

# Memory-Efficient LLM Training with Online Subspace Descent

Kaizhao Liang<sup>†</sup>, Bo Liu<sup>†</sup>, Lizhang Chen<sup>†</sup>, Qiang Liu<sup>†</sup>

<sup>†</sup>The University of Texas at Austin

{kaizhao1, bliu, lzchen, lqiang}@utexas.edu

## Abstract

Recently, a wide range of memory-efficient LLM training algorithms have gained substantial popularity. These methods leverage the low-rank structure of gradients to project optimizer states into a subspace using projection matrix found by singular value decomposition (SVD). However, convergence of these algorithms is highly dependent on the update rules of their projection matrix. In this work, we provide the *first* convergence guarantee for arbitrary update rules of projection matrix. This guarantee is generally applicable to optimizers that can be analyzed with Hamiltonian Descent, including most common ones, such as LION, Adam. Inspired by our theoretical understanding, we propose Online Subspace Descent, a new family of subspace descent optimizer without SVD. Instead of updating the projection matrix with eigenvectors, Online Subspace Descent updates the projection matrix with online PCA. Online Subspace Descent is flexible and introduces only minimum overhead to training. We show that for the task of pretraining LLaMA models ranging from 60M to 7B parameters on the C4 dataset, Online Subspace Descent achieves lower perplexity and better downstream tasks performance than state-of-the-art low-rank training methods across different settings and narrows the gap with full-rank baselines.<sup>1</sup>

## 1 Introduction

The continual advancement in training large language models (LLMs) presents a compelling challenge in balancing computational efficiency with model performance. As the scope and complexity of these models grow, so does the necessity for innovative strategies that optimize memory usage without compromising the learning capabilities of the model. Recent approaches in low-rank adaptation strategies, including Stochastic Subspace Descent [13], LoRA [11], ReLoRA [15], Gradient Low-Rank Projection (GaLore) [25] and Sketchy [9], have paved the way for memory-efficient training by utilizing a periodically updated low-rank projection matrix to manage parameter updates. In particular, GaLore and Sketchy both utilize expensive singular value decomposition to determine the projection matrix, whereas stochastic subspace descent suggests using random matrices as projection matrices and provides convergence analysis on convex functions and objectives. However, to the best of our knowledge, no one has offered any guarantee of convergence for this class of methods on non-convex functions and objectives.

In this work, we provide the first convergence guarantee for arbitrary update rules of the projection matrix. This guarantee is significant because it is broadly applicable to a wide range of optimizers that can be analyzed within the Hamiltonian descent framework [18]. By establishing this convergence guarantee, we demonstrate that our approach is not limited to specific or narrowly defined update rules, but can be extended to include many commonly used optimizers in the field. In particular, this

<sup>1</sup>Code is available at <https://github.com/kyleliang919/Online-Subspace-Descent>.

---

**Algorithm 1** Online Subspace Descent

---

- 1: Required: Optimizer  $\text{OptimizerW}$ , learning rate  $\epsilon_t^W$ , weight decay  $\lambda^W$  for model weights  $\mathbf{W}_t$ ; and  $\{\text{OptimizerP}, \epsilon_t^P, \lambda^P\}$  for the projection matrix  $\mathbf{P}_t$ . Proper initialization.
  - 2: **for** iteration  $t$  **do**
  - 3:   Calculate gradient  $\mathbf{G}_t = \nabla L(\mathbf{W}_t)$ ; Update model weights  $\mathbf{W}_t$  by
$$(\hat{\Delta}_t, \hat{\mathbf{S}}_t) = \text{OptimizerW}(\mathbf{P}_t^\top \mathbf{G}_t, \hat{\mathbf{S}}_{t-1}), \quad \mathbf{W}_{t+1} = \mathbf{W}_t + \epsilon_t^W (\mathbf{P}_t \hat{\Delta}_{t+1} - \lambda^W \mathbf{W}_t)$$
  - 4:   Calculate  $\mathbf{G}_t^P = \nabla L_{\mathbf{G}_t}(\mathbf{P}_t)$  for  $L_{\mathbf{G}_t}(\cdot)$  in Eq (6); Update the projection  $\mathbf{P}_t$  by
$$(\Delta_t^P, \mathbf{S}_t^P) = \text{OptimizerP}(\mathbf{G}_t^P, \mathbf{S}_{t-1}^P), \quad \mathbf{P}_{t+1} = \mathbf{P}_t + \epsilon_t^P (\Delta_t^P - \lambda^P \mathbf{P}_t)$$
  - 5: **end for**
  - 6: Remark: We added weight decay as a common heuristic. We recommend using Adam for both optimizers, and set  $\epsilon_t^P = \alpha \epsilon_t^W$  with a constant  $\alpha$  (e.g.,  $\alpha = 5$ ), and  $\lambda^W = \lambda^P$ .<sup>2</sup>
- 

includes popular algorithms such as LION [4] and Adam [12], which are widely used in various machine learning and optimization tasks. Our results therefore offer a robust theoretical foundation for understanding and analyzing the behavior of these optimizers, ensuring their effectiveness and reliability in diverse applications.

Inspired by our theoretical understanding, we introduce a novel family of memory-efficient optimizers named Online Subspace Descent, which incorporates a dynamically changing projection matrix, replacing the conventional periodic updating approach (SVD) with online PCA. By allowing the projection matrix to evolve in response to the changing gradient landscape, Online Subspace Descent enhances the model’s ability to navigate the parameter space more effectively. This dynamic adaptation aligns more closely with the natural progression of learning in deep neural networks, which is characterized by changes in the importance of different characteristics and interactions as training progresses. Through extensive experiments and comparative analysis, we demonstrate that our approach presents lower perplexity in pretraining LLaMA models (ranging from 60M to 1B parameters) on the C4 dataset compared to state-of-the-art low-rank training methods, closing the perplexity gap with full-rank baselines on language model pretraining.

## 2 Optimization Background

The training of deep learning models reduces to an optimization problem

$$\min_{\mathbf{W}} L(\mathbf{W}),$$

where  $\mathbf{W}$  is the set of weight matrices of the model. For simplicity, we assume  $\mathbf{W} \in \mathbb{R}^{n \times m}$  is a single matrix of size  $(n, m)$  without loss of generality. For notation, we write  $\langle \mathbf{A}, \mathbf{B} \rangle = \text{tr}(\mathbf{A}^\top \mathbf{B})$  for inner products of matrices, and  $\|\mathbf{A}\|^2 = \text{tr}(\mathbf{A}^\top \mathbf{A})$  the Frobenius norm. We use  $\mathbf{A} \odot \mathbf{B}$  to denote the elementwise product, and  $\mathbf{A}^{\odot 2} = \mathbf{A} \odot \mathbf{A}$ .

**Example 2.1.** Update rules of common optimizers:

*Gradient Descent* :  $\mathbf{W}_{t+1} = \mathbf{W}_t - \epsilon_t \nabla L(\mathbf{W}_t)$ ,

*Momentum* :  $\mathbf{W}_{t+1} = \mathbf{W}_t - \epsilon_t \mathbf{M}_t$ ,  $\mathbf{M}_t = (1 - \beta) \nabla L(\mathbf{W}_t) + \beta \mathbf{M}_{t-1}$ ,

*Lion-K* [4] :  $\mathbf{W}_{t+1} = \mathbf{W}_t - \epsilon_t \nabla \mathcal{K}(\mathbf{N}_t)$ ,  $\mathbf{N}_t = (1 - \beta_1) \nabla L(\mathbf{W}_t) + \beta_1 \mathbf{M}_t$   
 $\mathbf{M}_t = (1 - \beta_2) \nabla L(\mathbf{W}_t) + \beta_2 \mathbf{M}_{t-1}$ ,

*Adam* [12] :  $\mathbf{W}_{t+1} = \mathbf{W}_t - \epsilon_t \frac{\mathbf{M}_t}{\sqrt{\mathbf{V}_t} + e}$ ,  $\mathbf{M}_t = (1 - \beta_{1t}) \nabla L(\mathbf{W}_t) + \beta_{1t} \mathbf{M}_{t-1}$ ,

$\mathbf{V}_t = (1 - \beta_{2t}) \nabla L(\mathbf{W}_t)^{\odot 2} + \beta_{2t} \mathbf{V}_{t-1}$ ,

where  $\epsilon_t$  are step sizes, and  $\mathbf{M}_t, \mathbf{V}_t$  are the first and second order momentum, and  $\beta, \beta_1, \beta_2$  are momentum coefficients in  $(0, 1)$ , with  $\beta_{it} = \frac{\beta_i - \beta_i^{t+1}}{1 - \beta_i^{t+1}}$ ,  $i = 1, 2$  for  $\beta_1, \beta_2 \in (0, 1)$  in Adam, and  $\mathcal{K}$  is any convex function with  $\nabla \mathcal{K}(\mathbf{0}) = \mathbf{0}$  for Lion-K [4], and Lion [5] uses  $\mathcal{K}(\mathbf{X}) = \|\mathbf{X}\|_{1,1}$  and  $\nabla \mathcal{K}(\mathbf{X}) = \text{sign}(\mathbf{X})$ .

