

648 A USE OF LARGE LANGUAGE MODELS (LLMs)  
649  
650

651 In preparing this manuscript, we employed a Large Language Model (LLM) as a general-purpose  
652 writing assistant. Specifically, the LLM was used to polish the language, improve clarity and flow,  
653 and enhance the presentation of the text. All technical content, experimental design, data analysis,  
654 and model development were performed independently by the authors. The LLM was not used to  
655 generate any novel scientific ideas, experimental results, or interpretations.

656  
657 B DERIVATION OF PROJECTION-AWARE MOMENT UPDATES  
658  
659

660 The projection-aware update rule for the first-moment estimate arises from the multiplications with the  
661 matrices  $U_{t,L}^T U_{t-1,L} \in \mathbb{R}^{r_1 \times r_1}$  and  $U_{t-1,R}^T U_{t,R} \in \mathbb{R}^{r_2 \times r_2}$ . This enables an efficient shift between  
662 subspaces without generating any intermediate high-dimensional matrices. Since both matrices  
663  $U_{t,L}^T U_{t-1,L}$  and  $U_{t-1,R}^T U_{t,R}$  are orthogonal, they represent the change of basis between the two  
664 subspaces. Let  
665

$$666 B_{t-1,L} = [b_{t-1,L}^1, \dots, b_{t-1,L}^{r_1}] \quad \text{and} \quad B_{t,L} = [b_{t,L}^1, \dots, b_{t,L}^{r_1}],$$

667 denote orthonormal bases for the subspaces  $U_{t-1,L}$  and  $U_{t,L}$  at time steps  $t-1$  and  $t$ , respectively.  
668 Similarly, let  
669

$$670 B_{t-1,R} = [b_{t-1,R}^1, \dots, b_{t-1,R}^{r_2}] \quad \text{and} \quad B_{t,R} = [b_{t,R}^1, \dots, b_{t,R}^{r_2}],$$

671 denote orthonormal bases for the subspaces  $U_{t-1,R}$  and  $U_{t,R}$  at time steps  $t-1$  and  $t$ , respectively.  
672 Let  $A_{t-1} \in \mathbb{R}^{r_1 \times r_2}$  and we want to change its basis to  $A_t = U_{t,L}^T U_{t-1,L} A_{t-1} U_{t-1,R}^T U_{t,R}$ , then the  
673  $p$ -th row and the  $q$ -th column of a matrix  $A_{t-1}$  transforms under the change of basis to time step  $t$  is  
674

$$675 A_t^{ij} = \sum_{p=1}^{r_1} \sum_{q=1}^{r_2} \langle b_{t,L}^i, b_{t-1,L}^p \rangle A_{t-1}^{pq} \langle b_{t-1,R}^q, b_{t,R}^j \rangle, \quad (13)$$

676 where we use superscripts to denote the elements of a matrix. Followed by analysis by (Robert et al.,  
677 2024), the first and second moments in Adam represent the exponentially time-weighted expectation  
678 at time  $t$  with decay parameter  $\beta$ , i.e.,  $M_t = \mathbb{E}_{t,\beta_1}[S_t]$  and  $V_t = \mathbb{E}_{t,\beta_2}[S_t^2]$ . The first-moment  
679 estimate can be expressed under a change of basis through the transformation matrices  $U_{t,L}^T U_{t-1,L}$   
680 and  $U_{t-1,R}^T U_{t,R}$ . In particular, if  $b_{t,L}^i$  denote the  $i$ -th row basis vector at step  $t$ , then the  $ij$ -th entry of  
681 the matrix  $S_t$  at time step  $t$  is denoted by  $\langle b_{t,L}^i, S_t b_{t,R}^j \rangle$ , when it has  $B_{t,L}$  and  $B_{t,R}$  as left and right  
682 subspaces at time  $t$ . With this notation and eqn. equation 13

$$683 M_t^{ij} = \mathbb{E}_{t,\beta_1}[S_t^{ij}] = \mathbb{E}_{t,\beta_1}[\langle b_{t,L}^i, S_t b_{t,R}^j \rangle]$$

