

1                   **STOCHASTIC ASPECTS OF DYNAMICAL LOW-RANK**  
2                   **APPROXIMATION IN THE CONTEXT OF MACHINE LEARNING\***

3                   ARSEN HNATIUK<sup>†</sup>, JONAS KUSCH<sup>‡</sup>, LISA KUSCH<sup>§</sup>, NICOLAS R. GAUGER<sup>¶</sup>, AND  
4                   ANDREA WALTHER<sup>||</sup>

5                   **Abstract.** The central challenges of today’s neural network architectures are the prohibitive  
6 memory footprint and the training costs associated with determining optimal weights and biases. A  
7 large portion of research in machine learning is therefore dedicated to constructing memory-efficient  
8 training methods. One promising approach is dynamical low-rank training (DLRT) which represents  
9 and trains parameters as a low-rank factorization. While DLRT is equipped with several beneficial  
10 properties, analytic results are currently limited to deterministic gradient flows. In this work, we show  
11 that dynamical low-rank training in combination with stochastic gradient and momentum methods  
12 fulfills descent guarantees and prove its convergence to an optimal point.

13                   **Key words.** Dynamical Low-Rank Approximation (DLRA), Dynamical Low-Rank Training  
14 (DLRT), Machine Learning, Stochastic Gradient Descent, Deep Neural Networks

15                   **MSC codes.** 68Q25, 68T05, 68T07, 65K10, 93E35, 65L20, 65L70

16                   **1. Introduction.** In recent years, deep neural networks (DNNs) have consis-  
17 tently and radically redefined the state-of-the-art in tasks related to computer vision,  
18 natural language processing, and many others. This was in part made possible by  
19 the rapid increase in the size of the networks, with some of the newest DNN models  
20 having trillions of parameters [23]. Despite this, there is evidence that large model  
21 sizes are not necessary for good performance; large DNNs are known to contain pa-  
22 rameter redundancies [11, 21, 8, 2, 9] and experiments show that it is often possible  
23 to substantially reduce the model size with little loss in output quality. However, the  
24 task of finding such smaller yet well-performing models is highly challenging [11].

25                   The process of removing redundant parameters is called network pruning and  
26 various techniques have been proposed to achieve this [2]. One possibility to compress  
27 neural networks during training is low-rank pruning [24, 12, 26, 21]. This work will  
28 focus on the method proposed in [21], which applies the principles of Dynamical Low-  
29 Rank Approximation (DLRA) [15] to the task of training artificial neural networks.  
30 The resulting Dynamical Low-Rank Training (DLRT) offers two main advantages over  
31 conventional low-rank pruning methods: First, it allows for a dynamic adaptation of  
32 the approximation rank during training. Second, its convergence is not slowed down  
33 due to the curvature of the manifold containing low-rank matrices.

34                   The core idea is to train the network while dynamically restricting the rank of its  
35 parameter matrices. This method has been experimentally shown to significantly re-  
36 duce the model size and computational costs while sacrificing little accuracy for fully  
37 connected and convolutional layers [21]. Although much effort has been invested in the  
38 study of error bounds and robustness properties of the DLRA algorithm [15, 19, 5, 13],  
39 its convergence behavior when using stochastic gradients remains unexplored. Since

---

\*Submitted to the editors 03.17.2024.

**Funding:** The research was funded partly by the DFG under Germany’s Excellence Strategy – The Berlin Mathematics Research Center MATH+ (EXC-2046/1, project ID:390685689).

<sup>†</sup>Humboldt-Universität zu Berlin (arsen.hnatiuk@hu-berlin.de).

<sup>‡</sup>Norges miljø- og biovitenskapelige universitet (jonas.kusch@nmbu.no).

<sup>§</sup>Technische Universiteit Eindhoven (l.kusch@tue.nl).

<sup>¶</sup>University of Kaiserslautern-Landau (RPTU) (nicolas.gauger@scicomp.uni-kl.de).

<sup>||</sup>Humboldt-Universität zu Berlin (andrea.walther@math.hu-berlin.de).

the training of machine learning models such as DNNs generally relies on stochastic gradients [3], the study of DLRA in the stochastic setting is fundamental for developing a theoretical understanding of the method proposed in [21]. This work will attempt to fill this gap in the theory of DLRA. The main findings are

1. *Robustness of stochastic gradients.* We prove that the robust error bound of DLRA holds in combination with stochastic gradient descent algorithms.
2. *Descent direction.* We prove that for sufficiently small learning rates, the method in combination with stochastic gradient descent will retain the descent guarantee from the deterministic and time-continuous setting.
3. *Convergence.* We show that DLRA in combination with stochastic gradient and momentum methods will converge to a local minimum if the basis reaches equilibrium.