These optimizers can be unifiedly viewed as updating  $\mathbf{W}_t$  together with an optimizer state  $\mathbf{S}_t$ :

$$\mathbf{W}_{t+1} = \mathbf{W}_t + \phi_t(\mathbf{S}_t), \quad \mathbf{S}_t = \psi_t(\mathbf{S}_{t-1}, \nabla L(\mathbf{W}_t)), \quad (1)$$

with some mapping  $\phi_t, \psi_t$ . We have  $\mathbf{S}_t = \mathbf{M}_t$  for momentum and  $\mathbf{S}_t = \{\mathbf{M}_t, \mathbf{V}_t\}$  for Adam. Note that both  $\mathbf{M}_t, \mathbf{V}_t$  are of the same size as the model weights  $\mathbf{W}_t$ , resulting in high memory consumption for large models. This issue is particularly pronounced for Adam, which typically yields the best performance for large language models (LLMs) but incurs the highest memory cost due to the need to maintain both  $\mathbf{M}_t$  and  $\mathbf{V}_t$ . One key challenge is to retain the high performance of Adam while enhancing its memory efficiency.

**Hamiltonian+Descent** One powerful approach to studying the dynamic properties of optimizers is to examine their continuous-time ODE forms in the limit of infinitesimal step size. The continuous-time forms provide clearer insights into the asymptotic convergence of the algorithm, abstracting away the choices of step size, discretization, and stochastic errors. The underlying logic is that a "sound" optimizer should be guaranteed to converge to local optima of the loss, at least when using sufficiently small step sizes.

Inspired by [4, 18], we observe that the continuous-time form of many common optimizers yields a *Hamiltonian+Descent* structure,

$$\begin{aligned} \frac{d}{dt} \mathbf{W}_t &= \partial_{\mathbf{S}} H(\mathbf{W}_t, \mathbf{S}_t) - \Phi(\partial_{\mathbf{W}} H(\mathbf{W}_t, \mathbf{S}_t)) \\ \frac{d}{dt} \mathbf{S}_t &= -\partial_{\mathbf{W}} H(\mathbf{W}_t, \mathbf{S}_t) - \Psi(\partial_{\mathbf{S}} H(\mathbf{W}_t, \mathbf{S}_t)), \end{aligned} \quad (2)$$

where  $H(\mathbf{W}, \mathbf{S})$  is a Hamiltonian (or Lyapunov) function that satisfies

$$\min_{\mathbf{S}} H(\mathbf{W}, \mathbf{S}) = L(\mathbf{W}), \quad \forall \mathbf{W},$$

so that minimizing  $L(\mathbf{W})$  reduces to minimizing  $H(\mathbf{W}, \mathbf{S})$ ; and  $\Phi(\cdot), \Psi(\cdot)$  are two monotonic mappings satisfying

$$\|\mathbf{X}\|_{\Phi}^2 := \langle \mathbf{X}, \Phi(\mathbf{X}) \rangle \geq 0, \quad \|\mathbf{X}\|_{\Psi}^2 := \langle \mathbf{X}, \Psi(\mathbf{X}) \rangle \geq 0, \quad \forall \mathbf{X}.$$

With  $\Phi(\mathbf{X}) = \Psi(\mathbf{X}) = 0$ , the system in (2) reduces to the standard Hamiltonian system that keeps  $H(\mathbf{W}_t, \mathbf{S}_t) = \text{const}$  along the trajectory. When adding the descending components with  $\Phi$  and  $\Psi$ , the system then keeps  $H(\mathbf{W}, \mathbf{S})$  monotonically non-increasing:

$$\begin{aligned} \frac{d}{dt} H(\mathbf{W}_t, \mathbf{S}_t) &= \left\langle \partial_{\mathbf{W}} H_t, \frac{d}{dt} \mathbf{W}_t \right\rangle + \left\langle \partial_{\mathbf{S}} H_t, \frac{d}{dt} \mathbf{S}_t \right\rangle \\ &= \langle \partial_{\mathbf{W}} H_t, \partial_{\mathbf{S}} H_t - \Phi(\partial_{\mathbf{W}} H_t) \rangle + \langle \partial_{\mathbf{S}} H_t, -\partial_{\mathbf{W}} H_t - \Psi(\partial_{\mathbf{S}} H_t) \rangle \\ &= -\|\partial_{\mathbf{W}} H_t\|_{\Phi}^2 - \|\partial_{\mathbf{S}} H_t\|_{\Psi}^2 \leq 0, \end{aligned} \quad (3)$$

where we write  $\partial_{\mathbf{W}} H_t = \partial_{\mathbf{W}} H(\mathbf{W}_t, \mathbf{S}_t)$  and similarly for  $\partial_{\mathbf{S}} H_t$ . The main idea is that the cross terms  $\langle \partial_{\mathbf{W}} H_t, \partial_{\mathbf{S}} H_t \rangle$  are canceled, leaving only the negative terms.

**Example 2.2.** The momentum method yields following continuous-time form and Hamiltonian:

$$\frac{d}{dt} \mathbf{W}_t = -\mathbf{M}_t, \quad \frac{d}{dt} \mathbf{M}_t = a(\nabla L(\mathbf{W}_t) - \mathbf{M}_t), \quad \text{with} \quad H(\mathbf{W}, \mathbf{M}) = L(\mathbf{W}) + \frac{\|\mathbf{M}\|^2}{2a}.$$

**Example 2.3.** Adam [12] yields the following continuous-time form and Hamiltonian,

$$\begin{aligned} \frac{d}{dt} \mathbf{W}_t &= -\frac{\mathbf{M}_t}{\sqrt{\mathbf{V}_t + e}}, \quad \frac{d}{dt} \mathbf{M}_t = a(\nabla L(\mathbf{W}_t) - \mathbf{M}_t), \quad \frac{d}{dt} \mathbf{V}_t = b(\nabla L(\mathbf{W}_t)^{\odot 2} - \mathbf{V}_t), \\ \text{with} \quad H(\mathbf{W}, \mathbf{M}, \mathbf{V}) &= L(\mathbf{W}) + \frac{1}{2a} \left\langle \frac{\mathbf{M}}{\sqrt{\mathbf{V} + e}}, \mathbf{M} \right\rangle, \end{aligned}$$

for which we can show that  $\frac{d}{dt} H(\mathbf{W}_t, \mathbf{M}_t, \mathbf{V}_t) \leq 0$  when  $a \geq b/4$ .

**Example 2.4.** The Lion-K optimizer [5, 4] (without weight decay) can be written into

$$\frac{d}{dt} \mathbf{W}_t = \nabla \mathcal{K}((1-b)\mathbf{M}_t - b\nabla L(\mathbf{W}_t)), \quad \frac{d}{dt} \mathbf{M}_t = -a(\nabla L(\mathbf{W}_t) + \mathbf{M}_t),$$

where  $a \geq 0, b \in [0, 1]$  and  $\mathcal{K}(\mathbf{X})$  is any convex function that attains the minimum at  $\mathbf{X} = 0$ . One of its Hamiltonians that yields the Hamiltonian+descent structure (Eq (13) in Chen et al. [4]) is

$$H(\mathbf{W}, \mathbf{M}) = aL(\mathbf{W}) + \frac{1}{1-b} \mathcal{K}((1-b)\mathbf{M}).$$

### 3 Memory-Efficient Optimizers via Online Subspace Descent

We introduce the idea of constructing memory efficient optimizers by descending in the subspaces that dynamically changes across iterations as motivated by GaLore [25] and Sketchy [9]. We first derive *static* subspace descent by restricting the whole optimization on a subspace (Section 3.1), and then propose to dynamically change the subspace across iterations as a heuristic to attain the optimization in the full space while only using subspace descent (Section 3.2). In particular, we propose to update the subspaces via continuous online PCA like updates to avoid the need of exact SVD like in GaLore and Sketchy (Section 3.2). Finally, we remark in Section 3.3 the heuristic nature of the derivation of the method and highlight the difficulty in theoretical understanding, which motivates our analysis based on Hamiltonian dynamics in Section 4.