$$684 = \sum_{p=1}^{r_1} \sum_{q=1}^{r_2} \langle b_{t,L}^i, b_{t-1,L}^p \rangle \mathbb{E}_{t,\beta_1}[\langle b_{t-1,L}^p, S_t b_{t-1,R}^q \rangle] \langle b_{t-1,R}^q, b_{t,R}^j \rangle$$

$$685 = (U_{t,L}^T U_{t-1,L} M_{t-1} U_{t-1,R}^T U_{t,R})^{ij}.$$

702 With the same approach, we can change the basis for the second moment as well:  
 703

$$\begin{aligned}
 704 \quad V_t^{ij} &= \mathbb{E}_{t,\beta_2}[(S_t^{ij})^2] = \mathbb{E}_{t,\beta_2}[\langle b_{t,L}^i, S_t b_{t,R}^j \rangle^2] \\
 705 &= \left( \sum_{p=1}^{r_1} \sum_{q=1}^{r_2} \langle b_{t,L}^i, b_{t-1,L}^p \rangle \mathbb{E}_{t,\beta_2}[\langle b_{t-1,L}^p, S_t b_{t-1,R}^q \rangle] \langle b_{t-1,R}^q, b_{t,R}^j \rangle \right)^2 \\
 706 &= \sum_{p=1}^{r_1} \sum_{q=1}^{r_2} \langle b_{t,L}^i, b_{t-1,L}^p \rangle^2 \mathbb{E}_{t,\beta_2}[\langle b_{t-1,L}^p, S_t b_{t-1,R}^q \rangle^2] \langle b_{t-1,R}^q, b_{t,R}^j \rangle^2 \\
 707 &+ \sum_{k \neq l}^{r_1} \sum_{k' \neq l'}^{r_2} \langle b_{t,L}^i, b_{t-1,L}^k \rangle \mathbb{E}_{t,\beta_2}[\langle b_{t-1,L}^k, S_t b_{t-1,R}^{k'} \rangle] \langle b_{t-1,R}^{k'}, b_{t,R}^j \rangle \\
 708 &\quad \langle b_{t,L}^i, b_{t-1,L}^l \rangle \mathbb{E}_{t,\beta_2}[\langle b_{t-1,L}^l, S_t b_{t-1,R}^{l'} \rangle] \langle b_{t-1,R}^{l'}, b_{t,R}^j \rangle \\
 709 &= \sum_{p=1}^{r_1} \sum_{q=1}^{r_2} \langle b_{t,L}^i, b_{t-1,L}^p \rangle^2 V_{t-1}^{pq} \langle b_{t-1,R}^q, b_{t,R}^j \rangle^2 \\
 710 &+ \sum_{k \neq l}^{r_1} \sum_{k' \neq l'}^{r_2} \langle b_{t,L}^i, b_{t-1,L}^k \rangle M_{t-1}^{k,k'} \langle b_{t-1,L}^{k'}, b_{t-1,L}^j \rangle \\
 711 &\quad \langle b_{t,L}^i, b_{t-1,L}^l \rangle M_{t-1}^{l,l'} \langle b_{t-1,R}^{l'}, b_{t,R}^j \rangle. \tag{14}
 712 \\
 713 &
 \end{aligned}$$

