mathbf{z}_1 = \mathbf{x},$$
(2)

where L is the total number of layers, x is the input, $f_{W_{\ell}}^{\ell}$ is the ℓ -th layer's transform.

152 Limitation I: Ineffective utilization of intermediate derivatives during backpropagation. Due 153 to the compositional structure of deep neural networks, the gradients of the trainable parameters are 154 calculated according to the chain rule. For example, taking the stochastic gradient of the ℓ -th layer's 155 parameters W_{ℓ} requires computing the partial derivatives with respect to all the activation values of 156 deeper layers, as shown in the following equation:

$$\frac{\partial H}{\partial \boldsymbol{W}_{\ell}} = \underbrace{\frac{\partial H}{\partial \boldsymbol{z}_{L+1}} \frac{\partial \boldsymbol{z}_{L+1}}{\partial \boldsymbol{z}_{L}} \cdots \frac{\partial \boldsymbol{z}_{\ell+2}}{\partial \boldsymbol{z}_{\ell+1}}}_{I_{\ell+1}} \frac{\partial \boldsymbol{z}_{\ell+1}}{\partial \boldsymbol{W}_{\ell}}, \tag{3}$$

159 160

157

158

where *H* is the objective function of the neural network training problem. During the backpropagation process in optimization method like Adam, the intermediate partial derivatives $I_{\ell+1}$ of the 162 activations in (3) are properly utilized for computing the gradients of the $L, L-1, \dots, \ell+1$ -th 163 layers' weight parameters as well. However, since BCD only updates the active block W_{ℓ} , the term 164 $I_{\ell+1}$ is merely used for calculating the gradient of W_{ℓ} , resulting in ineffective utilization of the 165 computed partial derivatives of the activations during backpropagation.

166 Limitation II: Suboptimal landscape of BCD's sub-problem. To tackle a training problem, op-167 timization methods such as Adam optimize over all the trainable parameters W, while the BCD 168 optimization scheme (1) minimizes the objective over only the current active block, keeping the 169 others fixed. Intuitively, BCD appears to narrow the optimization landscape of the training problem 170 by freezing most of the parameters in each of its sub-problems, potentially eliminating better search 171 directions that can lead to rapid decrease of the objective function. To establish such an intuition 172 formally, we consider the following simple regression problem modeled by a 3-layer neural network: $\min_{\boldsymbol{W}_1, \boldsymbol{W}_2 \boldsymbol{W}_3} H(\boldsymbol{W}_1, \boldsymbol{W}_2, \boldsymbol{W}_3) := \|\boldsymbol{y} - \hat{\boldsymbol{y}}\|_2^2,$ 170

175 176

177

178

179

180 181 182

where $\boldsymbol{z}_2 = \sigma(\boldsymbol{W}_1 \boldsymbol{x}), \quad \boldsymbol{z}_3 = \sigma(\boldsymbol{W}_2 \boldsymbol{z}_2), \quad \hat{\boldsymbol{y}} = \boldsymbol{W}_3 \boldsymbol{z}_3.$ Here, x is a non-zero input vector, y is the true label vector, and $\sigma(z) := \max(0, z)$ is the ReLU activation function. The following proposition states that the optimization landscape of one of the BCD sub-problems is suboptimal in terms of minimizing H.

Proposition 1. (suboptimal landscape of BCD's sub-problem) Define

$$\widetilde{H}^* := \min_{m{W}_2} \|m{y} - \hat{m{y}}\|_2^2, \quad and \quad H^* := \min_{m{W}_1, m{W}_2, m{W}_3} \|m{y} - \hat{m{y}}\|_2^2.$$

Here, \hat{H}^* is the optimal value returned BCD that minimizes H only over block W_2 , while fixing W_1 183 and W_3 . Suppose that the fixed W_3 in BCD has full column rank. If $z_3^* := (W_3^\top W_3)^{-1} W_3^\top y$ has 184 at least one negative entry, then $H^* > H^*$. 185

186 The proof of Proposition 1 is presented in Appendix C.2. We note that the full column rank assump-187 tion of the fixed weight matrix W_3 is mild. In LLMs, the last matrix is typically the tall LM head 188 matrix, which has more rows (vocabulary size) than columns (embedding dimension). In practice, a 189 tall matrix often has full column rank. In addition, it is always feasible to have negative entries in z_3^* 190 by choosing a specific y. Moreover, although we prove such a result for a simple regression prob-191 lem modeled by a 3-layer neural network, we expect it to occur more frequently in the training and 192 finetuning of LLMs, as the training objective of LLMs is much more complicated. Finally, the result in Proposition 1 can be readily generalized to any nonnegative or bounded activation functions, e.g., 193 Swish Ramachandran et al. (2017), SiLU Elfwing et al. (2018), Sigmoid, etc. 194

195 When fixing W_1 and W_3 and training only on the intermediate layer's weight matrix W_2 , the 196 sub-problem of BCD represents a partial landscape of the full optimization problem. Proposition 1 197 demonstrates that this sub-problem might be incapable of finding the optimal value 0. It is important 198 to note that this result does not necessarily imply that BCD cannot find an optimal solution. Recall that BCD is a full parameter optimization method, and the weight matrices W_1 and W_3 will indeed 199 be updated in subsequent block iterations. Therefore, BCD may eventually converge to an optimal 200 solution that achieves 0 function value. However, our analysis reveals that the sub-problem of BCD 201 potentially excludes parts of the optimization landscape that provide search directions toward the 202 optimal solution. This observation suggests that BCD might slow down the convergence speed 203 compared to Adam. 204

Based on these limitations of BCD, we will design a new method to correct the landscape of BCD's 205 sub-problem and simultaneously utilize the computed partial derivatives of the activations properly. 206

207 208

209

4 ACCELERATING BCD VIA LANDSCAPE CORRECTION

210 In this section, we propose a new BCD method with landscape correction to solve the two limitations 211 revealed in Section 3 simultaneously.

212

- 213 4.1 THE BREAD METHOD
- In Section 3, our analysis indicates that the incomplete landscape of BCD's sub-problem may slow 215 down the convergence speed. To address this issue, one immediate approach is to apply Adam to





Figure 1: When W_3 is frozen, the projection error $||\boldsymbol{y} - \hat{\boldsymbol{y}}||_2^2 > 0$. When W_3 is trainable, the learnable space rotates to cover the label \boldsymbol{y} , the error $||\boldsymbol{y} - \hat{\boldsymbol{y}}||_2^2 = 0$.

Figure 2: With rank-1 landscape correction, BREAD converges significantly faster than BCD.

other blocks as well. However, this essentially reverts to using Adam and undermines the memory
 efficient property of the BCD optimization scheme. Fortunately, for the regression problem we
 analyzed in Section 3, we show in the following proposition that a *low-rank landscape correction* is
 sufficient to compensate for the incompleteness of BCD's sub-problem.

Proposition 2 (rank-1 landscape correction). For any non-zero input x and label y, there exists a rank-1 matrix C such that the following problem has optimal value 0:

235

236

227

228

229 230 231

239 240

241

242

243

244

245

246

247

251

where
$$z_2 = \sigma(W_1 z), \quad z_3 = \sigma(W_2 z_2), \quad \hat{y} = (W_3 + C)z_3.$$

min $\|\boldsymbol{u} - \hat{\boldsymbol{u}}\|^2$

The proof of Proposition 2 is shown in Appendix C.2. Motivated by this proposition, we propose to introduce additional *trainable low-rank correction matrices* to the matrices in the frozen inactive blocks $\{W_{\ell'}\}_{\ell' \neq \ell}$, where W_{ℓ} is the current active block. For simplicity, let us assume that each block is a matrix, and our design applies to general layer-wise or other types of partitions as well.

Mathematically, all the inactive block matrices are corrected by low-rank correction matrices as follows:

$$W_{\ell'} + C_{\ell'}$$
 with rank $(C_{\ell'}) \le r, \quad \forall \ell' \ne \ell.$ (4)

To maintain memory efficiency, we use the Burer-Monteiro factorization representation of a low-rank matrix Burer & Monteiro (2003):

$$C = UV, \quad U \in \mathbb{R}^{m \times r}, V \in \mathbb{R}^{r \times n}.$$
 (5)

Following the spirit of adapter Houlsby et al. (2019) and LoRA Hu et al. (2021), we only store the low dimensional matrices U and V rather than C. This approach allows us to store the gradients and optimizer states for these much smaller sized matrices, resulting in only negligible additional memory consumption. We initialize U to zero and initialize V from uniform distribution; see Appendix A for detailed setup.