This article is structured as follows. In Section 2 and Section 3, we provide an overview of neural network training and Dynamical Low-Rank Approximation. Section 4 discusses the robust error bound for DLRA in the presence of stochastic gradients. In Section 5, we investigate the descent direction with stochastic gradients and provide a convergence proof. Lastly, we provide a conclusion and outlook in Section 6.

**2. Recap: Training of deep neural networks.** Deep neural networks (DNNs) are a special type of machine learning models. In their simplest fully connected form, DNNs with  $N$  layers are functions

$$\mathcal{N} : \mathbb{R}^{D_1} \times \mathbb{R}^p \rightarrow \mathbb{R}^{D_2}, \quad (x; \mathcal{W}) \mapsto y,$$

where the vector of trainable parameters  $\mathcal{W} \in \mathbb{R}^p$  is arranged in a sequence of matrices  $W^1 \in \mathbb{R}^{n_{1,1} \times n_{1,2}}, \dots, W^N \in \mathbb{R}^{n_{N,1} \times n_{N,2}}$  of corresponding dimensions  $n_{k,1}, n_{k,2} \in \mathbb{N}$  with  $p = \sum_{i=1}^N n_{i,1} n_{i,2}$ ,  $x = z_0 \in \mathbb{R}^{D_1}$  is the input vector of dimension  $D_1 \in \mathbb{N}$ , and the output  $y \in \mathbb{R}^{D_2}$  of dimension  $D_2 \in \mathbb{N}$  is calculated by an  $N$ -fold nesting of intermediary steps

$$z_k = \phi_k(W^k z_{k-1})$$

for some non-linear functions  $\phi_k$  [21]. Each intermediary step represents a layer of the network and its trainable parameters are the matrices  $W^k$ .

The training of a DNN is the task of minimizing a real-valued cost function  $\mathcal{L}$  over the parameter space  $\mathbb{R}^p$ . A common way of performing this optimization is to apply the gradient descent algorithm

$$\mathcal{W}_{t+1} = \mathcal{W}_t - h \nabla_{\mathcal{W}} \mathcal{L}(\mathcal{W}_t),$$

where  $h$  is the step size, or, as it is called in the context of machine learning, the *learning rate* [3, 25]. This algorithm starts in some point  $\mathcal{W}_0 \in \mathbb{R}^p$  and successively moves in the direction of steepest descent  $-\nabla_{\mathcal{W}} \mathcal{L}(\mathcal{W}_t)$  of the objective function until some optimality criterion is reached. It is also possible to consider a setting where the parameters  $\mathcal{W}(t)$  evolve smoothly in time. In this case,  $\mathcal{W}(t)$  can be expressed as a solution to the differential equation (also called the gradient flow)

$$(2.1) \quad \dot{\mathcal{W}}(t) = -\nabla_{\mathcal{W}} \mathcal{L}(\mathcal{W}(t)), \quad \mathcal{W}(t_0) = \mathcal{W}_0.$$

81 Typically, the cost function  $\mathcal{L}$  has the form

$$82 \quad \mathcal{L} : \mathbb{R}^p \rightarrow \mathbb{R}, \quad \mathcal{W} \mapsto \frac{1}{d} \sum_{i=1}^d l(x_i, \mathcal{W}),$$

83 where  $\{x_1, \dots, x_d\}$  is a training data set, and  $l$  is a so-called loss function, which  
 84 quantifies the difference between the model output for a single point  $x_i$  and the desired  
 85 output, given the parameters  $\mathcal{W}$  [3]. We usually assume that  $l$  is differentiable, so the  
 86 gradient of  $\mathcal{L}$  can be written as

$$87 \quad \nabla_{\mathcal{W}} \mathcal{L}(\mathcal{W}) = \frac{1}{d} \sum_{i=1}^d \nabla_{\mathcal{W}} l(x_i, \mathcal{W}).$$

88 Since, as in (2.1), we are mostly interested in the negative gradient, we will use the  
 89 shorthand notation  $F(\mathcal{W}) = -\nabla_{\mathcal{W}} \mathcal{L}(\mathcal{W})$ .