714 We employ the following equation to rewrite the second term in the last equality:  
 715

$$\begin{aligned}
 716 \quad &\sum_{k,l} \sum_{k',l'} \langle b_{t,L}^i, b_{t-1,L}^k \rangle M_{t-1}^{k,k'} \langle b_{t-1,L}^{k'}, b_{t-1,L}^j \rangle \langle b_{t,L}^i, b_{t-1,L}^l \rangle M_{t-1}^{l,l'} \langle b_{t-1,R}^{l'}, b_{t,R}^j \rangle \\
 717 &= \sum_{k,k'} \langle b_{t,L}^i, b_{t-1,L}^k \rangle^2 (M_{t-1}^{k,k'})^2 \langle b_{t-1,L}^{k'}, b_{t-1,L}^j \rangle^2 + \\
 718 &\quad \sum_{k \neq l} \sum_{k' \neq l'} \langle b_{t,L}^i, b_{t-1,L}^k \rangle M_{t-1}^{k,k'} \langle b_{t-1,L}^{k'}, b_{t-1,L}^j \rangle \langle b_{t,L}^i, b_{t-1,L}^l \rangle M_{t-1}^{l,l'} \langle b_{t-1,R}^{l'}, b_{t,R}^j \rangle. \\
 719 &
 \end{aligned}$$

720 With this, we can rewrite eqn. equation 14:  
 721

$$\begin{aligned}
 722 \quad V_t^{ij} &= \mathbb{E}_{t,\beta_2}[(S_t^{ij})^2] = \sum_{p=1}^{r_1} \sum_{q=1}^{r_2} \langle b_{t,L}^i, b_{t-1,L}^p \rangle^2 V_{t-1}^{pq} \langle b_{t-1,R}^q, b_{t,R}^j \rangle^2 \\
 723 &+ \sum_{k,l} \sum_{k',l'} \langle b_{t,L}^i, b_{t-1,L}^k \rangle M_{t-1}^{k,k'} \langle b_{t-1,L}^{k'}, b_{t-1,L}^j \rangle \langle b_{t,L}^i, b_{t-1,L}^l \rangle M_{t-1}^{l,l'} \langle b_{t-1,R}^{l'}, b_{t,R}^j \rangle \\
 724 &= \sum_{p=1}^{r_1} \sum_{q=1}^{r_2} \langle b_{t,L}^i, b_{t-1,L}^p \rangle^2 V_{t-1}^{pq} \langle b_{t-1,R}^q, b_{t,R}^j \rangle^2 \\
 725 &+ \sum_{k,l} \sum_{k',l'} \langle b_{t,L}^i, b_{t-1,L}^k \rangle M_{t-1}^{k,k'} \langle b_{t-1,L}^{k'}, b_{t-1,L}^j \rangle \langle b_{t,L}^i, b_{t-1,L}^l \rangle M_{t-1}^{l,l'} \langle b_{t-1,R}^{l'}, b_{t,R}^j \rangle \\
 726 &- \sum_{k,k'} \langle b_{t,L}^i, b_{t-1,L}^k \rangle^2 (M_{t-1}^{k,k'})^2 \langle b_{t-1,R}^{k'}, b_{t,R}^j \rangle^2 \\
 727 &= \sum_{p=1}^{r_1} \sum_{q=1}^{r_2} \langle b_{t,L}^i, b_{t-1,L}^p \rangle^2 (V_{t-1}^{pq} - (M_{t-1}^{pq})^2) \langle b_{t-1,R}^q, b_{t,R}^j \rangle^2 \\
 728 &+ \left( \langle b_{t,L}^i, b_{t-1,L}^p \rangle M_{t-1}^{p,q} \langle b_{t-1,L}^q, b_{t-1,L}^j \rangle \right)^2 \\
 729 &= (U_{t,L}^T U_{t-1,L})^2 (V_{t-1} - M_{t-1}^2) (U_{t-1,R}^T U_{t,R})^2 + (U_{t,L}^T U_{t-1,L} M_{t-1} U_{t-1,R}^T U_{t,R})^2. \\
 730 &
 \end{aligned}$$

## 731 C CONVERGENCE ANALYSIS

732 **Theorem C.1. (Convergence of Clean with fixed projections).** Suppose that the gradient has  
 733 the parametric form  $G_t = \sum_{i=1}^N A_i - \sum_{i=1}^N B_i W_t C_i$  where  $N$  is a batch size and the func-