We note that one can add the landscape correction matrix C to each matrix in the inactive blocks or to part of the matrices, leading to two different variants. Additionally, we may also add a *full-rank* correction matrix C, training it using an on-the-fly SGD method Lv et al. (2023) to maintain the memory efficient feature of BCD. We refer to Section 4.3 for details.

261 Algorithm design. Based on the above developments, we propose accelerating block coordinate 262 descent via landscape correction (BREAD). We present the detailed procedure in Algorithm 1. Sup-263 pose we can add a total of P landscape correction matrices. BREAD first splits the model into D264 blocks, which can be partitioned either in a layer-wise or matrix-wise manner. Then, each block sub-265 problem is approximately solved using K steps of landscape corrected updates. For each landscape 266 corrected update (Algorithm 1, line 9-16), we sample a batch of data in a random reshuffled man-267 ner, and calculate the gradient of both active block and the correction matrices in one backward pass. Then, we update the active block and correction matrices with a single Adam step. Notably, opti-268 mizer states of the correction matrices are accumulated throughout the entire algorithm execution, 269 as they occupy only negligible memory space.

Al	gorithm 1: BREAD: Block coordinate descent with landscape correction.
1	input: model parameters $\{W_{\ell}^0\}_{\ell=1}^L$, number of blocks D, iterations per block K, training
	dataset: $\mathcal{D} = \{(x_j, y_j)\}_{j=1}^n$, batch size B.
2	initialization: block-epoch index $t \leftarrow 0$, the correction matrices $C_i^0 \leftarrow 0, \forall j \in [P]$, and
	the corresponding optimizer states $\tilde{s}_i^0 \leftarrow 0, \forall j \in [P]$.
3	while stopping criterion not meet do
4	generate a block partition $\pi = \{\pi_1, \ldots, \pi_D\}$;
5	repeat for one <i>block-epoch</i> $i \leftarrow 1, \ldots, D$
6	select correction matrices' indices $J \subset [P]$ as in (7);
7	$s_{\pi_i}^{t,0} \leftarrow 0;$ // Re-initialize optimizer states for the active block
8	$egin{array}{c} egin{array}{c} egin{array}$
9	repeat for <i>landscape corrected block updates</i> $k \leftarrow 1, \ldots, K$
10	sample a data batch in random reshuffled manner $\mathcal{D}_B = \{(x_j, y_j)\}_{j=1}^B \sim \mathcal{D};$
11	within one backward pass on the data batch \mathcal{D}_B
12	calculate the active block's grad. $g_i^{t,k}$ and correction matrices' grad. $\tilde{g}_J^{t,k}$;
13	end
14	
	// Update the active block and correction matrices
15	$W^{t,k}_{\pi_i}, s^{t,k}_{\pi_i} \leftarrow AdamStep(W^{t,k-1}_{\pi_i}, g^{t,k}_{\pi_i}, s^{t,k-1}_{\pi_i});$
16	$C_J^{t,k}, \tilde{s}_J^{t,k} \leftarrow AdamStep(C_{\pi_i}^{t,k-1}, \tilde{g}_J^{t,k}, \tilde{s}_J^{t,k-1});$
17	end
18	$W^{t+1}_{\pi^{t+1}} \leftarrow W^{t,K}_{\pi^{t+1}}; C^{t+1}_{I} \leftarrow C^{t,K}_{I}; \tilde{s}^{t+1}_{I} \leftarrow \tilde{s}^{t,K}_{I}; s^{t,K}_{\pi^{t+1}} \leftarrow \text{None};$
19	end
20	$t \leftarrow t+1;$
21	end
22	return parameters $\{W_{\ell}^t\}_{\ell=1}^L$ and correction matrices $\{C_i^t\}_{i=1}^P$.

We now carry out a simple experiment to validate the effectiveness of Algorithm 1 and the insights in Proposition 1 and Proposition 2. Specifically, we train a 3-layer neural network with rank-1 BREAD, with 100 Adam optimization steps for each sub-problem. As shown in Figure 2, rank-1 BREAD (orange) achieves significantly faster convergence than BCD (blue). The loss difference between rank-1 BREAD and BCD is eventually 0.586. This is a significant amount, as the loss is only decreased by 0.06, if we train C only (rank-1 update only, green). These empirical observations verify that the proposed combination of BCD and landscape correction accelerates the convergence of the individual scheme (BCD alone or landscape correction alone).

Convergence result. Below, we show that BREAD is a convergent method under some common assumptions; see Appendix C.1 for the detailed assumptions and analysis.

Theorem 1. (informal) Under certain common assumptions, BREAD with deterministic gradient is a descent method with the following property

$$H(\boldsymbol{W}^{t+1}) - H(\boldsymbol{W}^{t}) \leq -\mathcal{O}(\alpha K) \|\nabla H(\boldsymbol{W}^{t})\|^{2}.$$

Consequently, BREAD finds ε -approximate stationary point with $\mathcal{O}(\varepsilon^{-2})$ iterations.

4.2 MEMORY AND COMPUTATIONAL EFFICIENCY OF BREAD

298 299

300

301

302

303

304

305

306

307

308

309

310

315 316

317

In this section, we analyze the memory and computational cost of BREAD, showing that BREAD only introduces marginal additional costs compared to vanilla BCD.

321 **Memory cost analysis.** To simplify the analysis, we consider a *D*-layer neural network where each 322 layer consists of one matrix with dimensions $\mathbb{R}^{m \times m}$. The rank *r* correction matrix introduces addi-323 tional 2*Dmr* parameters. Since the correction matrix is primarily used for coarse-grained landscape correction, the rank is set to be small, e.g., $r \in [1, 8]$, making the additional memory required almost 324 negligible. In the scenario where r = 4, D = 32, and m = 4096, BREAD only increases memory 325 cost by $\frac{Dr(m+m)}{m^2} \approx 2.6\%$ compared to BCD. 326

Computational cost analysis. We now show that the additional backward cost is also cheap, since 327 the intermediate partial derivatives used for computing the active block's gradient can be directly 328 used for computing correction matrices' gradients, as we have identified in (3). Specifically, the 329 gradient of the correction matrix C_j can be expressed as 330

$$\frac{\partial H}{\partial C_{j}} = \underbrace{\frac{\partial H}{\partial \boldsymbol{z}_{L+1}} \frac{\partial \boldsymbol{z}_{L+1}}{\partial \boldsymbol{z}_{L}} \cdots \frac{\partial \boldsymbol{z}_{j+2}}{\partial \boldsymbol{z}_{j+1}}}_{\text{Computed in (3), }\forall j > \ell} \frac{\partial \boldsymbol{z}_{j+1}}{\partial C_{j}}.$$
(6)

Clearly, when $j \ge \ell$, computing the (stochastic) gradient of C_i only requires additional computation of $\frac{\partial(z_{j+1})}{\partial C_i}$, which is cheap given the low dimensionality after low-rank factorization representation, i.e., $\vec{U}_i = U_i V_i$. We empirically measure the memory and epoch training time in Table 1.

4.3 PRACTICAL VARIANTS OF BREAD

331

332 333

334

335

336

337 338

339

341

349 350 351

360

361 362

363 364

365

366

367 368

369

340 A computational efficient variant. For simplicity, we consider an L-layer neural network where each layer consists of one weight matrix, and our block partition is layer-wise. Based on the deriva-342 tion of (6), evaluating the gradients of the correction matrices is inexpensive for layers $\ell + 1, \ldots, L$. 343 However, the gradient evaluation for layers $1, \ldots, \ell - 1$ is more costly, as it requires calculating 344 $\frac{\partial z_{j+1}}{\partial z_i}$ for $j = 1, \dots, \ell - 1$. These intermediate partial derivatives of the activations are not com-345 puted during the backpropagation to the active layer ℓ . Therefore, one computationally efficient 346 variant of BREAD is to add correction matrices only for layers $\ell + 1, \ldots, L$. This leads to two 347 strategies of selecting correction matrices: 348

$$J = \begin{cases} [P], & \text{if use BREAD} \\ \{j \mid j \in [P], \ C_j \text{ corrects layers } \ell + 1, \dots, L\}, & \text{if use BREAD-partial} \end{cases}$$
(7)

A full-rank memory efficient variant. The previous implementation uses low-rank matrices for 352 landscape correction. Alternatively, one can apply a potentially full-rank linear transformation to 353 correct features. This can be achieved through a memory efficient on-the-fly SGD update on the 354 inactive blocks. Specifically, due to the compositional structure of neural networks, the gradient of 355 the model is computed from the deep layers to the shallow layers. The strategy is to perform an SGD 356 update on a matrix whenever its (stochastic) gradient is available, and then immediately discard 357 the corresponding gradient after the update. We term this approach BREAD-SGD. It introduces 358 additional memory cost for storing the gradient of the largest matrix, but this overhead is usually 359 negligible. We formally present this variant in Algorithm 2.