90 As the size  $d$  of the training data set gets larger, the cost of computing the full  
 91 gradient  $F$  becomes prohibitively high [3]. It is thus usual to only use a small subset  
 92 of the data to compute the so-called *minibatch* or *stochastic* gradient

$$93 \quad f(\mathcal{W}) := -\frac{1}{s} \sum_{i=1}^s \nabla_{\mathcal{W}} l(\xi^i, \mathcal{W}),$$

94 where  $s \ll d$  is fixed and the  $\xi^i$  are i.i.d. random variables that follow a uniform  
 95 distribution over the training data set  $\{x_1, \dots, x_d\}$  [3]. We can also write  $f(\mathcal{W}, \xi)$  to  
 96 underline the presence of randomness, where  $\xi = (\xi^1, \dots, \xi^s)$ . In a setting like that  
 97 in (2.1), one would in practice use the stochastic gradient  $f$  instead of  $F$  [3].

98 By construction, it is clear that  $\mathbb{E}_{\xi}[f(\mathcal{W}, \xi)] = F(\mathcal{W})$ . By the law of large num-  
 99 bers, for a large enough  $s$ , we can expect  $f$  to come arbitrarily close to  $F$ .

100 **3. Recap: Dynamical Low-Rank Approximation.** The fact that the para-  
 101 meters of a neural network naturally appear as matrices is central to the pruning  
 102 strategy proposed in [21]. If we can approximate a parameter matrix  $W \in \mathbb{R}^{m \times n}$  by  
 103 a matrix  $Y \in \mathbb{R}^{m \times n}$  of rank  $q \ll \min\{m, n\}$  while maintaining good performance of  
 104 the network, we can significantly reduce the number of trainable parameters and the  
 105 associated computational costs, both for training and for inference. We can calculate  
 106 that while the matrix  $W$  has  $nm$  entries, we can encode  $Y$  in only  $mq + nq + q^2$  en-  
 107 tries since the singular value decomposition allows us to write any rank- $q$  matrix  $Y$  as  
 108  $Y = USV^T$ , with  $S \in \mathbb{R}^{q \times q}$ . The main motivation behind the approach proposed in  
 109 [21] is to represent and train such a low-rank approximation without computing and  
 110 storing full-rank parameter matrices. Such an efficient and robust training method  
 111 is derived by the use of Dynamical Low-Rank Approximation [15], which is a model  
 112 order reduction technique for time-dependent matrices. In this section, we present  
 113 the overall principle of [15] and explore practical implementations.

114 Throughout this work, let  $\|\cdot\|$  and  $\langle \cdot, \cdot \rangle$  refer to the Frobenius norm and scalar  
 115 product. Also, for a matrix  $U$ , let  $P_U = UU^T$  be the projection onto the space  
 116 spanned by the columns of  $U$ . Lastly, we will generally denote full-rank matrices by  
 117 the letter  $W$ , while low-rank matrices will be denoted by  $Y$ .

118 **3.1. Rank- $q$  Approximation.** For  $q \in \mathbb{N}$  and  $m, n \in \mathbb{N}$  such that  $m, n \geq q$ ,  
 119 the space of rank- $q$   $\mathbb{R}^{m \times n}$  matrices is a smooth manifold, as seen in Example 8.14 of  
 120 [17], and we will denote it with  $\mathcal{M}_q = \mathcal{M}_q^{m \times n}$ . Let us further denote with  $\mathcal{T}_Y \mathcal{M}_q$

121 the tangent space of  $\mathcal{M}_q$  at  $Y \in \mathcal{M}_q$  and with  $P(Y)$  the orthogonal projection onto  
 122  $\mathcal{T}_Y \mathcal{M}_q$ .

123 The starting point for Dynamical Low-Rank Approximation (DLRA), first intro-  
 124 duced in [15], is the task of approximating time-dependent matrices  $W(t) \in \mathbb{R}^{m \times n}$ ,  
 125 smooth in  $t$ , by matrices  $Y(t) \in \mathcal{M}_q$  of rank  $q < \min\{m, n\}$ . **why  $<$  and not  $\leq$  as in**  
 126 **other places?** This task can be solved by finding elements of

$$127 \quad \arg \min_{Y(t) \in \mathcal{M}_q} \|Y(t) - W(t)\|.$$

128 This simple approach faces many challenges, such as the need to calculate a costly  
 129 singular value decomposition for each time value  $t$  and the fact that it yields a solution  
 130  $Y(t)$  that is generally not smooth in  $t$  [15].