756 *tions  $A_i, B_i$  and  $C_i$  have  $L_A, L_B$  and  $L_C$  continuity, respectively with respect to  $W$ . Let*  
 757  $\|W\|_F \leq M$  *with  $M$  constant. Define  $\widehat{B_{i,t}} = (\Lambda_{t,L}^i)^{-1/2} (U_{t,L}^i)^\top B_i(W_t) U_{t,L}^i (\Lambda_{t,L}^i)^{-1/2}$  and*  
 758  $\widehat{C_{i,t}} = (\Lambda_{t,R}^i)^{-1/2} (U_{t,R}^i)^\top C_i(W_t) U_{t,R}^i (\Lambda_{t,R}^i)^{-1/2}$  *where the  $\Lambda_{t,L}^i, U_{t,L}^i, \Lambda_{t,R}^i$  and  $U_{t,R}^i$  are the out-*  
 759 *puts of Algorithm ??.* Also let  $S_t = \Lambda_{t,L}^{-1/2} U_{t,L}^\top G_t U_{t,R} \Lambda_{t,R}^{-1/2}$ ,  $\kappa_t = \frac{1}{N} \sum_i \lambda_{\min}(\widehat{B_{i,t}}) \lambda_{\min}(\widehat{C_{i,t}})$   
 760 *and  $\lambda_{\min}(\Lambda_L), \lambda_{\min}(\Lambda_R) \geq \lambda_0 \geq 0$ . Assuming that the projection matrices remain constant during*  
 761 *the training. Then for the learning rate  $\eta$  and  $\min(\kappa_t) > (L_A + 2L_B L_C M^2)$ , the Clean satisfies*  
 762

$$\begin{aligned}\|S_t\|_F &\leq \frac{\eta}{\lambda_0} (L_A + 2L_B L_C M^2) \|S_{t-1}\|_F + \left(1 - \frac{\eta \kappa_{t-1}}{\lambda_0}\right) \|S_{t-1}\|_F \\ &= \left[1 - \frac{\eta}{\lambda_0} (\kappa_{t-1} - L_A - 2L_B L_C M^2)\right] \|S_{t-1}\|_F.\end{aligned}$$

768 *Proof.* During the proof we use the Sylvester equality. Let  $\otimes$  presents the Kronecker product then  
 769 for arbitrary matrices  $A, B$  and  $X$ ,  $\text{vec}(AXB) = (B^\top \otimes A)\text{vec}(X)$ . By vectorizing the gradient  
 770 parametric form we have

$$g_t = \text{vec}(G_t) = \text{vec}\left(\sum_{i=1}^N A_i - \sum_{i=1}^N B_i W_t C_i\right) = a_t - R_t w_t,$$

775 where  $w_t = \text{vec}(W_t)$ ,  $a_t = \frac{1}{N} \sum_i \text{vec}(A_{i,t})$  and  $R_t = \frac{1}{N} \sum_i C_{i,t} \otimes B_{i,t}$ . Using the Sylvester  
 776 equation, the vectorized form of  $S_t = \Lambda_{t,L}^{-1/2} U_{t,L}^\top G_t U_{t,R} \Lambda_{t,R}^{-1/2}$  is  
 777

$$s_t = \text{vec}(\Lambda_{t,L}^{-1/2} U_{t,L}^\top G_t U_{t,R} \Lambda_{t,R}^{-1/2}) = (\Lambda_{t,R}^{-1/2} U_{t,R}^\top) \otimes (\Lambda_{t,L}^{-1/2} U_{t,L}^\top) \text{vec}(G_t) \\ = (\Lambda_{t,R}^{-1/2} U_{t,R}^\top) \otimes (\Lambda_{t,L}^{-1/2} U_{t,L}^\top) g_t. \quad (15)$$

Moreover,  $\tilde{G}_t = U_{t,L} \Lambda_{t,L}^{-1/2} U_{t,L}^\top G_t U_{t,R} \Lambda_{t,R}^{-1/2} U_{t,R}^\top$  can be written as

$$\begin{aligned}\tilde{g}_t &= \text{vec}(U_{t,L}\Lambda_{t,L}^{-1/2}U_{t,L}^\top G_t U_{t,R}\Lambda_{t,R}^{-1/2}U_{t,R}^\top) = (U_{t,R} \otimes U_{t,L})(\Lambda_{t,R}^{-1/2}U_{t,R}^\top) \otimes (\Lambda_{t,L}^{-1/2}U_{t,L}^\top)g_t \\ &= (U_{t,R} \otimes U_{t,L})s_t.\end{aligned}\quad (16)$$