We compare the performance of these two variants of BREAD in Section 5.4.

NUMERICAL EXPERIMENTS 5

We evaluate the proposed BREAD in finetuning Llama 3.1-8B and Llama 3.1-70B model on math finetuning and instruction tuning tasks, comparing its memory cost, time cost and downstream performance with full training algorithm and memory efficient baselines.

5.1 Setup

370 We begin by introducing the experimental setup. 371

Baselines. We compare BREAD with 1) BAdam Luo et al. (2024), which applies vanilla BCD 372 algorithm with Adam as the inner solver; 2) LoRA Hu et al. (2021), which freezes the pre-trained 373 weight and only updates the injected low-rank adapters; 3) Galore Zhao et al. (2024), which projects 374 the gradient into low-rank spaces for reducing the memory cost; 4) Adam Kingma (2014), which 375 serves as the full parameter training baseline. 376

Math finetuning. We finetune the Llama 3.1-70B and Llama 3.1-8B models on MathInstruct 377 dataset Yue et al. (2023) for 3 epochs, which contains 260K questions that covers wide range of fields

Model Method		Peak Memory Cost	Epoch GPU Hour	GPU #	
	Adam (estimated)	1260 GB+	_	16+ A100-80G	
	LoRA	296.8 GB	213.1	8 A100-40G	
Llama 3.1-70B	B BAdam	276.2 GB	119.0	8 A100-40G	
	BREAD	288.6 GB	212.4	8 A100-40G	
	BREAD-partial	288.6 GB	152.7	8 A100-40G	
	Adam	208.2 GB	37.3	8 A100-40G	
	Galore	40.5 GB	10.3	1 A100-80G	
Llomo 2 1 PD	LoRA	25.0 GB	6.0	1 A100-40G	
Liallia 5.1-6D	BAdam	21.8 GB	3.3	1 A100-40G	
	BREAD	23.2 GB	5.8	1 A100-40G	
	BREAD-partial	23.2 GB	4.0	1 A100-40G	

Table 1: Memory footprint and time cost for finetuning models on MathInstruct.

in mathematics. The finetuned models are evaluated on 4 in-domain mathematical benchmarks, i.e.,
GSM8K, MATH, NumGLUE, and AQuA Cobbe et al. (2021); Hendrycks et al. (2021); Mishra et al.
(2022); Ling et al. (2017), and 1 out-of-domain mathematical benchmarks, i.e., SimulEq Koncel-Kedziorski et al. (2016). The evaluations are based on 0-shot prompt and 4-shot chain-of-thought
prompt, respectively. Due to the limited computational resource, we do not include the Adam's
results for 70B model. Since there is no model parallel implementation released for Galore by the
finish of the manuscript, we are unable to report its 70B results as well.

Instruction tuning. We perform supervise finetuning on the Llama 3.1-8B model using Alpaca-GPT4 dataset Peng et al. (2023), which contains 52K questions and corresponding GPT-4 generated answers. The model is evaluated on MT-bench Zheng et al. (2023) for examining the model's instruction-following capability.

Preference optimization. After the instruction tuning, we further align the tuned model using direct
 preference optimization (DPO) Rafailov et al. (2024) on Ultrafeedback dataset Cui et al. (2023). To
 compare with the baseline optimization methods more comprehensively, we report the evaluation
 results of using Adam instruction tuned model as base model, and using the same optimization
 method for both phases.

All the experiments are run for 3 epochs. The reported scores are the best one among checkpoints at epoch 1, 2, 3. The detailed hyper-parameters are presented in Appendix A.

412

414

391 392

413 5.2 MEMORY AND TIME COST MEASURE

In Table 1, we empirically measure the peak memory cost and one epoch's time cost of finetuning models on the MathInstruct dataset. The GPU hour is calculated as the training time × GPU number.

We set the LoRA rank to 64 to keep its number of trainable parameters (0.83 billion) close to a single
block of BREAD (0.86 billion). Compared with LoRA, the proposed BREAD consumes slightly
lower memory and costs about the same amount of training time. The computational-efficient variant
BREAD-partial requires significantly less training time than BREAD, and is comparable to BAdam.

Remark. The reported memory cost is higher than the theoretical value, especially for the 70B model's experiments which requires distributed training. This additional memory cost arises from storing activation values and computational buffers, e.g. the gradient buffer for performing reduce scatter operation. Furthermore, the training time may have slight fluctuations for different runs.

425 426

427

5.3 FINETUNING PERFORMANCE

Math finetuning. The evaluation results on math benchmarks are shown in Table 2. For the 8B model's finetuning, BREAD beats all the other 3 memory efficient baselines in both 0-shot and 4-shot average score. Under the 0-shot setting, BREAD even outperforms Adam baseline by 1.1. For the finetuning of 70B model, BREAD outperforms BAdam in all tasks, demonstrating the effective-ness of landscape correction. Furthermore, BREAD beats LoRA in 8 out of 10 tasks.

Base model: Llama 3.1-8B												
Method	GSM8K		MA	MATH		NumGLUE		SimulEq		AQuA		vg.
	0-shot	4-shot	0-shot	4-shot	0-shot	4-shot	0-shot	4-shot	0-shot	4-shot	0-shot	4-shot
Base model	17.8	52.5	8.6	23.2	25.7	40.6	12.2	28.8	19.3	43.7	16.7	37.8
Adam	62.3	64.9	17.4	22.9	56.4	56.8	28.6	33.5	44.9	52.8	41.9	46.2
Galore	46.7	57.2	16.2	22.9	42.8	45.0	28.7	32.3	47.8	48.4	36.4	41.2
LoRA	48.7	58.1	13.7	23.0	34.6	54.4	29.6	29.0	47.3	50.3	34.8	43.0
BAdam	53.9	58.3	17.2	23.6	53.7	57.2	32.5	32.8	50.4	49.6	41.5	44.3
BREAD	57.0	57.6	20.0	23.7	55.9	58.2	32.5	32.8	49.6	50.0	43.0	44.5
				Base	e model:	Llama 3	3.1-70B					
Method	GSN	M8K	MA	TH	Num	GLUE	Sim	ulEq	AÇ	JuA	A	vg.
	0-shot	4-shot	0-shot	4-shot	0-shot	4-shot	0-shot	4-shot	0-shot	4-shot	0-shot	4-shot
Base model	58.8	79.4	24.9	41.4	43.7	55.8	26.3	38.1	52.0	64.2	41.1	51.2
LoRA	83.8	82.0	41.7	44.2	70.4	69.0	40.3	48.8	61.4	65.8	59.5	62.0
BAdam	81.4	82.9	40.3	43.8	68.1	69.7	50.0	52.7	65.3	70.1	61.0	63.8
BREAD	83.4	84.2	41.4	44.7	73.1	74.4	51.3	56.8	68.3	70.5	63.5	66.1

Table 2: Math evaluation results for models finetuned on MathInstruct dataset.

Method	S	FT	DPO			
	GPT-4	GPT-40	GPT-4	GPT-40		
Base model	6.07	5.18	6.63	5.66		
Adam	6.63	5.66	7.83	6.18		
LoRA	6.52	5.60	7.48	5.95		
Galore	6.33	5.48	6.99	5.93		
BAdam	6.53	5.57	7.63	6.14		
BREAD	6.77	5.82	7.68	6.31		



Table 3: MT-bench scores of different methods finetuning Llama 3.1-8B.



Instruction tuning and DPO. We report the MT-bench score evaluated by both GPT-4 and GPT-40 models in Figure 3. After SFT, the MT-bench score of all baseline approaches improves over the base model. BREAD achieves the highest scores in both evaluations, which are even higher than Adam, demonstrating the effectiveness of landscape correction. Based on the model finetuned by Adam, we further align the model using direct preference optimization (DPO). Notably, BREAD achieves the highest evaluation score by GPT-40 model.

5.4 ABLATION ON VARIANTS OF BREAD

We present 1-epoch training loss of BREAD and its variants in Figure 3. For reference, we also display the loss of BAdam. Here, BREAD-partial means that only deeper layers of the active block will be updated; the BREAD-SGD applies on-the-fly SGD update for all the frozen blocks; the BREAD-SGD-partial incorporates the feature of both, updating frozen blocks that are deeper than the active block using SGD. Since SGD usually requires higher learning rate, we scale the SGD's learning rate up by 10 times compared to the update of the active block. As we can see, all the variants converge faster than the baseline method BAdam, justifying the effectiveness of landscape correction. The BREAD outperforms BREAD-partial since it optimizes all the correction matrices in each iteration. The BREAD-SGD variants converge faster than the low-rank ones, which may attribute to the higher learning rate of the correction matrices and the high-rank update. Notably, the BREAD-SGD-partial exhibits similar convergence as BREAD-SGD.