131 An alternative method of finding rank- $q$  approximations that avoids these issues  
 132 consists of approximating the initial value  $W(t_0)$  and the derivative  $\dot{W}(t)$  instead of  
 133  $W(t)$  itself. The task is thus to find a solution  $Y(t)$  of

$$134 \quad \dot{Y}(t) \in \arg \min_{\dot{Y}(t) \in \mathcal{T}_{Y(t)} \mathcal{M}_q} \|\dot{Y}(t) - \dot{W}(t)\|, \quad Y(t_0) = Y_0$$

135 or, equivalently,

$$136 \quad (3.1) \quad \dot{Y}(t) = P(Y(t))\dot{W}(t), \quad Y(t_0) = Y_0$$

137 with a  $Y_0 \in \mathcal{M}_q$  that approximates  $W(t_0)$  [15].

138 A common situation and one that we will explore from now on is where  $W(t)$  is  
 139 a solution of the matrix differential equation

$$140 \quad (3.2) \quad \dot{W}(t) = F(W(t)), \quad W(t_0) = W_0$$

141 for some smooth function  $F : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}^{m \times n}$ . The approach (3.1) fits naturally  
 142 in this setting, and if  $\dot{W}(t) = F(W(t))$  is not known, we can replace it with its  
 143 approximation  $F(Y(t))$  [15].

144 Approximating  $W(t)$  by the rank- $q$  solution  $Y(t)$  of

$$145 \quad (3.3) \quad \dot{Y}(t) = P(Y(t))F(Y(t)) \in \mathcal{T}_{Y(t)} \mathcal{M}_q, \quad Y(t_0) = Y_0 \in \mathcal{M}_q,$$

146 where  $Y_0$  is a rank- $q$  approximation of  $W_0$ , is the formulation of DLRA which we will  
 147 use in this work.

148 **3.2. Robust Numerical Integrators.** As already discussed, it is computation-  
 149 ally advantageous to work with rank- $q$  matrices  $Y \in \mathbb{R}^{m \times n}$  in their decomposed form  
 150

$$151 \quad (3.4) \quad Y = USV^\top,$$

152 where  $U \in \mathbb{R}^{m \times q}$  and  $V \in \mathbb{R}^{n \times q}$  have orthonormal columns and  $S \in \mathbb{R}^{q \times q}$  is invertible.

153 Therefore, and since computing singular value decompositions at each time step  
 154 is expensive, numerical integrations of (3.3) attempt to integrate the  $U$ ,  $S$ , and  $V$   
 155 matrices separately. Proposition 2.1 in [15] uses (3.3) to derive this as the following  
 156 system of equations.

$$157 \quad \begin{aligned} \dot{S}(t) &= U(t)^\top F(Y(t))V(t), \\ (3.5) \quad \dot{U}(t) &= (I - P_{U(t)})F(Y(t))V(t)S(t)^{-1}, \\ 159 \quad \dot{V}(t) &= (I - P_{V(t)})F(Y(t))^\top U(t)S(t)^{-\top}. \end{aligned}$$

160 From this, Lemma 4.1 in [15] yields an expression for the projector  $P(Y)$ .

$$161 \quad (3.6) \quad P(Y)Z = ZP_V - P_U ZP_V + P_U Z$$

162 for any  $Z \in \mathbb{R}^{m \times n}$ . This representation is unique since  $P_U$  and  $P_V$  project onto the  
163 range and co-range of  $Y$  and thus are uniquely determined by  $Y$ .

164 The numerical issues appearing in (3.5) when  $S$  contains small singular values  
165 have pushed researchers to develop more robust integrators [19, 7, 5, 6, 14, 4, 16].

---

**Algorithm 3.1** (Rank-Adaptive) Basis-Update and Galerkin (BUG) Integrator [5],  
as used in [21].

---

**Input:**  $Y_0 = U_0 S_0 V_0^\top \in \mathcal{M}_{q_0}$  as in (3.4), an initial rank  $q_0$ , and a truncation tolerance  
 $\vartheta > 0$ .

**for**  $k = 0, 1, \dots$  **and**  $t_0 < t_1 < \dots < t_k < \dots$  **do**

**K-step:** solve the  $\mathbb{R}^{m \times q_k}$  differential equation

$$\dot{K}(t) = F(K(t)V_k^\top)V_k, \quad K(t_k) = U_k S_k \text{ over } [t_k, t_{k+1}].$$

**if rank-adaptive then**

        | Set  $K(t_{k+1}) \leftarrow [K(t_{k+1}) \mid U_k]$ .