Suppose that the Nyström subspaces remain fixed over a window of iterations. That is, for the left projections we have  $U_{t,L} = U_R$  and  $\Lambda_{t,L} = \Lambda_L$  and analogously for the right projections. Then  $s_t$  and  $\tilde{q}_t$  can be restated as:

$$s_t = (\Lambda_B^{-1/2} U_B^\top) \otimes (\Lambda_L^{-1/2} U_L^\top) q_t \quad \text{and} \quad \tilde{q}_t = (U_B \otimes U_L) s_t.$$

Now we derive the recursive form of  $a_i$ :

$$\begin{aligned}
g_t &= a_t - R_t w_t = a_t - R_t w_t - g_{t-1} + g_{t-1} \\
&= a_t - R_t w_t - a_{t-1} + R_{t-1} w_{t-1} + a_{t-1} - R_{t-1} w_{t-1} \\
&= a_t - R_t w_t - a_{t-1} + R_{t-1} (w_t - \eta \tilde{g}_t) + a_{t-1} - R_{t-1} w_{t-1} \\
&= (a_t - a_{t-1}) + (R_{t-1} - R_t) w_t + a_{t-1} - R_{t-1} w_{t-1} - \eta R_{t-1} \tilde{g}_t \\
&\equiv e_t + q_{t-1} - \eta R_{t-1} \tilde{g}_t
\end{aligned}$$

where  $e_t = (a_t - a_{t-1}) + (R_{t-1} - R_t)w_t$  and we use  $w_t = w_{t-1} + \eta g_{t-1}$ . By left-multiplying  $\Lambda_R^{-1/2}U_R^\top \otimes (\Lambda_L^{-1/2}U_L^\top)$ , we obtain

$$\begin{aligned}
s_t &= (\Lambda_R^{-1/2} U_R^\top) \otimes (\Lambda_L^{-1/2} U_L^\top) g_t = (\Lambda_R^{-1/2} U_R^\top) \otimes (\Lambda_L^{-1/2} U_L^\top) (e_t + g_{t-1} - \eta R_{t-1} \tilde{g}_t) \\
&= (\Lambda_R^{-1/2} U_R^\top) \otimes (\Lambda_L^{-1/2} U_L^\top) e_t + (\Lambda_R^{-1/2} U_R^\top) \otimes (\Lambda_L^{-1/2} U_L^\top) g_{t-1} \\
&\quad - \eta (\Lambda_R^{-1/2} U_R^\top) \otimes (\Lambda_L^{-1/2} U_L^\top) R_{t-1} \tilde{g}_{t-1} \\
&= (\Lambda_R^{-1/2} U_R^\top) \otimes (\Lambda_L^{-1/2} U_L^\top) e_t + s_{t-1} - \eta (\Lambda_R^{-1/2} U_R^\top) \otimes (\Lambda_L^{-1/2} U_L^\top) R_{t-1} (U_R \otimes U_L) s_{t-1} \\
&= (\Lambda_R^{-1/2} U_R^\top) \otimes (\Lambda_L^{-1/2} U_L^\top) e_t + s_{t-1} - \eta (\Lambda_R^{-1/2} \otimes \Lambda_L^{-1/2}) (U_R \otimes U_L)^\top R_{t-1} (U_R \otimes U_L) s_{t-1}
\end{aligned}
\tag{17}$$