#### 3.1 Static Subspace Descent

One popular approach to improving memory efficiency is to confine the optimization to a low-dimensional space. To do this, we impose a low rank structure of  $\mathbf{W} = \mathbf{P}\hat{\mathbf{W}}$ , where  $\mathbf{P} \in \mathbb{R}^{n \times k}$  is a projection matrix to be determined later and  $\hat{\mathbf{W}} \in \mathbb{R}^{k \times m}$  is a dimension-reduced parameter. When  $k \ll n$ ,  $\mathbf{P}$  and  $\hat{\mathbf{W}}$  are much smaller in size compared to  $\mathbf{W}$ . Now consider

$$\min_{\hat{\mathbf{W}}} L(\mathbf{P}\hat{\mathbf{W}}).$$

Applying the optimizer from (1) to update  $\hat{\mathbf{W}}$  along with an optimizer state  $\hat{\mathbf{S}}$ , and mapping the update rule  $\hat{\mathbf{W}}_{t+1} = \hat{\mathbf{W}}_t + \phi_t(\hat{\mathbf{S}}_t)$  to that of  $\mathbf{W} = \mathbf{P}\hat{\mathbf{W}}_t$ , we get

$$\mathbf{W}_{t+1} = \mathbf{W}_t + \mathbf{P}\phi_t(\hat{\mathbf{S}}_t), \quad \hat{\mathbf{S}}_t = \psi_t(\hat{\mathbf{S}}_{t-1}, \mathbf{P}^\top \nabla L(\mathbf{W}_t)), \quad (4)$$

where we used the fact that  $\nabla_{\hat{\mathbf{W}}} L(\mathbf{P}\hat{\mathbf{W}}) = \mathbf{P}^\top \nabla_{\mathbf{W}} L(\mathbf{W})$ . This yields a more memory-efficient optimizer, as the size of  $\hat{\mathbf{S}}_t$  is proportional to that of  $\hat{\mathbf{W}}_t$ , much smaller than  $\mathbf{S}_t$  in (1) when  $k \ll n$ .

#### 3.2 Online Subspace Descent

With a static  $\mathbf{P}$ , regardless of its values, the parameter  $\mathbf{W}$  is restricted to have a low rank structure. Although low rank assumption is proved to be useful for fine-tuning with LoRA-like methods [11], it is often too limited for pre-training or when the desirable model weights are not inherently low-rank.

To address this problem, Zhao et al. [25] suggested to keep the projected updated in (4), but use different  $\mathbf{P}$  across the iterations:

$$\mathbf{W}_{t+1} = \mathbf{W}_t + \mathbf{P}_t \phi_t(\hat{\mathbf{S}}_t), \quad \hat{\mathbf{S}}_t = \psi_t(\hat{\mathbf{S}}_{t-1}, \mathbf{P}_t^\top \nabla L(\mathbf{W}_t)), \quad \mathbf{P}_{t+1} = \chi_t(\mathbf{P}_t, \mathbf{W}_t, \hat{\mathbf{S}}_t),$$

where  $\chi_t$  is a update rule of  $\mathbf{P}_t$  that will be determined in the sequel. The intuition is to open up different projection directions at different iterations, so that optimization can be conducted in different subspaces across different iterations. This is similar to the update of coordinate descent, except in a continuous fashion. Note that the update of  $\mathbf{P}_t$  can be done in parallel with that of  $(\mathbf{W}_t, \hat{\mathbf{S}}_t)$ , and incurs no slowdown once it is fast enough to not cause a speed bottleneck.

**Example 3.1.** *Examples of common optimizers equipped with online subspace descent:*

$$\begin{aligned} \text{Gradient Descent : } & \mathbf{W}_{t+1} = \mathbf{W}_t - \epsilon_t \mathbf{P}_t \mathbf{P}_t^\top \mathbf{G}_t, & \mathbf{G}_t &= \nabla L(\mathbf{W}_t), \\ \text{Momentum : } & \mathbf{W}_{t+1} = \mathbf{W}_t - \epsilon_t \mathbf{P}_t \hat{\mathbf{M}}_t, & \hat{\mathbf{M}}_t &= (1 - \beta) \mathbf{P}_t^\top \mathbf{G}_t + \beta \hat{\mathbf{M}}_{t-1}, \\ \text{Lion-K : } & \mathbf{W}_{t+1} = \mathbf{W}_t - \epsilon_t \mathbf{P}_t \nabla \mathcal{K}(\hat{\mathbf{N}}_t), & \hat{\mathbf{G}}_t &= \mathbf{P}_t^\top \mathbf{G}_t \\ & \hat{\mathbf{N}}_t = (1 - \beta_1) \hat{\mathbf{G}}_t + \beta_1 \hat{\mathbf{M}}_t, & \hat{\mathbf{M}}_t &= (1 - \beta_2) \hat{\mathbf{G}}_t + \beta_2 \hat{\mathbf{M}}_{t-1}, \\ \text{Adam : } & \mathbf{W}_{t+1} = \mathbf{W}_t - \epsilon_t \mathbf{P}_t \frac{\hat{\mathbf{M}}_t}{\sqrt{\hat{\mathbf{V}}_t + e}}, & \hat{\mathbf{G}}_t &= \mathbf{P}_t^\top \mathbf{G}_t, \\ & \hat{\mathbf{M}}_t = (1 - \beta_{1t}) \hat{\mathbf{G}}_t + \beta_{1t} \hat{\mathbf{M}}_{t-1}, & \hat{\mathbf{V}}_t &= (1 - \beta_{2t}) \hat{\mathbf{G}}_t^{\odot 2} + \beta_{2t} \hat{\mathbf{V}}_{t-1}. \end{aligned}$$

How Should  $\mathbf{P}_t$  be Updated? It is useful to draw intuition from the projected gradient descent rule

$$\mathbf{W}_{t+1} = \mathbf{W}_t - \epsilon_t \mathbf{P}_t \mathbf{P}_t^\top \mathbf{G}_t, \quad \mathbf{G}_t = \nabla L(\mathbf{W}_t), \quad (5)$$

in which  $\mathbf{P}_t \mathbf{P}_t^\top$  can be viewed as a low rank preconditioning of  $\mathbf{G}_t$ . To make it follow the exact gradient descent, we hope to make  $\mathbf{P}_t \mathbf{P}_t^\top \mathbf{G}_t$  approximate  $\mathbf{G}_t$  as much as possible. In Galore, this is achieved by performing singular value decomposition (SVD) on  $\mathbf{G}_t$  periodically every  $T$  iterations:

$$\mathbf{P}_t, \_, \_ = \text{torch.linalg.svd}(\mathbf{G}_{T\lfloor t/T \rfloor}),$$

where  $T\lfloor t/T \rfloor$  is the largest multiple of  $T$  less than or equal to  $t$ . However, numerical SVD incurs a large computational cost for very large models. Also, since  $\mathbf{P}_t$  is fully determined by  $\mathbf{G}_{T\lfloor t/T \rfloor}$  calculated from a single mini-batch at the last periodic point, it does not incorporate the gradient information from all data in a timely fashion.

In this work, we propose to update  $\mathbf{P}_t$  in a continuous online fashion that incorporates the most recent gradient information in a timely fashion, without calling `torch.linalg.decompositions` routines. We view the update of  $\mathbf{P}_t$  as conducting an online principal component analysis (PCA) based on the streaming of  $\{\mathbf{G}_t\}$ . In particular, we propose to update  $\mathbf{P}_t$  at time  $t$  by minimizing the following PCA objective:

$$L_{G_t}(\mathbf{P}) = \left\| \mathbf{P} \mathbf{P}^\top \tilde{\mathbf{G}}_t - \tilde{\mathbf{G}}_t \right\|^2 + \lambda \left\| \mathbf{P}^\top \mathbf{P} - \mathbf{I}_{k \times k} \right\|^2, \quad \tilde{\mathbf{G}}_t = \frac{\mathbf{G}_t}{\|\mathbf{G}_t\|}, \quad (6)$$

where  $\|\mathbf{A}\| = \text{tr}(\mathbf{A}^\top \mathbf{A})^{1/2}$  and  $\mathbf{I}_{k \times k}$  is the  $k \times k$  identity matrix; we introduced an auxiliary loss to encourage the columns of  $\mathbf{P}$  to be orthonormal and normalizes  $\mathbf{G}_t$  to increase stability.