    Using the QR-decomposition, obtain  $U_{k+1}^* R_{k+1} = K(t_{k+1})$ , where the columns of  
 $U_{k+1}^*$  form an orthonormal basis of the range of  $K(t_{k+1})$ .

    Set  $M \leftarrow U_{k+1}^{*\top} U_k$ .

**L-step:** solve the  $\mathbb{R}^{n \times q_k}$  differential equation

$$\dot{L}(t) = F(U_k L(t)^\top)^\top U_k, \quad L(t_k) = V_k S_k^\top \text{ over } [t_k, t_{k+1}].$$

**if rank-adaptive then**

        | Set  $L(t_{k+1}) \leftarrow [L(t_{k+1}) \mid V_k]$ .

    Using the QR-decomposition, obtain  $V_{k+1}^* \tilde{R}_{k+1} = L(t_{k+1})$ , where the columns of  
 $V_{k+1}^*$  form an orthonormal basis of the range of  $L(t_{k+1})$ .

    Set  $N \leftarrow V_{k+1}^{*\top} V_k$ .

**S-step:** solve the  $\mathbb{R}^{\tilde{q}_U \times \tilde{q}_V}$  differential equation (for  $\tilde{q}_U, \tilde{q}_V \in [q_k, 2q_k]$ , depending  
on the sizes of  $U_{k+1}^*$  and  $V_{k+1}^*$  respectively)

$$\dot{S}(t) = U_{k+1}^{*\top} F(U_{k+1}^* S(t) V_{k+1}^{*\top}) V_{k+1}^*, \quad S(t_k) = M S_k N^\top \text{ over } [t_k, t_{k+1}].$$

    Set  $S_{k+1}^* \leftarrow S(t_{k+1})$ .

**if rank-adaptive then**

        | **Truncation step:** Compute the singular value decomposition  $S_{k+1}^* = P \Sigma Q^\top$   
with  $\Sigma = \text{diag}(\sigma_i)$ .

        | Determine the maximal set of singular values  $\sigma_i$  of  $S_{k+1}^*$  satisfying  $\sum \sigma_i^2 \leq \vartheta^2$   
and define  $\Sigma_1, P_1$  and  $Q_1$  by removing the rows and columns from  $\Sigma, P$   
and  $Q$  corresponding to those  $\sigma_i$ .

        | Set  $S_{k+1} \leftarrow \Sigma_1, U_{k+1} \leftarrow U_{k+1}^* P_1$  and  $V_{k+1} \leftarrow V_{k+1}^* Q_1$ .

**else**

        | Set  $U_{k+1} \leftarrow U_{k+1}^*, V_{k+1} \leftarrow V_{k+1}^*$  and  $S_{k+1} \leftarrow S_{k+1}^*$ .

    Set  $Y_{k+1} \leftarrow U_{k+1} S_{k+1} V_{k+1}^\top$ .

---

166 The method that will be studied in this work is the Rank-Adaptive (often referred  
167 to as augmented) Basis-Update and Galerkin (BUG) Integrator (Algorithm 3.1), first

168 proposed in [5], which is the rank-adaptive version of the fixed-rank BUG integrator  
 169 [7]. The rank-adaptive modification allows the algorithm to determine the optimal  
 170 rank based on a given threshold automatically. The resulting matrix  $Y_1 = U_1 S_1 V_1^\top$   
 171 after one step of the algorithm is an approximation of  $Y(t_1)$  from (3.3).

172 Notice that in Algorithm 3.1 it is possible to perform the K- and L-steps in  
 173 parallel. Also, notice that in the rank-adaptive setting, the number of singular values  
 174  $\sigma_i$  of  $S_{k+1}^*$  remaining after the truncation step determines the (adaptive) rank  $q_{k+1}$   
 175 of the matrix  $Y_{k+1}$  resulting from the  $k$ -th pass of Algorithm 3.1.

176 An important fact for the subsequent theory is that the spans of the matrices  
 177  $U_{k+1}^*$  and  $V_{k+1}^*$  obtained in the rank-adaptive method also contain the spans of  $U_k$   
 178 and  $V_k$ . Following [6], thus we can write

$$179 \quad (3.7) \quad U_{k+1}^* = [U_k \mid U_k^+],$$

180 where  $U_k^+$  is composed of columns that expand the orthonormal basis spanning  $U_k$ .  
 181 Much in the same manner,  $U_{k+1}^*$  can also be thought of as including and expanding  
 182 upon the spans of  $U_{k+1}$  obtained with either the adaptive or non-adaptive methods.  
 183 Thus, we can also write  $U_{k+1}^* = [U_{k+1} \mid U_{k+1}^+]$ . Analogous expressions hold for  $V_{k+1}^*$ .  
 184

Another important property of Algorithm 3.1 is

$$185 \quad (3.8) \quad \|U_{k+1}^* S_{k+1}^* V_{k+1}^{*\top} - U_{k+1} S_{k+1} V_{k+1}^\top\| \leq \vartheta,$$

186 which follows directly from the construction of the truncation step.