810 Let

$$\begin{aligned}
 812 \quad \hat{R}_t &= (U_R \otimes U_L)^\top R_t (U_R \otimes U_L) = \frac{1}{N} \sum_i (U_R \otimes U_L)^\top (C_{i,t} \otimes B_{i,t}) (U_R \otimes U_L) \\
 813 \\
 814 &= \frac{1}{N} \sum_i (U_R^\top C_{i,t} U_R) \otimes (U_L^\top B_{i,t} U_L),
 815
 \end{aligned}$$

816 then eqn. equation 17 can be presented as

$$818 \quad s_t = (\Lambda_R^{-1/2} \otimes \Lambda_L^{-1/2})(U_R \otimes U_L)^\top e_t + (I - \eta(\Lambda_R^{-1/2} \otimes \Lambda_L^{-1/2})\hat{R}_{t-1})s_{t-1}. \quad (18)$$

819 Now we need to bound  $s_t$  in above equation. Since  $U_L$  and  $U_R$  are orthonormal matrices we have  
 820  $U_L^\top U_L = I$  and  $U_R^\top U_R = I$ . Therefore by using Sylvester identity and submultiplicativity property of  
 821 the norm,

$$\begin{aligned}
 822 \quad \left\| (\Lambda_R^{-1/2} \otimes \Lambda_L^{-1/2})(U_R \otimes U_L)^\top e_t \right\|_2 &= \left\| \text{vec} \left[ (U_R \otimes U_L)^\top E_t (\Lambda_R^{-1/2} \otimes \Lambda_L^{-1/2}) \right] \right\|_2 \\
 823 \\
 824 &= \left\| (U_R \otimes U_L)^\top E_t (\Lambda_R^{-1/2} \otimes \Lambda_L^{-1/2}) \right\|_F \\
 825 \\
 826 &\leq \left\| (U_R \otimes U_L)^\top E_t \right\|_2 \left\| \Lambda_R^{-1/2} \otimes \Lambda_L^{-1/2} \right\|_F \\
 827 \\
 828 &\leq \frac{\|E_t\|_F}{\lambda_0}, \\
 829
 \end{aligned} \quad (19)$$

830 where  $E_t = \frac{1}{N} \sum_i A_{i,t} - A_{i,t-1} + \frac{1}{N} \sum_i B_{i,t-1} W_t C_{i,t-1} - B_{i,t} W_t C_{i,t}$  is the matrix form of  $e_t$ .  
 831 Therefore, we need to bound  $\|E_t\|_F$ . By Lipschitz-continuity property of  $A_i$ ,  $B_i$  and  $C_i$  we have:

$$\begin{aligned}
 833 \quad \|A_t - A_{t-1}\|_F &\leq L_A \|W_t - W_{t-1}\|_F = \eta L_A \left\| \tilde{G}_{t-1} \right\|_F \leq \eta L_A \|S_{t-1}\|_F \\
 834 \\
 835 \quad \|(B_{t-1} - B_t)W_t C_{t-1}\|_F &\leq L_B \|W_t - W_{t-1}\|_F \|W_t\|_F \|C_{t-1}\|_F \leq \eta L_B L_C M^2 \|S_{t-1}\|_F \\
 836 \\
 837 \quad \|B_t W_t (C_{t-1} - C_t)\|_F &\leq L_C \|W_t - W_{t-1}\|_F \|W_t\|_F \|B_t\|_F \leq \eta L_B L_C M^2 \|S_{t-1}\|_F.
 \end{aligned}$$

838 Thus the upper bound for  $\|E_t\|_F$  can be derived:

$$\begin{aligned}
 839 \quad \frac{\|E_t\|_F}{\lambda_0} &\leq \frac{\eta L_A \|S_{t-1}\|_F + \eta L_B L_C M^2 \|S_{t-1}\|_F + \eta L_B L_C M^2 \|S_{t-1}\|_F}{\lambda_0} \\
 840 \\
 841 &\leq \frac{\eta}{\lambda_0} (L_A + 2L_B L_C M^2) \|S_{t-1}\|_F. \\
 842
 \end{aligned} \quad (20)$$