The key property of  $L_{G_t}(\mathbf{P})$  in (6) is that all its stable local minimum is a global minimum, and  $\mathbf{P}$  is a global minimum iff  $\mathbf{P} \mathbf{P}^\top \tilde{\mathbf{G}}_t$  forms the optimal rank- $k$  approximation of  $\tilde{\mathbf{G}}_t$  [e.g., 3]; moreover, we have  $\mathbf{P}^\top \mathbf{P} = \mathbf{I}_{k \times k}$  at optima when  $\lambda > 0$ .

Instead of minimizing  $L_{G_t}(\mathbf{P})$  exactly, to retain computational efficiency, we propose to update  $\mathbf{P}_t$  by only performing one step of optimization on  $L_{G_t}(\mathbf{P})$ :

$$\mathbf{P}_{t+1} = \text{OptimizerP.step}(\mathbf{P}_t, \nabla_{\mathbf{P}} L_{G_t}(\mathbf{P}_t)),$$

where `OptimizerP.step` can be a favorite optimizer, such as gradient descent or Adam. Note that when using Adam, we introduce a copy of optimizer state  $\mathbf{S}_t^P$  for  $\mathbf{P}_t$ . See Algorithm 1. Compared to the exact SVD, each online update of  $\mathbf{P}_t$  here is fast and can be executed in parallel with the  $(\mathbf{W}_t, \hat{\mathbf{S}}_t)$  updates to avoid slowdown.

### 3.3 Difficulty in Theoretical Understanding

The idea above of projecting an arbitrary optimizer with a dynamically changing  $\mathbf{P}_t$  is heuristically motivated and lacks an immediate rigorous theoretical justification. The main challenge lies in the complex interaction between the update of  $\mathbf{U}_t$  and the optimization state  $\mathbf{S}_t$ , which could potentially degrade the convergence and other theoretical properties of the original optimizer. A key question is whether we can develop a theoretical framework to understand how  $\mathbf{P}_t$  impacts the optimizer's convergence behavior and provide guidance for the design of the update rules of  $\mathbf{P}_t$ .

To gain understanding, it is useful to first exam the simple case of projected gradient descent in (5) which does not have an optimizer state ( $\mathbf{S}_t = \emptyset$ ). In this case, since  $\mathbf{P}_t \mathbf{P}_t^\top$  is positive semi-definite, the update  $\mathbf{P}_t \mathbf{P}_t^\top \mathbf{G}_t$  is always non-increasing direction of  $L(\mathbf{W})$  for any  $\mathbf{P}_t$ . The algorithm is essentially a variant of coordinate or subspace descent, where  $\mathbf{P}_t$  defines the subspace on which one step of gradient descent is conducted at iteration  $t$ . To ensure that (5) finds a local optimum, we mainly need to ensure that  $\mathbf{P}_t \mathbf{G}_t = 0$  only if  $\mathbf{G}_t = 0$  to prevent the optimizer from stopping prematurely; this is a mild condition that can be satisfied e.g. when  $\mathbf{P}_t$  is updated by (online) PCA on  $\mathbf{G}_t$ .

Unfortunately, this coordinate-descent-like interpretation does not apply to more advanced optimizers that track a momentum state  $\mathbf{S}_t$ . This is because  $\mathbf{S}_t$  accumulates the information from the projected gradient  $\mathbf{P}_\tau \mathbf{G}_\tau$  at all earlier iterations  $\tau \leq t$ . As  $\mathbf{P}_t$  changes across time, it is unclear whether the gradient projected to different subspaces  $\mathbf{P}_\tau$  would be coherent with each other, and useful for future updates that are conducted in different subspaces  $\mathbf{P}_t$  for  $t > \tau$ . The difficulty is the inertia effect of  $\mathbf{S}_t$  that entangles the different subspaces, making the dynamic behavior fundamentally more complicated than naive coordinate descent where the descent in different subspaces is uncoupled. This is what we address in Section 4 via the Hamiltonian descent framework.

## 4 Hamiltonian Descent Meets Subspace Descent: A Lyapunov Analysis

In this section, we show a surprising result that the complication outlined above in Section 3.3 *is not* a problem for optimizers that yields the Hamiltonian+descent structure in (2). Our result is two-fold:

- Section 4.1: When applying Online Subspace Descent on systems in (2), the Hamiltonian+descent structure is preserved once the update rule of  $P_t$  has a smooth continuous-time limit. Hence, under very mild conditions, Online Subspace Descent equipped with common optimizers like Adam and Lion automatically yield a Lyapunov function and hence benign continuous-time convergence. Moreover,  $P_t$  can, in fact, be generalized to an arbitrary linear operator as shown in Section 4.3.
- Section 4.2: For any smooth  $P_t$  update rules that eliminates the degenerate case of  $P_t^\top G_t = 0$  while  $G_t = 0$  at convergence, the online subspace optimizer guarantees to converge in continuous time to a stationary point of the loss  $L(W)$ . This mild condition is satisfied, for example, when  $P_t$  is updated by a typical optimizer on the online PCA objective  $L_{G_t}(P)$ .

### 4.1 Online Subspace Descent Preserves the Hamiltonian+Descent Structure

Applying dynamic projection to Hamiltonian descent in (2), we obtain the following systems:

$$\begin{aligned}\frac{d}{dt}W_t &= P_t \partial_{\hat{S}} H(W_t, \hat{S}_t) - \Phi(\partial_W H(W_t, \hat{S}_t)) \\ \frac{d}{dt}\hat{S}_t &= -P_t^\top \partial_W H(W_t, \hat{S}_t) - \Psi(\partial_{\hat{S}} H(W_t, \hat{S}_t)) \\ \frac{d}{dt}P_t &= \Gamma(P_t, \nabla L(W_t)),\end{aligned}\tag{7}$$

where  $\Gamma$  specifies the update rule of  $P_t$ . Following essentially the same derivation as (3), one can show that  $H(W, S)$  remains a Lyapunov function of (7), regardless of the choice of  $\Gamma$ :

$$\begin{aligned}\frac{d}{dt}H(W_t, \hat{S}_t) &= -\|\partial_W H_t\|_\Phi^2 - \|\partial_{\hat{S}} H_t\|_\Psi^2 + \langle \partial_W H_t, P_t \partial_{\hat{S}} H_t \rangle - \langle \partial_{\hat{S}} H_t, P_t^\top \partial_W H_t \rangle \\ &= -\|\partial_W H_t\|_\Phi^2 - \|\partial_{\hat{S}} H_t\|_\Psi^2 \leq 0,\end{aligned}\tag{8}$$

where the key is to use the *adjoint* property of  $P$  and  $P^\top$  that  $\langle P_t X, Y \rangle = \langle X, P_t^\top Y \rangle$ , which cancels the crossing terms, independent of the values of  $P_t$ . There is no requirement on  $\Gamma$  here, besides that the derivative in (8) should exist. As shown in Section 4.3, we can generalize (8) by replacing  $P_t$  and  $P_t^\top$  with a general linear operator  $\mathcal{P}_t$  and its adjoint  $\mathcal{P}_t^*$ .

The following are examples of continuous-time Momentum, Lion- $\mathcal{K}$  and Adam with subspace descent and their Hamiltonian functions.