187 **3.3. DLRA for Machine Learning.** As seen in Section 2, the trainable pa-  
 188 rameters in neural networks naturally appear as matrices. The central idea of [21]  
 189 is to apply Algorithm 3.1 to (2.1), thereby leveraging the machinery of DLRA for  
 190 training individual layers of DNNs. By doing so, the optimization parameters are  
 191 restricted to the space of low-rank matrices and network pruning (that is, the task of  
 192 finding accurate models of smaller size) is performed during training itself.

193 The general compatibility between DLRA and the training of DNNs can be seen  
 194 when comparing the equations (2.1) and (3.2). In the notation of (3.2), the  $W$  becomes  
 195 the matrix of trainable parameters of one layer of the network, while  $F(W)$  becomes  
 196 the negative gradient  $-\nabla_W \mathcal{L}(W)$  of the objective function.

197 One major difference between the two methods, however, exists. As seen in  
 198 Section 2, during network training, one generally uses the stochastic gradient  $f$  instead  
 199 of the full gradient  $F$ . This difference gives rise to the central question of this work:  
 200 how does Algorithm 3.1 behave when  $F$  is replaced by  $f$ ?

201 We will now specify some notation. A DNN, by definition, contains multiple lay-  
 202 ers and the implementation in [21] applies Algorithm 3.1 to the parameter matrices  
 203  $W^1, W^2, \dots$  of each layer separately, cycling between the layers at each training step.  
 204 In this work, we focus on training a single layer of a DNN and present results corre-  
 205 sponding to this approach. Because of this, although the objective function  $\mathcal{L}$  depends  
 206 on all the parameter matrices  $W^1, W^2, \dots$ , we will omit writing  $\mathcal{L}(W^1, W^2, \dots)$  and  
 207 instead write  $\mathcal{L}(W)$ . Since the optimization of one layer requires only those gradi-  
 208 ents corresponding to the single parameter matrix in question, we will also write  $\nabla \mathcal{L}$   
 209 instead of  $\nabla_W \mathcal{L}$  and consider the other parameters fixed.

210 It is also worth noting that although we write  $F(W) = -\nabla_W \mathcal{L}(W)$  and theoret-  
 211 ically treat it as such, considerable computational gains can be made if we notice that,  
 212 for example, as in the K-step of Algorithm 3.1,

$$213 \quad F(K(t)V_k^\top)V_k = -\nabla_W \mathcal{L}(K(t)V_k^\top)V_k = -\nabla_K \mathcal{L}(K(t)V_k^\top),$$

214 making the computation of the gradient with respect to the full parameter matrix  $W$   
 215 in the K-, L-, and S-steps unnecessary [21].

216 **3.4. Optimality over the Manifold  $\mathcal{M}_q$ .** Since we are performing optimiza-  
 217 tion over manifolds of low-rank matrices, we should take an interest in the optimality  
 218 conditions that exist there.

219 We know that for a point  $Y \in \mathcal{M}_q$ , a necessary first-order condition for optimality  
 220 over  $\mathcal{M}_q$  is given by

$$221 \quad (3.9) \quad P(Y)\nabla\mathcal{L}(Y) = 0.$$

222 This follows from Proposition 3.3 and Theorem 3.4 in [20].

223 Thanks to the explicit expression (3.6) of  $P(Y)$ , we can directly obtain a some-  
 224 what weaker, but in the context of this work more accessible necessary condition of  
 225 optimality

$$226 \quad (3.10) \quad P_U\nabla\mathcal{L}(Y)P_V = 0.$$

227 We call this condition more accessible because expressions of the form  $P_U\nabla\mathcal{L}(Y)P_V$   
 228 appear very naturally in the theory that we will tackle in Section 5.

229 In Section 5, we show that in the stochastic setting, the DLRA optimization  
 230 algorithm can come arbitrarily close to satisfying these conditions in expectation.