486 5.5 Additional Experimental Results

488 We conduct additional ablation study on the hyperparameters of BREAD, and display the results in Appendix D. Below, we summarize the experimental results: 1) Effect of K. We observe a 489 consistent improved convergence speed when increasing K, which may due to that Adam can ag-490 gregate more historical information with larger K. 2) Effect of r. BREAD with rank-2 correction 491 significantly outperforms BAdam, and can be further accelerated with larger choices of r. 3) Effect 492 of ordering strategies. Different block ordering yields similar convergence behavior. 4) Conver-493 gence versus time. In terms of wall-clock time, BREAD-SGD-partial yields the fastest convergence 494 among the BREAD variants. 495

- 496 6 RELATED WORKS
- 497

498 Block coordinate descent method. Block coordinate descent (BCD) is a classic optimization 499 paradigm that dates back at least to Hildreth (1957). It has gained popularity in recent years, due 500 to its scalability and efficiency for many machine learning applications (Nesterov, 2012; Richtárik 501 & Takáč, 2014; Peng & Vidal, 2023; Ding et al.; Peng & Yin, 2024). The community seems to 502 converge to a consensus that, in order for BCD to be efficient, the problem it optimizes needs to 503 possess the so-called coordinate-friendly structure (Shi et al., 2016). Nevertheless, deep networks are of a compositional nature and not coordinate-friendly, which is perhaps why recent surveys or 504 books have never mentioned training deep networks as an application of BCD (Wright, 2015; Shi 505 et al., 2016; Beck, 2017; Wright & Recht, 2022; Sayed, 2022). Recently, BAdam Luo et al. (2024) 506 was proposed to finetune LLMs based on the BCD framework, where each block sub-problem is 507 approximately solved using several Adam steps. Although BAdam achieved preliminary success, it 508 is based on the vanilla BCD framework and shares the fundamental limitations we revealed in this 509 work. In light of these, we believe identifying the limitations of BCD for LLM fintuning and fixing 510 them entail certain insights, and this is what makes our contributions non-trivial and valuable. 511

Memory efficient finetuning. To address memory issue, multiple variants have been proposed. Pa-512 rameter efficient finetuning (PEFT) methods achieve memory efficiency by only training small por-513 tion of (possibly extra) parameters while freezing most of the others, such as Adapter tuning Houlsby 514 et al. (2019), prompt tuning, and prefix tuning Lester et al. (2021); Li & Liang (2021). Low-rank 515 adaptation (LoRA) is perhaps the most popular technique that approximates model updates using 516 two smaller, trainable low-rank matrices Hu et al. (2021). LoRA' variants have been proposed to 517 address its rank constraints and further reducing the memory cost Lialin et al. (2024); Xia et al. 518 (2024); Dettmers et al. (2023). Galore Zhao et al. (2024) projects the gradient into low-rank space 519 so that it does not need to store the full gradient and optimizer states in the memory. LOMO updates 520 parameters in real time during the backpropagation process Lv et al. (2023), so that one can perform SGD without store stochastic gradients. MeZO offers an alternative by approximating SGD using 521 only forward passes Malladi et al. (2023), drawing from zeroth-order optimization that estimates 522 stochastic gradients through the difference in function values. While this paper addresses the same 523 application as these methods, they remain orthogonal to the proposed approaches. They can function 524 as lightweight updates in the frozen layers for landscape correction. 525

7 CONCLUSIONS AND DISCUSSIONS

526 527 528

This paper investigates the application of a classic optimization method, known as BCD, to the finetuning of LLMs. We pinpoint two primary shortcomings of the standard BCD approach when applied to deep neural networks: the unnecessary computational overhead during backpropagation, and the misguiding optimization landscape caused by frozen blocks. To overcome these challenges, we introduce a new method termed BREAD, which unfreezes the inactive blocks and updates them in a lightweight manner. Our experimental results demonstrate that BREAD significantly enhances downstream task performance while maintaining the original BCD algorithm's memory and computational efficiency.

For future research, it would be intriguing to explore the potential of BCD in the (continual) pretraining phase of LLMs. Moreover, it is not necessary to apply feature correction at each iteration.
Exploring the frequency of updating correction matrices is another meaningful direction, which can
further save the computational cost.

540 REFERENCES

542 543	Amir Beck. <i>First-Order Methods in Optimization</i> . Society for Industrial and Applied Mathematics, 2017.
544 545 546	Samuel Burer and Renato DC Monteiro. A nonlinear programming algorithm for solving semidefinite programs via low-rank factorization. <i>Mathematical programming</i> , 95(2):329–357, 2003.
547 548 549	Kai-Wei Chang and Dan Roth. Selective block minimization for faster convergence of limited mem- ory large-scale linear models. In <i>Proceedings of the 17th ACM SIGKDD international conference</i> <i>on Knowledge discovery and data mining</i> , pp. 699–707, 2011.
550 551 552	Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to solve math word problems. <i>arXiv preprint arXiv:2110.14168</i> , 2021.
554 555 556	Ganqu Cui, Lifan Yuan, Ning Ding, Guanming Yao, Wei Zhu, Yuan Ni, Guotong Xie, Zhiyuan Liu, and Maosong Sun. Ultrafeedback: Boosting language models with high-quality feedback. <i>arXiv</i> preprint arXiv:2310.01377, 2023.
557 558 559	Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. QLoRA: Efficient finetuning of quantized LLMs. <i>Advances in Neural Information Processing Systems</i> , 36, 2023.
560 561	Lisang Ding, Ziang Chen, Xinshang Wang, and Wotao Yin. Efficient algorithms for sum-of- minimum optimization. In <i>Forty-first International Conference on Machine Learning</i> .
562 563 564	Stefan Elfwing, Eiji Uchibe, and Kenji Doya. Sigmoid-weighted linear units for neural network function approximation in reinforcement learning. <i>Neural networks</i> , 107:3–11, 2018.
565 566 567	Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In <i>Proceedings of the IEEE international conference on computer vision</i> , pp. 1026–1034, 2015.
568 569 570	Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset. <i>arXiv</i> preprint arXiv:2103.03874, 2021.
572 573	Clifford Hildreth. A quadratic programming procedure. <i>Naval Research Logistics Quarterly</i> , 4(1): 79–85, 1957.
574 575 576 577	Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe, An- drea Gesmundo, Mona Attariyan, and Sylvain Gelly. Parameter-efficient transfer learning for NLP. In <i>International Conference on Machine Learning</i> , pp. 2790–2799. PMLR, 2019.
578 579 580	Cho-Jui Hsieh, Kai-Wei Chang, Chih-Jen Lin, S Sathiya Keerthi, and Sellamanickam Sundarara- jan. A dual coordinate descent method for large-scale linear svm. In <i>Proceedings of the 25th</i> <i>international conference on Machine learning</i> , pp. 408–415, 2008.
581 582 583 584	Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. LoRA: Low-rank adaptation of large language models. <i>arXiv preprint arXiv:2106.09685</i> , 2021.
585 586	Diederik P Kingma. Adam: A method for stochastic optimization. <i>arXiv preprint arXiv:1412.6980</i> , 2014.
587 588 589 590 591	Rik Koncel-Kedziorski, Subhro Roy, Aida Amini, Nate Kushman, and Hannaneh Hajishirzi. Mawps: A math word problem repository. In <i>Proceedings of the 2016 conference of the north</i> <i>american chapter of the association for computational linguistics: human language technologies</i> , pp. 1152–1157, 2016.
592 593	Brian Lester, Rami Al-Rfou, and Noah Constant. The power of scale for parameter-efficient prompt tuning. In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pp. 3045–3059, 2021.