**Example 4.1.** *Momentum + Online Subspace Descent is*

$$\frac{d}{dt}W_t = -P_t \hat{M}_t, \quad \hat{G}_t = P_t^\top \nabla L(W_t), \quad \frac{d}{dt}\hat{M}_t = a(\hat{G}_t - \hat{M}_t),$$

with Lyapunov function  $H(W, \hat{M}) = L(W) + \frac{\|\hat{M}\|^2}{2a}$ , for which

$$\frac{d}{dt}H(W_t, \hat{M}_t) = -\nabla L(W_t)^\top P_t \hat{M}_t + \hat{M}_t^\top (P_t^\top \nabla L(W_t) - \hat{M}_t) = -\|\hat{M}_t\|_2^2 \leq 0.$$

**Example 4.2.** *Adam + Online Subspace Descent is*

$$\frac{d}{dt}W_t = P_t \frac{\hat{M}_t}{\sqrt{\hat{V}_t} + e}, \quad \hat{G}_t = P_t^\top \nabla L(W_t), \quad \frac{d}{dt}\hat{M}_t = a(\hat{G}_t - \hat{M}_t), \quad \frac{d}{dt}\hat{V}_t = b(\hat{G}_t^2 - \hat{V}_t).$$

with Lyapunov function  $H(W, M, V) = L(W) + \frac{1}{2a} \left\langle \frac{\hat{M}}{\sqrt{\hat{V}} + e}, \hat{M} \right\rangle$ , for which

$$\begin{aligned}
& \frac{d}{dt} H(\mathbf{W}_t, \hat{\mathbf{M}}_t, \hat{\mathbf{V}}_t) \\
&= - \left\langle \mathbf{G}_t, \mathbf{P}_t \frac{\hat{\mathbf{M}}_t}{\sqrt{\hat{\mathbf{V}}_t + e}} \right\rangle + \frac{1}{a} \left\langle \frac{\hat{\mathbf{M}}_t}{\sqrt{\hat{\mathbf{V}}_t + e}}, a(\mathbf{P}_t^\top \mathbf{G}_t - \hat{\mathbf{M}}_t) \right\rangle - \frac{b}{4a} \left\langle \frac{\hat{\mathbf{M}}_t^{\odot 2}}{\sqrt{\hat{\mathbf{V}}_t} \odot (\sqrt{\hat{\mathbf{V}}_t + e})^{\odot 2}}, (\hat{\mathbf{G}}_t^{\odot 2} - \hat{\mathbf{V}}_t) \right\rangle \\
&= - \left\langle 1 - \frac{b}{4a} \frac{\sqrt{\hat{\mathbf{V}}_t}}{\sqrt{\hat{\mathbf{V}}_t + e}}, \frac{\hat{\mathbf{M}}_t^{\odot 2}}{\sqrt{\hat{\mathbf{V}}_t + e}} \right\rangle - \frac{b}{4a} \left\langle \frac{\hat{\mathbf{M}}_t^{\odot 2}}{\sqrt{\hat{\mathbf{V}}_t} \odot (\sqrt{\hat{\mathbf{V}}_t + e})^{\odot 2}}, \hat{\mathbf{G}}_t^{\odot 2} \right\rangle \\
&\leq - \left(1 - \frac{b}{4a}\right) \left\| \frac{\hat{\mathbf{M}}_t}{\sqrt{\sqrt{\hat{\mathbf{V}}_t + e}}} \right\|^2 - \frac{b}{4a} \left\| \frac{\hat{\mathbf{M}}_t \hat{\mathbf{G}}_t}{\sqrt[4]{\hat{\mathbf{V}}_t} (\sqrt{\hat{\mathbf{V}}_t + e})} \right\|^2 \leq 0,
\end{aligned}$$

where we assume  $a \geq b/4$ .

**Example 4.3.** The Lion- $\mathcal{K}$  + Online Subspace Descent is

$$\frac{d}{dt} \mathbf{W}_t = \mathbf{P}_t \nabla \mathcal{K}((1-b)\hat{\mathbf{M}}_t - b\hat{\mathbf{G}}_t), \quad \frac{d}{dt} \mathbf{M}_t = -a(\hat{\mathbf{G}}_t + \hat{\mathbf{M}}_t), \quad \hat{\mathbf{G}}_t = \mathbf{P}_t^\top \nabla L(\mathbf{W}_t)$$

Consider the Hamiltonian function in Eq (13) of [4]:

$$H(\mathbf{W}, \hat{\mathbf{M}}) = aL(\mathbf{W}) + \frac{1}{1-b} \mathcal{K}((1-b)\hat{\mathbf{M}}).$$

$$\begin{aligned}
\frac{d}{dt} H(\mathbf{W}_t, \mathbf{M}_t) &= a \langle \mathbf{G}_t, \mathbf{P}_t \nabla \mathcal{K}((1-b)\hat{\mathbf{M}}_t - b\hat{\mathbf{G}}_t) \rangle - a \langle \nabla \mathcal{K}((1-b)\hat{\mathbf{M}}_t), \hat{\mathbf{G}}_t + \hat{\mathbf{M}}_t \rangle \\
&= a \langle \hat{\mathbf{G}}_t, \nabla \mathcal{K}((1-b)\hat{\mathbf{M}}_t - b\hat{\mathbf{G}}_t) - \nabla \mathcal{K}((1-b)\hat{\mathbf{M}}_t) \rangle - a \langle \nabla \mathcal{K}((1-b)\hat{\mathbf{M}}_t), \hat{\mathbf{M}}_t \rangle \\
&= -\frac{a}{b} [(1-b)\hat{\mathbf{M}}_t; -b\hat{\mathbf{G}}_t]_{\nabla \mathcal{K}} - \frac{a}{(1-b)} [\mathbf{0}; (1-b)\hat{\mathbf{M}}_t]_{\nabla \mathcal{K}} \leq 0
\end{aligned}$$

where we defined  $[\mathbf{X}; \mathbf{Y}]_{\nabla \mathcal{K}} = \langle \mathbf{Y}, \nabla \mathcal{K}(\mathbf{X} + \mathbf{Y}) - \nabla \mathcal{K}(\mathbf{X}) \rangle$  and used the fact that  $[\mathbf{X}; \mathbf{Y}]_{\nabla \mathcal{K}} \geq 0$  by the convexity of  $\mathcal{K}$ ; we used  $\langle \mathbf{G}_t, \mathbf{P}_t \mathbf{X}_t \rangle = \langle \mathbf{P}_t^\top \mathbf{G}_t, \mathbf{X}_t \rangle = \langle \hat{\mathbf{G}}_t, \mathbf{X}_t \rangle$ .

In all examples above, although the form of Hamiltonian  $H(\mathbf{W}, \hat{\mathbf{S}})$  is independent of the update rule of  $\mathbf{P}_t$ , the decreasing rate  $\frac{d}{dt} H(\mathbf{W}_t, \hat{\mathbf{S}}_t)$  depends on  $\mathbf{P}_t$  in a complicate way through  $\hat{\mathbf{M}}_t, \hat{\mathbf{V}}_t, \hat{\mathbf{G}}_t$ . An interesting direction for future investigation is to find optimal rules of  $\mathbf{P}_t$  to maximize the decreasing rate as an optimal control problem.

## 4.2 Convergence to Local Optima

In addition to the Lyapunov structure, we need an additional mild condition on the update rule of  $\mathbf{P}_t$  to ensure the system converges to the local optimum of the loss  $L(\mathbf{W})$ . The main idea is to prevent the system from stopping prematurely before reaching zero gradient  $\mathbf{G}_t = 0$  by excluding the degenerate case of  $\mathbf{P}_t \mathbf{G}_t = 0$  while  $\mathbf{G}_t \neq 0$  in the invariant set of the system.

**Assumption 4.4.** Assume the functions in system (7) are continuously differentiable and

- i)  $\frac{d}{dt} H(\mathbf{W}_t, \hat{\mathbf{S}}_t) = 0$  implies  $\hat{\mathbf{G}}_t = \mathbf{P}_t^\top \nabla L(\mathbf{W}_t) = 0$  and  $\frac{d}{dt} \mathbf{W}_t = 0$ .
- ii) When  $\mathbf{G}_t \equiv \mathbf{G} \neq 0$ , the set  $\{\mathbf{P} : \mathbf{P}^\top \mathbf{G} = 0\}$  is not a positive invariant set of  $\frac{d}{dt} \mathbf{P}_t = \Gamma(\mathbf{P}_t, \mathbf{G}_t)$ .

This is a mild condition. Assumption i) says that the optimizer should stop at a point with  $\hat{\mathbf{G}}_t = 0$ , which is easy to verify for the common optimizers like momentum, Adam, Lion- $\mathcal{K}$ . Assumption ii) ensures  $\hat{\mathbf{G}}_t = 0$  would imply  $\mathbf{G}_t = 0$  in invariance sets, which is satisfied when for example,  $\mathbf{P}_t$  is updated by a reasonable optimizer of the online PCA loss that converges to a stable local minimum.

**Theorem 4.5.** Assume Assumption 4.4 holds. Let  $(\mathbf{W}_t, \mathbf{S}_t, \mathbf{P}_t)_t$  be a bounded solution of (7), then all the accumulation points  $\{\mathbf{W}_t\}$  as  $t \rightarrow +\infty$  are stationary points of  $L(\mathbf{W})$ .