231 **4. Robustness of Stochastic DLRA.** Using the notation of Section 2, con-  
 232 sider the algorithm that results from replacing every  $F(\cdot)$  in Algorithm 3.1 by  $f(\cdot, \xi)$ ,  
 233 where a single realization of  $\xi$  is used per pass in the for-loop. We call it the *Sto-*  
 234 *chastic* Algorithm 3.1. In this section, we explore how well the resulting matrix  $Y_1$   
 235 of one step of the Stochastic Algorithm 3.1 approximates the solution  $W(t_1)$  of (3.2)  
 236 at time  $t_1 = t_0 + h$ ,  $h > 0$ . This extends the results found in [21], where Theorem 1  
 237 states a similar error bound for the deterministic algorithm.

238 We can use the manifold structure of  $\mathcal{M}_q$  to investigate the stochastic gradients.  
 239 Since  $\mathcal{T}_Y\mathcal{M}_q$  is a subspace of  $\mathbb{R}^{m \times n}$  for any  $Y \in \mathcal{M}_q$ , we can decompose the gradients  
 240  $F(Y)$  and  $f(Y)$  into components  $F(Y) = M(Y) + R(Y)$  and  $f(Y) = m(Y) + r(Y)$ ,  
 241 where  $M(Y)$ ,  $m(Y) \in \mathcal{T}_Y\mathcal{M}_q$ . We make the following assumptions:

242 ASSUMPTION 4.1. *There exists an  $\varepsilon > 0$  such that  $\|R(Y)\|, \|r(Y)\| \leq \varepsilon$  for all*  
 243  *$Y \in \mathcal{M}_q$ .*

244 ASSUMPTION 4.2. *The two functions  $F$  and  $f$  are bounded by a constant  $B > 0$*   
 245 *and Lipschitz continuous with respect to  $\|\cdot\|$ . The corresponding Lipschitz constant*  
 246 *is denoted by  $L > 0$ .*

247 ASSUMPTION 4.3. *There exists a constant  $C > 0$  such that  $\|F(Y) - f(Y)\| \leq C$*   
 248 *for all  $Y \in \mathbb{R}^{m \times n}$ .*

249 Assumption 4.1 for  $F$  states that  $F(Y)$  is contained in  $\mathcal{T}_Y\mathcal{M}_q$  up to a small factor  
 250 and is based on empirical observations [21]. Assumption 4.2 for  $F$  is common in  
 251 DLRA theory [21, 13] and smooth optimization in general. The extension of these  
 252 assumptions to  $f$  can be justified by the structure of the gradients; both  $f$  and  $F$   
 253 directly inherit their properties from the loss function  $l$ , as seen in Section 2, so we  
 254 can expect them to behave similarly. The bound on  $\|F(Y) - f(Y)\|$  in Assumption 4.3  
 255 can be justified by the law of large numbers, as explained in Section 2. These three  
 256 assumptions hold for the rest of this section.

257 In the following, let  $W(t)$  be the solution of (3.2) with initial value  $W_0$  and  
 258  $Y_1 = U_1S_1V_1^\top$  be the resulting matrix after one step of the rank-adaptive Stochastic

259 Algorithm 3.1 with initial value  $Y_0 \in \mathcal{M}_q$  that approximates  $W_0$ . We wish to find a  
260 robust error bound for the local error

$$261 \quad (4.1) \quad \|W(t_1) - Y_1\| \leq \|W(t_1) - P_{U_1^*} W(t_1) P_{V_1^*}\| + \|P_{U_1^*} W(t_1) P_{V_1^*} - Y_1^*\| + \|Y_1^* - Y_1\|.$$

262 We know from (3.8) that the last term can be bounded by  $\vartheta$ . The bounds for the  
263 other terms are given by the following lemmas. In the following, let us assume that  
264  $\|W_0 - Y_0\| \leq \delta$  holds.

265 LEMMA 4.4. *There exists a constant  $C_1$  such that*

$$266 \quad \|W(t_1) - P_{U_1^*} W(t_1) P_{V_1^*}\| \leq C_1 h^2 + 2(L\delta + \varepsilon + C)h + 2\delta,$$

267 where  $C_1$  does not depend on the condition number of the coefficient matrix  $S$ .