843 To find an upper bound of  $\|S_t\|_F$  we also need to find an upper bound for  $I - \eta(\Lambda_L^{-1/2} \otimes \Lambda_R^{-1/2})\hat{R}_{t-1}$ .  
 844 This involves finding the minimum eigenvalue of  $\hat{R}_t$ . Let  $\lambda_{i,t} = \lambda_{\min}(U_R^\top C_{i,t} U_R)$  and  $\lambda'_{i,t} =$   
 845  $\lambda_{\min}(U_L^\top B_{i,t} U_L)$ , then

$$\begin{aligned}
 846 \quad \lambda_{\min} \left( (\Lambda_L^{-1/2} \otimes \Lambda_R^{-1/2})(U_R^\top C_{i,t} U_R) \otimes (U_L^\top B_{i,t} U_L) \right) &= \\
 847 \\
 848 \quad \lambda_{\min}(\Lambda_L^{-1/2}) \lambda_{\min}(\Lambda_R^{-1/2}) \lambda_{\min}(U_L^\top B_{i,t} U_L) \lambda_{\min}(U_R^\top C_{i,t} U_R) &\leq \frac{\lambda_{i,t} \lambda'_{i,t}}{\lambda_0}.
 849
 \end{aligned}$$

850 For any vector  $v$  we have

$$\begin{aligned}
 851 \quad v^\top \left[ (\Lambda_L^{-1/2} \otimes \Lambda_R^{-1/2}) \hat{R}_t \right] v &= \frac{1}{N} \sum_i v^\top \left[ (\Lambda_L^{-1/2} \otimes \Lambda_R^{-1/2}) (U_R^\top C_{i,t} U_R) \otimes (U_L^\top B_{i,t} U_L) \right] v \\
 852 \\
 853 &\geq \frac{1}{N} \sum_i \frac{\lambda_{i,t} \lambda'_{i,t}}{\lambda_0} = \frac{\kappa_t}{\lambda_0}
 854
 \end{aligned}$$

855 Therefore,  $\lambda_{\max}(I - \eta(\Lambda_L^{-1/2} \otimes \Lambda_R^{-1/2})\hat{R}_{t-1}) \leq (1 - \frac{\eta \kappa_t}{\lambda_0})$ . By combining eqn. equation 20 and  
 856 eqn. equation 18 we obtain:

$$\begin{aligned}
 857 \quad \|S_t\|_F &\leq \frac{\eta}{\lambda_0} (L_A + 2L_B L_C M^2) \|S_{t-1}\|_F + (1 - \frac{\eta \kappa_t}{\lambda_0}) \|S_{t-1}\|_F \\
 858 \\
 859 &= [1 - \frac{\eta}{\lambda_0} (\kappa_t - L_A - 2L_B L_C M^2)] \|S_{t-1}\|_F.
 860
 \end{aligned}$$

□

864 C.1 RANDOMIZED NYSTRÖM APPROXIMATION  
865866 In this section, we give the randomized Nyström approximation algorithms from (Frangella et al.,  
867 2023; Tropp et al., 2017)  
868870 **Algorithm 2** Randomized Nyström Approximation (RandNystromApprox)  
871

```

872 1: Input: Matrix  $G \in \mathbb{R}^{m \times n}$ , rank size  $r \leq m$ ,  $\Omega \in \mathbb{R}^{m \times r}$ .
873 2:  $\Omega \sim \mathcal{N}(0, I) \in \mathbb{R}^{m \times r}$                                 {random Gaussian matrix}
874 3:  $\Omega \leftarrow \text{QR}(\Omega, 0)$                                 {thin QR decomposition}
875 4:  $Y \leftarrow GG^\top \Omega$ 
876 5:  $\nu \leftarrow \text{eps}(\text{norm}(Y, 2))$                                 {compute shift}
877 6:  $Y_\nu \leftarrow Y + \nu\Omega$                                 {add shift for stability}
878 7:  $B \leftarrow \Omega^\top Y_\nu$ 
879 8:  $C \leftarrow \text{CHOL}((B + B^\top)/2)$                                 {Force symmetry and Cholesky decomposition}
880 9:  $B \leftarrow Y_\nu/C$                                 {triangular solve}
881 10:  $[U, \Sigma, \sim] = \text{SVD}(B, 0)$                                 {thin SVD}
882 11:  $\hat{\Lambda} \leftarrow \max\{0, \Sigma^2 - \nu I\}$                                 {remove shift, compute eigs}
883 12: Output:  $[U, \hat{\Lambda}]$ 
884
885
886
```