*Proof.* By LaSalle's invariance principle, the positive limit set of  $(\mathbf{W}_t, \mathbf{S}_t, \mathbf{P}_t)_t$  must be contained in  $\mathcal{I}$ , where  $\mathcal{I} = \{\text{the union of complete trajectories satisfying } \frac{d}{dt} H(\mathbf{W}_t, \hat{\mathbf{S}}_t) = 0, \forall t\}$ .



From the Assumption i), the trajectories contained in  $\mathcal{I}$  must satisfy  $\frac{d}{dt}\mathbf{W}_t = 0$ , which implies  $\frac{d}{dt}\mathbf{G}_t = \frac{d}{dt}\nabla L(\mathbf{W}_t) = 0$  and  $\hat{\mathbf{G}}_t = 0$  and hence  $\mathbf{G}_t \equiv \mathbf{G}$  is a constant with  $\mathbf{P}_t^\top \mathbf{G} = 0$ . Moreover, from Assumption ii), we must have  $\nabla L(\mathbf{W}_t) = \mathbf{G}_t \equiv 0$ , since otherwise the trajectory is not invariant. As a result, all trajectories in the limit set  $\mathcal{I}$  must have  $\nabla L(\mathbf{W}_t) = 0$ . Because  $\frac{d}{dt}\mathbf{W}_t = 0$ , these trajectories are static points of  $\mathbf{W}_t$ .  $\square$

### 4.3 Online Subspace Descent with General Linear Projection Operators

We can generalize the online subspace descent with general linear operators:

$$\begin{aligned}\frac{d}{dt}\mathbf{W}_t &= \mathcal{P}_t(\partial_{\hat{\mathbf{S}}}H(\mathbf{W}_t, \hat{\mathbf{S}}_t)) - \Phi(\partial_{\mathbf{W}}H(\mathbf{W}_t, \hat{\mathbf{S}}_t)) \\ \frac{d}{dt}\hat{\mathbf{S}}_t &= -\mathcal{P}_t^*(\partial_{\mathbf{W}}H(\mathbf{W}_t, \hat{\mathbf{S}}_t)) - \Psi(\partial_{\hat{\mathbf{S}}}H(\mathbf{W}_t, \hat{\mathbf{S}}_t)) \\ \frac{d}{dt}\mathcal{P}_t &= \Gamma(\mathcal{P}_t, \nabla L(\mathbf{W}_t)),\end{aligned}$$

where we generalize  $\mathbf{P}_t$  to be any linear operator  $\mathcal{P}_t$  with an adjoint operator  $\mathcal{P}_t^*$ , satisfying

$$\langle \mathbf{X}, \mathcal{P}_t(\mathbf{Y}) \rangle = \langle \mathcal{P}_t^*(\mathbf{X}), \mathbf{Y} \rangle, \quad \forall \mathbf{X}, \mathbf{Y}.$$

The derivation of Lyapunov follows a similar way:

$$\begin{aligned}\frac{d}{dt}H(\mathbf{W}_t, \hat{\mathbf{S}}_t) &= -\|\partial_{\mathbf{W}}H_t\|_{\Phi}^2 - \|\partial_{\hat{\mathbf{S}}}H_t\|_{\Psi}^2 + \langle \partial_{\mathbf{W}}H_t, \mathcal{P}_t(\partial_{\hat{\mathbf{S}}}H_t) \rangle - \langle \partial_{\hat{\mathbf{S}}}H_t, \mathcal{P}_t^*(\partial_{\mathbf{W}}H_t) \rangle \\ &= -\|\partial_{\mathbf{W}}H_t\|_{\Phi}^2 - \|\partial_{\hat{\mathbf{S}}}H_t\|_{\Psi}^2 \leq 0,\end{aligned}$$

where the crossing terms are again canceled due to the adjoint property.

As an example of the general framework, consider  $\mathcal{P}_t(\mathbf{X}) = \mathbf{P}_t\mathbf{X}\mathbf{Q}_t$ , where  $\mathbf{Q}_t$  is another projection matrix applied on the different dimension of  $\mathbf{X}$  (see also [25]). The adjoint operator of  $\mathcal{P}_t$  is  $\mathcal{P}_t^*(\mathbf{X}) = \mathbf{P}_t^\top \mathbf{X} \mathbf{Q}_t^\top$ . This can be verified by

$$\langle \mathbf{P}_t\mathbf{X}\mathbf{Q}_t, \mathbf{Y} \rangle = \text{tr}(\mathbf{P}_t\mathbf{X}\mathbf{Q}_t\mathbf{Y}^\top) = \text{tr}(\mathbf{X}\mathbf{Q}_t\mathbf{Y}^\top\mathbf{P}_t) = \text{tr}(\mathbf{X}(\mathbf{P}_t^\top\mathbf{Y}\mathbf{Q}_t^\top)^\top) = \langle \mathbf{X}, \mathbf{P}_t^\top\mathbf{Y}\mathbf{Q}_t^\top \rangle.$$

The subspace descent system of this operator is

$$\begin{aligned}\frac{d}{dt}\mathbf{W}_t &= \mathbf{P}_t\partial_{\hat{\mathbf{S}}}H(\mathbf{W}_t, \hat{\mathbf{S}}_t)\mathbf{Q}_t - \Phi(\partial_{\mathbf{W}}H(\mathbf{W}_t, \hat{\mathbf{S}}_t)) \\ \frac{d}{dt}\hat{\mathbf{S}}_t &= -\mathbf{P}_t^\top\partial_{\mathbf{W}}H(\mathbf{W}_t, \hat{\mathbf{S}}_t)\mathbf{Q}_t^\top - \Psi(\partial_{\hat{\mathbf{S}}}H(\mathbf{W}_t, \hat{\mathbf{S}}_t)) \\ \frac{d}{dt}\mathbf{P}_t &= \Gamma_P(\mathbf{P}_t, \mathbf{Q}_t, \nabla L(\mathbf{W}_t)) \\ \frac{d}{dt}\mathbf{Q}_t &= \Gamma_Q(\mathbf{P}_t, \mathbf{Q}_t, \nabla L(\mathbf{W}_t)),\end{aligned}$$

where  $\mathbf{P}_t, \mathbf{Q}_t$  can be updated jointly via an online SVD on  $\mathbf{G}_t$ .

Another linear operator that involves two matrices is  $\mathcal{P}_t(\mathbf{X}) = \mathbf{P}_t\mathbf{X} + \mathbf{X}\mathbf{Q}_t$ , which yields  $\mathcal{P}_t^*(\mathbf{X}) = \mathbf{P}_t^\top\mathbf{X} + \mathbf{X}\mathbf{Q}_t^\top$ .

## 5 Experiment

We answer a number of key questions with pretraining experiments of LLaMA [22] on the C4 dataset [20]. All experiments except for large 7B experiments are conducted on a *single* NVIDIA A100 GPU.

### 5.1 Why do we Need Online Subspace Descent?

Overall, Online Subspace Descent offers two major advantages over previous methods that rely on SVD, better convergence and lower overhead. In this section, we discuss both in detail.



First, Online Subspace Descent closes the gap between the state-of-the-art low-rank method and full rank baseline uniformly across different model sizes, as shown in figure 1. A highlight amongst these results is LLaMA 1B (SS 256). As shown in table 1, Online Subspace Descent attains significant improvement over GaLore in perplexity, while consuming a similar amount of GPU memory (8.64 GB v.s 9.01 GB). One additional observation in 1 shows as model size and sequence length grow, Online Subspace Descent becomes more effective. We hypothesize that this is due to the higher intrinsic rank of the underlying optimization problem in larger models. Hence, the positive impact on the convergence of the online update of  $P_t$  becomes more obvious. See more details in Appendix A.

Method	Perplexity(↓)		
	60M	350M	1B
8bit-AdamW (Full Rank)	32.75	30.43	29.40
GaLore (Rank = 512)	57.03	44.34	35.52
<b>Ours (Rank = 512)</b>	<b>56.12</b>	<b>43.67</b>	<b>31.30</b>

Table 1: Pretraining LLaMA 1B with a sequence length of 256 and for 10K steps, perplexity was reported as the training average of the last 10 steps. AdamW8bit serves as the base optimizer for both.