594 595 596	Xiang Lisa Li and Percy Liang. Prefix-tuning: Optimizing continuous prompts for generation. In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing</i> , pp. 4582–4597, 2021.
597 598 599 600	Vladislav Lialin, Sherin Muckatira, Namrata Shivagunde, and Anna Rumshisky. ReLoRA: High- rank training through low-rank updates. In <i>The Twelfth International Conference on Learning</i> <i>Representations</i> , 2024.
601 602 603	Wang Ling, Dani Yogatama, Chris Dyer, and Phil Blunsom. Program induction by rationale gener- ation: Learning to solve and explain algebraic word problems. <i>arXiv preprint arXiv:1705.04146</i> , 2017.
604 605 606	Qijun Luo, Hengxu Yu, and Xiao Li. Badam: A memory efficient full parameter training method for large language models. <i>arXiv preprint arXiv:2404.02827</i> , 2024.
607 608 609	Kai Lv, Yuqing Yang, Tengxiao Liu, Qinghui Gao, Qipeng Guo, and Xipeng Qiu. Full parameter fine-tuning for large language models with limited resources. <i>arXiv preprint arXiv:2306.09782</i> , 2023.
610 611 612	Sadhika Malladi, Tianyu Gao, Eshaan Nichani, Alex Damian, Jason D Lee, Danqi Chen, and Sanjeev Arora. Fine-tuning language models with just forward passes. <i>Advances in Neural Information Processing Systems</i> , 36, 2023.
613 614 615 616	Paulius Micikevicius, Sharan Narang, Jonah Alben, Gregory Diamos, Erich Elsen, David Garcia, Boris Ginsburg, Michael Houston, Oleksii Kuchaiev, Ganesh Venkatesh, et al. Mixed precision training. arXiv preprint arXiv:1710.03740, 2017.
617 618 619	Swaroop Mishra, Arindam Mitra, Neeraj Varshney, Bhavdeep Sachdeva, Peter Clark, Chitta Baral, and Ashwin Kalyan. Numglue: A suite of fundamental yet challenging mathematical reasoning tasks. <i>arXiv preprint arXiv:2204.05660</i> , 2022.
620 621 622	Yu Nesterov. Efficiency of coordinate descent methods on huge-scale optimization problems. <i>SIAM Journal on Optimization</i> , 22(2):341–362, 2012.
623 624	Baolin Peng, Chunyuan Li, Pengcheng He, Michel Galley, and Jianfeng Gao. Instruction tuning with gpt-4. <i>arXiv preprint arXiv:2304.03277</i> , 2023.
625 626 627	Liangzu Peng and René Vidal. Block coordinate descent on smooth manifolds: Convergence theory and twenty-one examples. Technical report, arXiv:2305.14744v3 [math.OC], 2023.
628 629	Liangzu Peng and Wotao Yin. Block acceleration without momentum: On optimal stepsizes of block gradient descent for least-squares. <i>arXiv preprint arXiv:2405.16020</i> , 2024.
630 631 632	Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. <i>Advances in Neural Information Processing Systems</i> , 36, 2024.
633 634 635	Prajit Ramachandran, Barret Zoph, and Quoc V Le. Searching for activation functions. <i>arXiv</i> preprint arXiv:1710.05941, 2017.
636 637	Peter Richtárik and Martin Takáč. Iteration complexity of randomized block-coordinate descent methods for minimizing a composite function. <i>Mathematical Programming</i> , 144(1):1–38, 2014.
638 639 640	Ali H Sayed. Inference and Learning from Data: Foundations, volume 1. Cambridge University Press, 2022.
641 642	Hao-Jun Michael Shi, Shenyinying Tu, Yangyang Xu, and Wotao Yin. A primer on coordinate descent algorithms. Technical report, arXiv:1610.00040v2 [math.OC], 2016.
643 644 645	Eran Treister and Javier S Turek. A block-coordinate descent approach for large-scale sparse inverse covariance estimation. <i>Advances in neural information processing systems</i> , 27, 2014.
646 647	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. <i>Advances in Neural Information Processing Systems</i> , 2017.

- Stephen J Wright. Coordinate descent algorithms. *Mathematical Programming*, 151(1):3–34, 2015.
- Stephen J Wright and Benjamin Recht. Optimization for data analysis. Cambridge University Press, 2022.
- Wenhan Xia, Chengwei Qin, and Elad Hazan. Chain of LoRA: Efficient fine-tuning of language models via residual learning. arXiv preprint arXiv:2401.04151, 2024.
- Hsiang-Fu Yu, Cho-Jui Hsieh, Kai-Wei Chang, and Chih-Jen Lin. Large linear classification when data cannot fit in memory. ACM Transactions on Knowledge Discovery from Data (TKDD), 5(4): 1-23, 2012.
- Xiang Yue, Xingwei Qu, Ge Zhang, Yao Fu, Wenhao Huang, Huan Sun, Yu Su, and Wenhu Chen. Mammoth: Building math generalist models through hybrid instruction tuning. arXiv preprint arXiv:2309.05653, 2023.
 - Jiawei Zhao, Zhenyu Zhang, Beidi Chen, Zhangyang Wang, Anima Anandkumar, and Yuandong Tian. Galore: Memory-efficient llm training by gradient low-rank projection. arXiv preprint arXiv:2403.03507, 2024.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. Judging llm-as-a-judge with mt-bench and chatbot arena. Advances in Neural Information Processing Systems, 36:46595-46623, 2023.
- Yaowei Zheng, Richong Zhang, Junhao Zhang, Yanhan Ye, Zheyan Luo, and Yongqiang Ma. LlamaFactory: Unified efficient fine-tuning of 100+ language models. arXiv preprint arXiv:2403.13372, 2024. URL http://arxiv.org/abs/2403.13372.

702 A DETAILED EXPERIMENTAL SETUP

We introduce the detailed hyperparameters and experimental setup in this section.

Global setup. For all the experiments in math finetuning, instruction tuning and direct preference 706 optimization, we fix the effective batch size to be 16 and train the model for 3 epochs. We use DeepSpeed ZeRO-3 to implement all the experiments that require distributed training (shown in 708 Table 1). For all the experiments, we apply gradient checkpointing to reduce the memory cost for storing activation values. We use mixed-precision training with BFloat 16 as the low-precision 710 datatype except for Galore, where we follow the setup in its paper, using pure BFloat 16 and 8-711 bit Adam optimizer for reducing the memory cost. Since the downstream tasks' performance of 712 the language model have high variability, we grid search learning rate from {1e-6, 1e-5, 5e-5} with 713 cosine learning rate schedule, and report the best result among the checkpoints at the end of epoch 1, 714 2, 3. The grid search is not extensive due to our limited computation resources. The implementation of BAdam, Galore, LoRA are based on LLama-Factory Zheng et al. (2024). 715

716 Math finetuning. We randomly select 100,000 samples from the MathInstruct dataset and finetune 717 all the models using the same samples. The benchmarks scores are evaluated using the MAm-718 moTH's repository¹ (without using program-of-thought). The rank of correction matrices U and V 719 for BREAD and BREAD-partial is set to 8. We initialize U as zero, and initialize V from the Kaim-720 ing uniform distribution He et al. (2015), i.e. $\left(-\mathcal{U}(\frac{\sqrt{6}}{r}), \mathcal{U}(\frac{\sqrt{6}}{r})\right)$. The rank of LoRA is set to 80 721 and 64 for finetuning Llama 3.1-8B and Llama 3.1-70B, respectively, so that the trainable parameter 722 number of LoRA is close to that of one BAdam/BREAD's active block. We follow the conventional 723 setup to set the LoRA scaling factor $\alpha = 2 \times$ LoRA rank. We set Galore's rank to be 256, with the 724 period of re-calculating the projection matrix being 256. We set K = 100 for BAdam and BREAD.

Instruction tuning and direct preference optimization. The evaluation of MT-bench score is based on FastChat Zheng et al. (2023) using both GPT-4 and GPT-40 API. The maximum sequence is set to 1024 and 2048 for the experiments on Alpaca-GPT4 and UltraFeedback, respectively. For BAdam and BREAD, we solve each block sub-problem for 128 steps, i.e. K = 128.

- 730
- 731 732
- 733

734

735 736 737

> 739 740

738

741 742

743

744 745 746

747 748 749

750



¹https://github.com/TIGER-AI-Lab/MAmmoTH

756 B BREAD-SGD ALGORITHM

We introduce a variant of the BREAD algorithm, termed BREAD-SGD, which employs Adam for updating the active block and on-the-fly SGD for the inactive block. The detailed procedure is outlined in Algorithm 2. Analogous to BREAD, BREAD-SGD partitions the model into *D* distinct blocks and combines the gradient computation and update steps into a singular operation. Specifically, gradients are calculated on a layer-by-layer basis; active layers are updated using Adam, while inactive layers undergo a single SGD step. Once an inactive layer is updated, its gradient is discarded to enhance memory efficiency.