887 D EXPERIMENTAL SETUP DETAIL  
888889 **Fine-tuning GLUE.** As outlined earlier, we evaluated **CLEAN** by fine-tuning RoBERTa-base  
890 model on the GLUE benchmark. The list of hyperparameters used in this experiment is also reported  
891 in Table 4.894 Table 4: Hyperparameters for fine-tuning RoBERTa-base on GLUE.  
895

	CLEAN	SOAP	LoRA	GaLore	AdamW
Epochs			3		
Warm-up			×		
Batch			16		
Max Len			128		
Data Type			bffloat32		
LR	$\{1e-4, 2e-4, \dots, 5e-4\} / \{1e-5, \dots, 5e-5\}$ linear to 0%				
LR Schedule					
$\beta_1$	0.9	0.95	0.9	0.9	0.9
$\beta_2$	0.999	0.95	0.999	0.999	0.999
Weight Decay	×	×	×	×	×
Dropout	×	×	×	×	×
Grad Clip	×	×	×	×	×
Subspace Update Interval $T$	32	32	×	500	×
Rank $r$	32	×	4/8	4/8	×
Subspace Proj	✓	×	×	✓	×
Accum Weight	✓	✓	×	×	×
Grad Proj	✓	×	×	✓	×

913 **Pre-training on C4.** To illustrate the feasibility of our approach within limited time, we pre-trained  
914 the Llama 130M model using the hyperparameters summarized in Table 5.  
915

918  
919  
920  
921  
922  
923  
924  
925  
926  
927  
928  
929  
930  
931  
932  
933  
934  
935  
936  
937  
938  
939  
940  
941  
942  
943  
944  
945  
946  
947  
948  
949  
950  
951  
952  
953  
954  
955  
956  
957  
958  
959  
960  
961  
962  
963  
964  
965  
966  
967  
968  
969  
970  
971  
Table 5: Hyperparameters for training Llama 130M across different optimizers.

		CLEAN	SOAP	AdamW
Training Steps		1300		
Warm-up Steps		130		
Max Len		1024		
Batch		2000		
Token Batch		2M		
Data Type		bf16		
LR		$\{5e-4, 1e-3, 5e-3, 1e-2\}$		
Warm-up Schedule		linear from 0%		
LR Schedule		cosine to 0%		
$\beta_1$	0.9	0.95	0.9	
$\beta_2$	0.999	0.95	0.999	
$\mu$	0.95	0.95		$\times$
Regularization factor $\rho$	$1e-2$	$\times$	$\times$	
Weight Decay	$\times$	$\times$	$\times$	
Dropout	$\times$	$\times$	$\times$	
Grad Clip	$\times$	$\times$		1.0
Subspace Update Interval $T$	10	10		$\times$
Rank $r$	256	$\times$	$\times$	
Subspace Proj	$\checkmark$	$\checkmark$	$\times$	
Accum Weight	$\checkmark$	$\times$	$\times$	
Grad Proj	$\checkmark$	$\checkmark$	$\times$	

972 USE OF LLMs  
973974 We verify that we use LLMs (e.g., chatGPT, Claude) to polish and rephrase some sentences of the  
975 paper text. We also used it for discovery purposes like finding related works.  
976

977

978

979

980

981

982

983

984

985

986

987

988

989

990

991

992

993

994

995

996

997

998

999

1000

1001

1002

1003

1004

1005

1006

1007

1008

1009

1010

1011

1012

1013

1014

1015

1016

1017

1018

1019

1020

1021

1022

1023

1024

1025