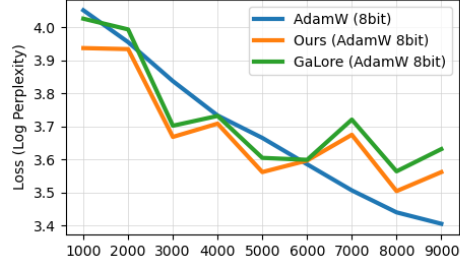


Figure 1: Validation perplexity of LLaMA 1B with sequence length 256, rank 512 for 10K steps.

Another favorable characteristic of Online Subspace Descent is its minimum overhead. In figure 2, we measure and analyze the execution time of SVD and online PCA on a popular data center GPU (A100) and a consumer GPU (RTX 3090). The typical Pytorch implementation of SVD can be up to 142 times slower than running a single-step online PCA on representative weight tensors from LLaMA architectures. Online PCA is fast because it is implemented as a single optimization step with respect to a simple loss function. Hence, each step of online PCA can be cleverly scheduled and hidden in the weight optimization step when executed in parallel, whereas SVD is too expensive to be hidden.

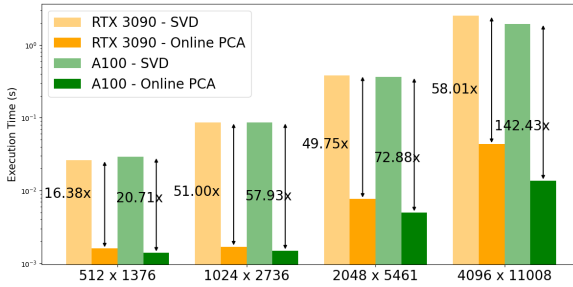


Figure 2: The execution time of `torch.svd` and that a single-step `backward()` call for online PCA in PyTorch, on matrices of typical shapes in linear layers in the LLaMA 60M to 7B. Thanks to the high speed of single-step online PCA,  $P_t$  updates can be executed in parallel with weight updates, adding no overhead to the training process. In contrast, SVD incurs significant overhead as the model and weight tensor sizes increase.

## 5.2 What Rank Should we Pick for *Online Subspace Descent*?

We conduct an ablation study on the rank of Online Subspace Descent. Figure 3 shows that the final perplexity is inversely correlated with rank: higher ranks result in lower convergent perplexity. However, the rate of reduction of perplexity decreases as the rank increases, eventually reaching a saturation point. We propose an intuitive explanation for this phenomenon. In language modeling, high-frequency tokens can be effectively learned with low-rank training. However, learning lower-frequency tokens requires higher ranks. Once these lower-frequency tokens are adequately learned, further increasing the rank does not significantly decrease perplexity. In conclusion, given sufficient time and resources, higher ranks yield better performance for Online Subspace Descent. It is recommended that the highest rank be selected until the perplexity reduction saturates.

## 5.3 What are the Best Hyperparameters?

$\alpha$  and  $\lambda$ : The parameter  $\alpha$  controls the update speed of  $P_t$ , while  $\lambda$  determines the regularization strength on the optimization objective of  $P_t$ . Empirically, we find that the result is not sensitive to  $\lambda$

for small models (60M). and set  $\lambda = 0.1$  for all subsequent experiments. We find that  $\alpha$  must be kept small to avoid instability (Figure 3), and we set  $\alpha = 5$  for all experiments.

**Learning rate:** For the small model (60M), learning rate choices are more flexible, producing similar results. However, for larger models (350M, 1B), we recommend using a learning rate that is 10 times smaller, specifically 0.001. Larger learning rates cause unrecoverable spikes and instability, a general characteristic observed across all methods. See additional hyperparameter choices in Appendix A.

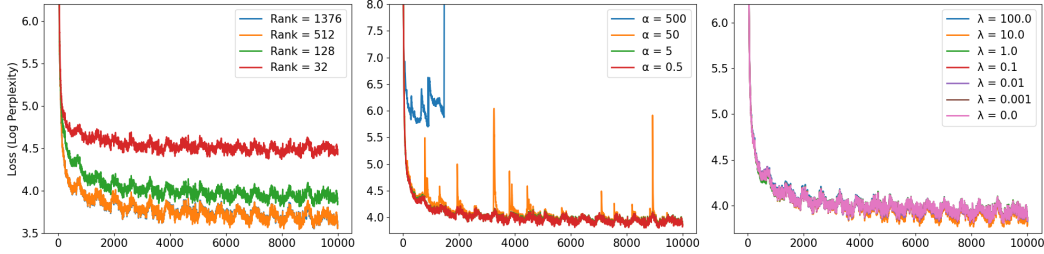


Figure 3: From left to right are loss curves of 10K steps on LLaMA 60M: leftmost is the sweep of rank, middle is the sweep of  $\alpha$  and rightmost is the sweep of  $\lambda$ .

#### 5.4 Can Online Subspace Descent be Applied to Different Optimizers?

One straightforward extension of Online Subspace Descent is to apply it to other base optimizers beyond AdamW8bit. We conduct ablation studies on LION [6] and Adafactor [21], finding that Online Subspace Descent behaves similarly to how it does with AdamW8bit. Despite the initial observation that updating  $P_t$  with AdamW8bit consistently yields better results, we discover that updating  $P_t$  with simple SGD can achieve similar performance.

Method	Galore		Ours			
	Lion	Adaf.	Lion+Lion	Adaf.+Adaf.	Lion+AdamW	Adaf.+AdamW
<b>Perplexity</b>	46.90	34.32	57.97	47.61	<b>44.76</b>	<b>34.15</b>

Table 2: LLaMA 60M on C4 with sequence length 1024, with optimizers on  $P_t$  and  $W_t$ , denote as "Ours ( $W_t$  optimizer) + ( $P_t$  optimizer)". Adaf., and Adam refer to Adafactor and 8bit-AdamW, respectively.

#### 5.5 Can Online Subspace Descent Scale to Larger Model?

We pretrain from scratch a 7B LLaMA model on the C4 dataset for 10K steps, where the  $P_t$  matrix is updated by SGD. The perplexity is the lower the better. The final perplexity and training wall-clock time are provided in Table 3. We further provide the downstream evaluation of the pretrained checkpoints using Galore and our method on the GLUE benchmark in Table 4. Our method consistently outperforms Galore when the model size scales up.

Method	Perplexity	Wall Clock Time (hours)
Galore	51.21	9.7439
Ours	<b>43.72</b>	<b>7.1428</b>

Table 3: Perplexity and Wall Clock Time for 7B models pretrained on C4 for 10K steps. Lower perplexity is better. Online Subspace Descent can be upto 1.3x faster than Galore.

## 6 Related Works

We discuss related works on memory-efficient optimization and low-rank adaptation techniques.

Method	MRPC	RTE	SST2	MNLI	QNLI	QQP	AVG
Galore	0.6838	<b>0.5018</b>	0.5183	0.3506	0.4946	0.3682	0.4862
Ours	<b>0.6982</b>	0.4901	<b>0.5233</b>	<b>0.3654</b>	<b>0.5142</b>	<b>0.3795</b>	<b>0.4951</b>

Table 4: Standardized GLUE evaluation for 7B model checkpoints using eval-harness. Results are reported for various downstream tasks.

**Low-Rank Adaptation** Low-Rank Adaptation (LoRA) [11] adds a low-rank adaptor to specific linear layers in a model, and finetune only the low-rank adaptor. As the adaptors are small, LoRA is widely applied for finetuning large models. Many variants have been proposed since LoRA, including support for multi-task learning Wang et al. [23] and further memory reductions Dettmers et al. [8]. Notably, Lialin et al. [15] proposed ReLoRA for pretraining, requiring a full-rank training warmup to match standard performance levels. It’s important to note that LoRA is fundamentally distinct from subspace descent. While subspace descent optimizes within the original model parameter space, LoRA focuses its optimization efforts within the space of the adaptors.