Algorithm 2: BREAD-SGD

765 766

767	1 1	input: model parameters $\{W_{\ell}^0\}_{\ell=1}^L$, number of blocks D, iterations per block K, step size
768	(of inactive blocks β .
769	2 i	initialization: block-epoch index $t \leftarrow 0$, and the corresponding optimizer states
770	į	$ ilde{m{s}}^0_j \leftarrow m{0}, \ orall j \in [P].$
771	3	while stopping criterion not meet do
772	4	generate a block partition $\pi = \{\pi_1, \ldots, \pi_D\}$;
773	5	repeat for one <i>block-epoch</i> $i \leftarrow 1, \dots, D$
774	6	select correction matrices' indices $J \subset [P]$;
775	7	$s_{\pi_i}^{t,0} \leftarrow 0;$ // Re-initialize Adam optimizer states
776	8	$egin{array}{c} W^{t,0}_{\pi_i} \leftarrow W^t_{\pi_i}; \end{array}$
777	9	repeat for landscape corrected block updates $k \leftarrow 1, \ldots, K$
778	10	sample a data batch in random reshuffled manner $\mathcal{D}_B = \{(x_j, y_j)\}_{j=1}^B \sim \mathcal{D};$
779	11	within one backward pass on the data batch \mathcal{D}_B
780		// Update the active block
781	12	$g_{\pi_i}^{t,k} \leftarrow \frac{\partial H}{\partial \mathbf{W}};$
782	12	$W^{t,k} \stackrel{\mathbf{s}^{t,k}}{\overset{\mathbf{s}^{t,k}}{\leftarrow}} \operatorname{AdamStep}(W^{t,k-1} \stackrel{\mathbf{a}^{t,k}}{\overset{\mathbf{s}^{t,k-1}}}{\overset{\mathbf{s}^{t,k-1}}}{\overset{\mathbf{s}^{t,k-1}}{\overset{\mathbf{s}^{t,k-1}}}{\overset{\mathbf{s}^{t,k-1}}}{\overset{\mathbf{s}^{t,k-1}}}{\overset{\mathbf{s}^{t,k-1}}}{\overset{\mathbf{s}^{t,k-1}}}{\overset{\mathbf{s}^{t,k-1}}}{\overset{\mathbf{s}^{t,k-1}}}{\overset{\mathbf{s}^{t,k-1}}}{\overset{t}}{\overset{t}^{t,k-1}}}}}}}}}}}}}}}}}}}}}}}}$
783	13	$(\gamma_{\pi_i}, \sigma_{\pi_i}), (\gamma_{\pi_i}, \gamma_{\pi_i}), (\gamma_{\pi_i}, \sigma_{\pi_i}), (\gamma_{\pi_i}), (\gamma_{\pi_i}, \sigma_{\pi_i}), (\gamma_{\pi_i}), (\gamma_$
784		// Correct inactive blocks
785	15	for $\ell \in J$ do
786	16	$egin{array}{c c c c c c c c c c c c c c c c c c c $
787	17	$oldsymbol{W}_{e}^{t,k} \leftarrow oldsymbol{W}_{e}^{t,k-1} - eta oldsymbol{g}_{e}^{t,k-1};$
788	18	$q_{\ell}^{t,k-1} \leftarrow \text{None};$
789	19	end
790	20	end
791	21	end
792	22	$W^{t+1}_{\pi_{\star}} \leftarrow W^{t,K}_{\pi_{\star}}; \ s^{t,K}_{\pi_{\star}} \leftarrow \text{None};$
793	23	end
794	24	$t \leftarrow t + 1;$
795	25	end
796	26 1	return parameters $\{\boldsymbol{W}_{\ell}^t\}_{\ell=1}^L$.
797		
798		
799		
800		
801		
802		

803

804 805

806

806

807 808

CONVERGENCE RESULT AND ADDITIONAL PROOFS C

C.1 CONVERGENCE OF BREAD

In this section, we establish a preliminary convergence result for the proposed algorithms. For brevity of expression, we analyze the BREAD-SGD. The analysis of BREAD and BREAD-partial follows the similar strategy. We begin by introducing the assumptions.

Assumption 1. The function H(W) is L-smooth on W and L_i -smooth on block W_i for blocks $i = 1, \cdots, D$:

$$\|\nabla_{\boldsymbol{W}}H(\boldsymbol{W}) - \nabla_{\boldsymbol{W}}H(\bar{\boldsymbol{W}})\| \le L\|\boldsymbol{W} - \bar{\boldsymbol{W}}\|.$$
(8)

The block-wise smoothness (9) can be naturally induced by the function smoothness (8).

Assumption 2. There block derivatives $\nabla_{\mathbf{W}_i} H(\mathbf{W})$ are uniformly bounded by a constant G:

$$\|\nabla_{\boldsymbol{W}_i} H(\boldsymbol{W})\| \leq G, i = 1, \cdots, D$$

 $\left\| \frac{\partial H}{\partial \boldsymbol{W}_i} \right\|_{\boldsymbol{W}_i^1} - \frac{\partial H}{\partial \boldsymbol{W}_i} \right\|_{\boldsymbol{W}_i^2} \leq L_i \| \boldsymbol{W}_i^1 - \boldsymbol{W}_i^2 \|, i = 1, \cdots, D.$

Theorem 2 (Descent of BREAD). Under Assumption 1 and Assumption 2, the Algorithm 2 with deterministic gradient achieves the following descent after each block-epoch of updates:

$$H(\boldsymbol{W}^{t+1}) - H(\boldsymbol{W}^{t}) \le -\mathcal{O}(\alpha K) \|\nabla H(\boldsymbol{W}^{t})\|^{2},$$
(10)

(9)

under the step size choice $\alpha \leq \min\{\frac{\lambda}{2LK^2}, \frac{\lambda^2}{24LKG}, \frac{LG}{2}, \frac{108D}{16G^2K^3}\}$.

By telescoping (10) from $t = 1, \dots, T$, and divide each side by T, we obtain the sample complexity of $\|\nabla H(\mathbf{W}^t)\|^2 = \mathcal{O}(K/T)$. The proof of Theorem 2 is based on the following Lemma, which establishes the descent property of one block sub-problem.

Lemma 1. Under Assumption 1 and Assumption 2, the Algorithm 2 update yields the following *approximate descent property:*

$$H(\boldsymbol{W}_{i}^{t}) - H(\boldsymbol{W}_{i-1}^{t}) \leq -\frac{\alpha K}{4G} \|\nabla H(\boldsymbol{W}_{i-1}^{t})\|_{2}^{2}$$

Proof. The proof of Lemma 1 follows the analysis framework of Luo et al. (2024), except that we need to analyze the descent property of the update in correction matrices. We define some notations for the ease of expression. At the beginning of block epoch t, block sub-problem i, let $W_{i,i}^t$ be the parameter of block j, W_i^t be the full parameter. Based on the smoothness property in Assumption 1, we have

$$H(\boldsymbol{W}_{i}^{t}) - H(\boldsymbol{W}_{i-1}^{t}) \leq \left\langle \nabla H(\boldsymbol{W}_{i-1}^{t}), \boldsymbol{W}_{i}^{t} - \boldsymbol{W}_{i-1}^{t} \right\rangle + \frac{L}{2} \left\| \boldsymbol{W}_{i}^{t} - \boldsymbol{W}_{i-1}^{t} \right\|_{2}^{2}$$

$$= \sum_{\substack{j \in [D] \setminus i \\ I_{1}}} \left\langle \nabla_{j} H(\boldsymbol{W}_{i-1}^{t}), \boldsymbol{W}_{j,i}^{t+1} - \boldsymbol{W}_{j,i-1}^{t+1} \right\rangle + \sum_{\substack{j \in [D] \setminus i \\ I_{2}}} \frac{L_{j}}{2} \left\| \boldsymbol{W}_{j,i}^{t+1} - \boldsymbol{W}_{j,i-1}^{t+1} \right\|_{2}^{2}$$

$$+ \left\langle \nabla_{i} H(\boldsymbol{W}_{i-1}^{t}), \boldsymbol{W}_{i,i}^{t+1} - \boldsymbol{W}_{i,i-1}^{t+1} \right\rangle + \frac{L_{i}}{2} \left\| \boldsymbol{W}_{i,i}^{t+1} - \boldsymbol{W}_{i,i-1}^{t+1} \right\|_{2}^{2}$$
(11)

The analysis of the last two terms directly follows the strategy in BAdam. We now analyze the first two terms. Let $\beta = c_{\beta}\alpha$ be the step size for the correction matrices. Define the gradient bias term of the block sub-problem as $e_{j,i}^t := \frac{1}{K} \sum_{k=1}^K \nabla_j H(W_{i-1}^{t,k}) - \nabla_j H(W_{i-1}^t)$, where $W_i^{t,k}$ is the parameter at block epoch *j*, sub-problem *i*, inner Adam step *k*, we have

$$I_1 = K\beta \sum_j \left\langle \nabla_j H(\boldsymbol{W}_{i-1}^t), -\nabla_j H(\boldsymbol{W}_{i-1}^t) + \boldsymbol{e}_{j,i}^t \right\rangle$$
$$\stackrel{(i)}{\leq} -\sum_j \frac{3}{4} K\beta \|\nabla_j H(\boldsymbol{W}_{i-1}^t)\|_2^2 + \sum_j 4\beta K \|\boldsymbol{e}_{j,i}^t\|_2^2$$

where (i) uses Young's inequality, (ii) applies the definition of $e_{j,i}^t$, (iii) uses the fact that $(\sum_{i=1}^{K} a_i)^2 \leq K a_i^2$, (iv) is due to the Assumption 1, (v) is due to

$$\sum_{k=1}^{K} \| \boldsymbol{W}_{j,i-1}^{t,k} - \boldsymbol{W}_{j,i-1}^{t} \|_{2}^{2} \leq \sum_{k=1}^{K} \left(\sum_{m=1}^{k-1} \| W_{j,i-1}^{t,m+1} - W_{j,i-1}^{t,m} \| \right)^{2}$$
$$\leq \sum_{k=1}^{K} ((k-1)\beta G)^{2}$$

and the last inequality is due to $\beta \leq \sqrt{\frac{3D}{16K^3}}$, which is ensured by $\alpha \leq \frac{108D}{16G^2K^3}$.