**Memory-Efficient Optimization** Several approaches aim to reduce memory costs associated with gradient statistics in adaptive optimization algorithms [21, 2, 7]. In particular, Adafactor [21] factorizes the second-order statistics by a row-column outer product and update the factorized bases on the fly, hence achieving a sub-linear memory cost. K-Fac [19] presents a factorized approximation of the Fisher information matrix which leads to a sublinear natural gradient method. More recently, Feinberg et al. [9] observes that the spectra of the Kronecker-factored gradient covariance matrix in deep learning (DL) training tasks are concentrated on a small leading eigenspace and propose to maintain a matrix preconditioner using the frequent directions sketch. However, their method requires conducting the eigendecomposition at every step, which can be costly for large models. Other than factorization methods, quantization techniques [7, 1, 24, 16] are also widely used, where the gradient (or the momentum and the preconditioner) are directly quantized to tradeoff performance for memory. Fused gradient computation method [17] have also been used to minimize memory costs during training. Galore [25] is the most relevant work to ours. Galore focuses on low-rank gradient structures, reducing memory costs for both first and second-order statistics. Our method can be viewed as a general extension to Galore where we replace the infrequent SVD by a continuous subspace descent [14, 10]. As a result, our method not only provides a more general framework to study memory-efficient subspace descent, but is also more performant than Galore in practice.

## 7 Conclusion

In conclusion, we provide the first convergence guarantee for arbitrary update rules of projection matrix, applicable to a range of optimizers that can be analyzed using Hamiltonian Descent, including common ones like LION, AdamW, and Adafactor. Inspired by this theoretical foundation, we introduce Dynamic Subspace Descent, a novel family of subspace descent optimizers that eschews SVD in favor of online PCA for updating projection matrix. Dynamic Subspace Descent is both flexible and minimally intrusive, and our experiments show that it achieves lower perplexity in pretraining LLaMA models (ranging from 60M to 1B parameters) on the C4 dataset compared to state-of-the-art low-rank training methods, while also closing the perplexity gap with full-rank baselines.

For future research, we propose several open and intriguing questions: (1) Are there alternative methods for updating projection matrix that could accelerate convergence? (2) What is the impact of weight decay on convergence in Dynamic Subspace Descent? (3) Can low-rank gradients and updates be combined with dynamic low-rank weights (e.g., Mixture of Experts) to further enhance training efficiency? (4) Can this method be applied to problems beyond language modeling? We hope that our work provides a strong foundation for exploring these questions.

## References

- [1] Dan Alistarh, Demjan Grubic, Jerry Li, Ryota Tomioka, and Milan Vojnovic. Qsgd: Communication-efficient sgd via gradient quantization and encoding. *Advances in neural information processing systems*, 30, 2017.
- [2] Rohan Anil, Vineet Gupta, Tomer Koren, and Yoram Singer. Memory efficient adaptive optimization. *Advances in Neural Information Processing Systems*, 32, 2019.
- [3] Pierre Baldi and Kurt Hornik. Neural networks and principal component analysis: Learning from examples without local minima. *Neural networks*, 2(1):53–58, 1989.
- [4] Lizhang Chen, Bo Liu, Kaizhao Liang, and Qiang Liu. Lion secretly solves constrained optimization: As lyapunov predicts. *arXiv preprint arXiv:2310.05898*, 2023.
- [5] Xiangning Chen, Chen Liang, Da Huang, Esteban Real, Kaiyuan Wang, Yao Liu, Hieu Pham, Xuanyi Dong, Thang Luong, Cho-Jui Hsieh, et al. Symbolic discovery of optimization algorithms. *arXiv preprint arXiv:2302.06675*, 2023.
- [6] Xiangning Chen, Chen Liang, Da Huang, Esteban Real, Kaiyuan Wang, Hieu Pham, Xuanyi Dong, Thang Luong, Cho-Jui Hsieh, Yifeng Lu, et al. Symbolic discovery of optimization algorithms. *Advances in Neural Information Processing Systems*, 36, 2024.
- [7] Tim Dettmers, Mike Lewis, Sam Shleifer, and Luke Zettlemoyer. 8-bit optimizers via block-wise quantization. *arXiv preprint arXiv:2110.02861*, 2021.
- [8] Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. Qlora: Efficient finetuning of quantized llms. *Advances in Neural Information Processing Systems*, 36, 2024.
- [9] Vladimir Feinberg, Xinyi Chen, Y Jennifer Sun, Rohan Anil, and Elad Hazan. Sketchy: Memory-efficient adaptive regularization with frequent directions. *Advances in Neural Information Processing Systems*, 36, 2024.
- [10] Guy Gur-Ari, Daniel A Roberts, and Ethan Dyer. Gradient descent happens in a tiny subspace. *arXiv preprint arXiv:1812.04754*, 2018.
- [11] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021.
- [12] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [13] David Kozak, Stephen Becker, Alireza Doostan, and Luis Tenorio. Stochastic subspace descent. *arXiv preprint arXiv:1904.01145*, 2019.
- [14] Brett W Larsen, Stanislav Fort, Nic Becker, and Surya Ganguli. How many degrees of freedom do we need to train deep networks: a loss landscape perspective. *arXiv preprint arXiv:2107.05802*, 2021.
- [15] Vladislav Lialin, Sherin Muckatira, Namrata Shivagunde, and Anna Rumshisky. Relora: High-rank training through low-rank updates. In *Workshop on Advancing Neural Network Training: Computational Efficiency, Scalability, and Resource Optimization (WANT@ NeurIPS 2023)*, 2023.
- [16] Bo Liu, Lemeng Wu, Lizhang Chen, Kaizhao Liang, Jiaxu Zhu, Chen Liang, Raghuraman Krishnamoorthi, and Qiang Liu. Communication efficient distributed training with distributed lion. *arXiv preprint arXiv:2404.00438*, 2024.
- [17] Kai Lv, Yuqing Yang, Tengxiao Liu, Qinghui Gao, Qipeng Guo, and Xipeng Qiu. Full parameter fine-tuning for large language models with limited resources. *arXiv preprint arXiv:2306.09782*, 2023.
- [18] Chris J Maddison, Daniel Paulin, Yee Whye Teh, Brendan O’Donoghue, and Arnaud Doucet. Hamiltonian descent methods. *arXiv preprint arXiv:1809.05042*, 2018.

- [19] James Martens and Roger Grosse. Optimizing neural networks with kronecker-factored approximate curvature. In *International conference on machine learning*, pages 2408–2417. PMLR, 2015.
- [20] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *arXiv e-prints*, 2019.
- [21] Noam Shazeer and Mitchell Stern. Adafactor: Adaptive learning rates with sublinear memory cost. In *International Conference on Machine Learning*, pages 4596–4604. PMLR, 2018.
- [22] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- [23] Yiming Wang, Yu Lin, Xiaodong Zeng, and Guannan Zhang. Multilora: Democratizing lora for better multi-task learning. *arXiv preprint arXiv:2311.11501*, 2023.
- [24] Wei Wen, Cong Xu, Feng Yan, Chunpeng Wu, Yandan Wang, Yiran Chen, and Hai Li. Terngrad: Ternary gradients to reduce communication in distributed deep learning. *Advances in neural information processing systems*, 30, 2017.
- [25] 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.

## A Experiments

### A.1 Hyperparameters

We sweep learning rate from [0.01, 0.005, 0.001]. For GaLore as well as Adam8bit, we follow the recommended hyperparameters as much as possible. For instance, GaLore update gap is set to recommended default, 200. Warmup is set to 10% of the total training steps. Batch size is set to 512 and gradient clipping is set to 1.0.

### A.2 Rank Sweep

In the following table, is a sweep on different ranks and their final perplexity of LLaMA 60M (SS = 1024) on C4. All other hyperparameters are fixed and using recommended default. Notice that as the rank increases, both Dynamic Subspace Descent and GaLore improve.

Rank	Perplexity (Ours)	Perplexity (GaLore)
32	85.90	86.16
128	49.01	48.05
512	37.41	36.93
Full	37.18	36.51

Table 5: On LLaMA 60M SS 1024, we sweep across different ranks, the trend is clear and intuitive that higher rank is preferred when it’s feasible.

### A.3 Optimizer Sweep

Method	Perplexity
Lion	52.65
Adafactor	33.45
Adam8bit	29.77
SGD	3469.14
Galore LION	46.90
Galore Adafactor	34.32
Galore AdamW8bit	48.05
Ours LION + LION	57.97
Ours Adafactor + Adafactor	47.61
Ours AdamW8bit + AdamW8bit	49.01
Ours LION + Adam8bit	44.76
Ours Adafactor + AdamW8bit	34.15
Ours AdamW8bit + SGD	53.53

Table 6: In this experiment, we train LLaMA 60M on C4 with sequence length of 1024. We combine different base optimizers to update both  $P_t$  and  $W_t$ , denote as "Ours weight optimizer +  $P_t$  optimizer".