Similar to the analysis above, we can bound I_2 based on Assumption 2:

$$I_{2} = \sum_{j} \frac{L_{j}}{2} \beta^{2} \left\| \sum_{k=1}^{K} \nabla_{j} H(\boldsymbol{W}_{i-1}^{t,k}) - \nabla_{j} H(\boldsymbol{W}_{i-1}^{t}) \right\|_{2}^{2}$$
(13)

 $=\beta^2 G^2 \frac{(K-1)K(2K-1)}{6} \leq \frac{\beta^2 G^2 K^3}{3},$

$$\leq \sum_{j} \frac{L_j}{2} \beta^2 K 4 G^2 \tag{14}$$

$$\leq 2G^2 K L \beta^2 \leq \frac{1}{4} K \beta G^2,\tag{15}$$

where the first inequality is due to Assumption 2, the second inequality is based on the fact that $L > L_j, \forall j \in [D]$, and the last inequality is due to $\beta \leq L/8$. Combine (12), (15) and (11), and given that $\beta = \alpha/4G$, we have

$$H(\boldsymbol{W}_{i}^{t}) - H(\boldsymbol{W}_{i-1}^{t}) \leq -\frac{1}{4} K \beta \|\nabla H(\boldsymbol{W}_{i-1}^{t})\|_{2}^{2} + \langle \nabla_{i} H(\boldsymbol{W}_{i-1}^{t}), \boldsymbol{W}_{i,i}^{t+1} - \boldsymbol{W}_{i,i-1}^{t+1} \rangle + \frac{L_{i}}{2} \|\boldsymbol{W}_{i,i}^{t+1} - \boldsymbol{W}_{i,i-1}^{t+1}\|_{2}^{2} \leq -\frac{\beta K}{2} \|\nabla H(\boldsymbol{W}_{i-1}^{t})\|_{2}^{2} - \frac{\alpha K}{8G} \|\nabla_{i} H(\boldsymbol{W}_{i-1}^{t})\|_{2}^{2} \leq -\frac{\alpha K}{8G} \|\nabla H(\boldsymbol{W}_{i-1}^{t})\|_{2}^{2}$$
(16)

910 where the second inequality follows (Luo et al., 2024)'s analysis. Specifically, by following exact 911 the same argument in (Luo et al., 2024) Lemma D.6 and Lemma D.7, the last two terms can be 912 decomposed into a descent term $-\mathcal{O}(\alpha \|\nabla H(\mathbf{W}_i^t)\|)$ plus an error term $\mathcal{O}(\alpha \|\tilde{\mathbf{e}}_i^t\|^2)$ that represents 913 the gradient bias, where the error $\|\tilde{\mathbf{e}}_i^t\|$ can be shown to be in the order of $\mathcal{O}(\alpha \|\nabla_i H(\mathbf{W}_i^t)\|)$, which 914 can be eliminated by controlling the step size α .

Proof of Theorem 2. Since (16) establishs exact the same descent property as in Luo et al. (2024)
 Corollary D.8, one can follow the same argument as in BAdam's analysis to prove the Theorem 2.

918 C.2 PROOF OF PROPOSITIONS 919

920 Proof of Proposition 1. We first show that $H^* = 0$. One can construct $z_3 = [1, 0, \dots, 0]^{\top}$ and 921 $W_3 = [y, 0, \dots, 0]$. Note that such a choice of z_3 is always achievable by choosing a specific W_2 . 922 Hence, 0 function value can be attained by the constructed feasible point. This yields $H^* = 0$ after 923 realizing that the objective function must be nonnegative. We illustrate this issue in Figure 1.

924 In BCD, W_1 and W_3 are fixed. We further assume that the fixed W_3 has full column rank. We split 925 our discussion into two cases. Case I: $y \notin \operatorname{range}(W_3)$. We trivially have $\tilde{H}^* > 0 = H^*$. Case 926 II: $y \in \operatorname{range}(W_3)$. In this case, $z_3^* := (W_3^\top W_3)^{-1} W_3^\top y$ is the unique point that can achieve 927 0 function value. However, since z_3^* has at least one negative entry and $z_3 \ge 0$ (due to the ReLU 928 activation), we have $||z_3 - z_3^*||_2^2 > 0$. Therefore, we have $||y - \hat{y}||_2^2 = ||W_3(z_3^* - z_3)||_2^2 > 0 = H^*$, 929 where the last inequality follows from the full column rankness of W_3 .

931 Proof of Proposition 2. We construct $z_3 = e_1 = [1, 0, \dots, 0]^{\top}$. Let $C = [y - W_3^{(1)}, 0, \dots, 0]$, 932 where $W_3^{(1)}$ is the first column of W_3 . Then, we have $||(W_3 + C)z_3 - y||_2^2 = ||Ce_1 - (y - W_3e_1)||_2^2 = 0$.

C.3 ANALYSIS OF MULTI-LAYER MODEL AND CROSS ENTROPY LOSS

 $\boldsymbol{z}_1 = \sigma(\boldsymbol{W}_1 \boldsymbol{x})$

In this section, we generalize the Proposition 1 to *L*-layer neural network model and cross entropy loss. The corresponding numerical verification are presented in Appendix D.4. Let us consider an *L*-layer model:

 $\boldsymbol{z}_i = \sigma(\boldsymbol{W}_i \boldsymbol{z}_{i-1}), \quad i = 2, \cdots, L-1$

935

936 937

938

939

942 943 944

945

946 947

952 953 954

955

956

957 958

959

966 967 $\hat{\boldsymbol{y}} = \boldsymbol{W}_L \boldsymbol{z}_{L-1},$

where $\sigma(\boldsymbol{x}) = \max(0, \boldsymbol{x})$ is the ReLU activation function and $\boldsymbol{z}_i \in \mathbb{R}^{d_i}$.

C.3.1 SUBOPTIMALITY ANALYSIS FOR L-LAYER MODEL

Let us first consider the general regression loss $\|\hat{y} - y\|_2^2$, where y is the target we aim to fit.

Effect of freezing W_L . When W_L is full column rank, the optimal z_{L-1} we seek to fit is the least square solution $z_{L-1}^* = (W_L^\top W_L)^{-1} W_L^\top y$. When z_{L-1}^* contains negative entries, it cannot be fit due to the non-negativity of the ReLU function, which induces the suboptimality:

$$\min_{oldsymbol{W}_1,\cdots,oldsymbol{W}_{L-1}} \|\hat{oldsymbol{y}} - oldsymbol{y}\|_2^2 > \min_{oldsymbol{W}_1,\cdots,oldsymbol{W}_L} \|\hat{oldsymbol{y}} - oldsymbol{y}\|_2^2$$

Effect of freezing intermediate layers. Each intermediate layer performs the transformation $z_i = \sigma(W_i z_{i-1}) := \mathcal{M}_i$. When W_i is trainable, we have range $(\mathcal{M}_i) = \mathbb{R}^{d_i+}$ when $z_{i-1} \neq 0$. However, when W_i is frozen, range (\mathcal{M}_i) is limited to the projected "restricted" column space of W_i , where "restricted" means that the column combination should be positive, due to the positivity of z_{i-1} .

C.3.2 SUBOPTIMALITY ANALYSIS FOR CROSS ENTROPY LOSS

960 Without loss of generality, let us assume that the ground truth label is the first class. The cross 961 entropy loss is given by $-\log(\exp \hat{y}_1/(\sum_{i=1}^m \exp \hat{y}_i)))$, where *m* is the number of classes. Consider 962 the case where the weight of the last layer W_L has a row with the same weight as the first row, i.e. 963 $\exists j$ such that $W_L^{(j)} = W_L^{(1)}$, we have $\hat{y}_j = W_L^{(j)} \boldsymbol{z}_{L-1} = W_L^{(1)} \boldsymbol{z}_{L-1} = \hat{y}_1$. In this case, we will 964 never be able to drive the loss down to $-\log \frac{1}{2}$:

$$-\log\left(\frac{\exp\hat{y}_1}{\left(\sum_{i=1}^m\exp\hat{y}_i\right)}\right) > -\log\left(\frac{\exp\hat{y}_1}{\exp\hat{y}_1 + \exp\hat{y}_i}\right) = -\log\frac{1}{2}$$

While it is not common for W_L to have exactly two same rows, one can expect large error when there are rows that form small angle, i.e.

970
971
$$\frac{\left\langle \boldsymbol{w}_{L}^{(i)}, \boldsymbol{w}_{L}^{(j)} \right\rangle}{\|\boldsymbol{w}_{L}^{(i)}\| \|\boldsymbol{w}_{L}^{(j)}\|} \approx 1$$

D ADDITIONAL EXPERIMENTS

972

973 974

975 976

977 978

979

980

981

982 983

984

985

986

987

988

989

990 991

1008

1009 1010

1011

1012

1013

1014

1015 1016

1017

1018

1020

1021

D.1 MORE ABLATION STUDY RESULTS

We present additional ablation studies on the hyperparameters of BREAD, including the block switch frequency K, the rank of correction matrices r, and the block ordering strategies.



Figure 4: Ablation study on the effect of Adam inner steps K, rank of correction matrices r, and block ordering strategies.

992Effect of K. We present the effect of sub-problem update steps K in Fig. 4a, which is by default993128 in our paper's experiments. The convergence of BAdam is provided for reference. Evidently,994increasing K consistently accelerates the convergence of BREAD for the examined range, where995BREAD with K = 512 takes only half of the iterations to reach the final training loss of BREAD996with K = 32. One possible explanation for the phenomenon is that when using larger K, the Adam997thereby finds better search direction and scaling. We leave the scientific study of K as a future998direction. Notably, BREAD outperforms BAdam under all choices of K.

Effect of r. The effect of correction matrices' rank r is shown in Fig. 4b, which is set to 8 in our paper. We note that by adding rank-2 correction matrices, BREAD converges significantly faster than BAdam, which corroborates our observation in Proposition 2. BREAD exhibits faster convergence as the rank increases, since larger rank offers higher freedom of search directions.

Effect of ordering strategies. We test the ordering strategies of ascending (from input layer to output layer), descending (from output layer to input layer), and random (select the layer in random reshuffling manner). As shown in Fig. 4c, different ordering strategy does not result in evident difference of convergence speed.

Convergence versus iteration Convergence versus time 0.65 0.65 BAdam BAdam BREAD 0.60 BREAD 0.60 BREAD-partial **BREAD**-partial 0.55 BREAD-SGD 0.55 BREAD-SGD BREAD-SGD-partial BREAD-SGD-partia sol 0.50 ^{SO} 0.50 fuining 0.45 Dujuju 0.45 0.40 0.40 0.35 0.35 0.30 0.30 60 ò 10000 20000 50000 ò 10 20 30 40 50 30000 40000 Time (s) Number of iterations (a) Training loss in terms of time. (b) Training loss in terms of iteration.

D.2 CONVERGENCE IN TERMS OF TIME.



1026 D.3 MAGNITUDE OF THE LEARNED PERTURBATION

1033 1034 1035

1036

1039

1040

1041

1042

1043

1044 1045

1046

1047

1053

1028 1029 Under the task of finetuning Llama 3.1-8B on MathInstruct, we plot the norm of the learned per-1030 turbation by BAdam and BREAD, where the magnitude is defined as the difference between the 1030 finetuned model and the pre-trained model: $\Delta W := W_T - W_0$. The results is shown in Fig. 6. 1031 Under the same step size choice, BREAD learns perturbation that has higher norm. This partially 1032 verifies that BREAD can expand the search space more efficiently than BAdam.



Figure 6: Norm of the learned perturbation. BREAD learns larger update than BAdam.

D.4 FEATURE CORRECTION STUDY ON L-LAYER MODEL AND CROSS ENTROPY LOSS

1054 In this section, we present the training loss of classification problem, where we use the cross entropy 1055 as the training objective. Specifically, we treat each layer as one block, and set the block switch 1056 frequency K = 100. We begin training with the first layer.

As shown in Fig. 7a, BREAD with rank-1 landscape correction converges dramatically faster than BAdam. In particular, BAdam only begins to converge rapidly after the 300 steps, when the final layer has been trained. This phenomenon supports our discussion in Appendix C.3.2, where we noted that a poorly trained final layer may hinder convergence. In Fig. 7b, we show that BREAD boosts the convergence for 8-layer neural network as well, which corroborates our *L*-layer analysis in Appendix C.3.1.





Base model: Llama 3.1-8B													
Method	GSI	GSM8K		MATH		NumGLUE		SimulEq		AQuA		Avg.	
	0-shot	4-shot	0-shot	4-shot	0-shot	4-shot	0-shot	4-shot	0-shot	4-shot	0-shot	4-shot	
Base model	17.8	52.5	8.6	23.2	25.7	40.6	12.2	28.8	19.3	43.7	16.7	37.8	
Adam	62.3	64.9	17.4	22.9	56.4	56.8	28.6	33.5	44.9	52.8	41.9	46.2	
Galore	46.7	57.2	16.2	22.9	42.8	45.0	28.7	32.3	47.8	48.4	36.4	41.2	
LOMO	31.9	53.9	15.6	22.7	39.5	39.8	22.6	28.2	36.2	44.9	29.2	37.9	
LoRA	48.7	58.1	13.7	23.0	34.6	54.4	29.6	29.0	47.3	50.3	34.8	43.0	
BAdam	53.9	58.3	17.2	23.6	53.7	57.2	32.5	32.8	50.4	49.6	41.5	44.3	
BREAD	57.0	57.6	20.0	23.7	55.9	58.2	32.5	32.8	49.6	50.0	43.0	44.5	
BREAD-partial	56.1	62.4	18.5	22.3	53.5	53.5	32.2	32.9	48.9	49.6	41.8	44.1	
BREAD-SGD	56.9	60.6	19.6	21.4	54.1	58.2	31.5	31.8	48.0	50.8	42.0	44.6	
BREAD-SGD-part	ial 58.0	60.1	17.8	20.8	53.9	54.3	32.5	31.5	47.3	47.0	41.9	42.7	

Table 3: Math evaluation results for models finetuned on MathInstruct dataset.

0-shot results											
Method	GSM8K	MATH	NumGLUE	SimulEq	AQuA	Avg.					
Base model	4.2	3.9	18.9	3.5	25.2	11.1					
BAdam	10.9 ± 1.27	4.6 ± 0.56	24.4 ± 3.00	5.9 ± 0.56	28.3 ± 1.50	14.8					
BREAD	$\textbf{14.0} \pm 1.19$	$\textbf{5.2}\pm0.31$	28.1 ± 1.31	$\textbf{6.8} \pm 1.35$	28.1 ± 2.46	16.4					
4-shot results											
Method	GSM8K	MATH	NumGLUE	SimulEq	AQuA	Avg.					
Base model	6.2	6.3	20.2	5.6	28.3	13.3					
BAdam	14.8 ± 2.28	5.9 ± 0.39	24.4 ± 3.07	5.9 ± 1.19	29.5 ± 2.05	16.1					
BREAD	$\textbf{17.2} \pm 1.40$	$\textbf{6.0}\pm0.29$	27.3 ± 1.53	$\textbf{5.9} \pm 1.18$	$\textbf{29.9} \pm 3.22$	17.3					

Table 4: Math benchmark scores for Llama 3.2-1B finetuned on MathInstruct dataset. Results are averaged over 3 independent trials.

D.5 ADDITIONAL EVALUATION ON MATH BENCHMARKS

In this section, we provide the math benchmark scores of BREAD variants in Table 3, as well as the averaged scores of finetuned Llama 3.2-1B models in Table 4. All the evaluation protocol follows the released code of MathInstruct.

As shown in Table 3, the variants with partial update have slight performance degrade compared with their full counterparts. Notably, All the BREAD variants outperform BAdam in 0-shot setting, and outperform other memory efficient baselines more significantly. The variants with partial update have slight performance degrade compared with their full counterparts.

Llama 3.2-1B experiments. We finetune Llama 3.2-1B on MathInstruct dataset for 3 independent trials, and report the averaged scores with standard deviation in Table 4. In both the 0-shot and 4-shot settings, BREAD outperforms BAdam across all tasks, supporting the effectiveness of the landscape correction.