268 *Proof.* We can write

$$269 \quad \|W(t_1) - P_{U_1^*} W(t_1) P_{V_1^*}\| \leq \|W(t_1) - P_{U_1^*} W(t_1)\| + \|P_{U_1^*} (W(t_1) - W(t_1) P_{V_1^*})\| \\ 270 \quad \leq \|W(t_1) - P_{U_1^*} W(t_1)\| + \|W(t_1) - W(t_1) P_{V_1^*}\|.$$

271 First notice that

$$272 \quad \|W_0 - P_{U_1^*} W_0\| \leq \|(I - P_{U_1^*})Y_0\| + \|(I - P_{U_1^*})(W_0 - Y_0)\| = \|(I - P_{U_1^*})(W_0 - Y_0)\|,$$

273 where  $(I - P_{U_1^*})Y_0 = 0$  holds by construction of  $U_1^*$ . Then, using the above and  
274 Assumptions 4.1-4.3 we have

$$275 \quad \|W(t_1) - P_{U_1^*} W(t_1)\| \leq \int_{t_0}^{t_1} \|F(W(t)) - P_{U_1^*} F(W(t))\| dt + \delta \\ 276 \quad \leq \int_{t_0}^{t_1} \|(I - P_{U_1^*})f(W(t))\| dt + Ch + \delta \\ 277 \quad \leq \int_{t_0}^{t_1} \|(I - P_{U_1^*})f(W_0)\| dt + LBh^2 + Ch + \delta \\ 278 \quad \leq \int_{t_0}^{t_1} \|(I - P_{U_1^*})f(Y_0)\| dt + (hL + 1)\delta + LBh^2 + Ch \\ 279 \quad \leq \int_{t_0}^{t_1} \|(I - P_{U_1^*})P(Y_0)f(Y_0)\| dt + h\varepsilon + (hL + 1)\delta + LBh^2 + Ch \\ 280 \quad = \int_{t_0}^{t_1} \|(I - P_{U_1^*})f(Y_0)P_{V_0}\| dt + h\varepsilon + (hL + 1)\delta + LBh^2 + Ch \\ 281 \quad \leq \frac{1}{h} \int_{t_0}^{t_1} \|(I - P_{U_1^*})(K(t_1)V_0^\top - Y_0)\| dt + c_1 h^2 + h\varepsilon + (hL + 1)\delta + Ch \\ 282 \quad = c_1 h^2 + h\varepsilon + (hL + 1)\delta + Ch$$

283 where in the last inequality we use the fact that  $f(Y_0)V_0 = \dot{K}(t_0)$ . The analogous  
284 derivation for the co-range proves the lemma.  $\square$

285 Lastly, we have

286 LEMMA 4.5. *It holds that*

$$287 \quad \|P_{U_1^*} W(t_1) P_{V_1^*} - Y_1^*\| \leq 2LBh^2 + Ch + (Lh + 1)\delta.$$

288 *Proof.* Let  $\bar{Y}(t) := U_1^* S(t) V_1^{*\top}$  for  $S(t)$  being the solution of the S-step of the  
 289 algorithm. From the construction of  $U_1^*$  and  $V_1^*$  it follows that  $\bar{Y}(t_0) = Y_0$ . From the  
 290 construction of the S-step, we know that

$$291 \quad Y_1^* = Y_0 + \int_{t_0}^{t_1} P_{U_1^*} f(\bar{Y}(t)) P_{V_1^*} dt.$$

292 Furthermore,

$$293 \quad P_{U_1^*} W_0 P_{V_1^*} - Y_0 = P_{U_1^*} (W_0 - Y_0) P_{V_1^*}.$$

294 With the above and Assumptions 4.2 and 4.3, we have

$$\begin{aligned} 295 \quad \|P_{U_1^*} W(t_1) P_{V_1^*} - Y_1^*\| &\leq \int_{t_0}^{t_1} \|P_{U_1^*} F(W(t)) P_{V_1^*} - P_{U_1^*} f(\bar{Y}(t)) P_{V_1^*}\| dt + \delta \\ 296 \quad &\leq \int_{t_0}^{t_1} \|P_{U_1^*} f(W(t)) P_{V_1^*} - P_{U_1^*} f(\bar{Y}(t)) P_{V_1^*}\| dt + Ch + \delta \\ 297 \quad &\leq L \int_{t_0}^{t_1} \|W(t) - \bar{Y}(t)\| dt + Ch + \delta \\ 298 \quad &\leq L \int_{t_0}^{t_1} \int_{t_0}^t \|F(W(s)) - P_{U_1^*} f(\bar{Y}(s)) P_{V_1^*}\| ds dt + (Lh + 1)\delta + Ch \\ 299 \quad &\leq 2LBh^2 + Ch + (Lh + 1)\delta. \quad \square \end{aligned}$$

300 Together, these lemmas yield the bound

$$301 \quad \|W(t_1) - Y_1\| \leq (C_1 + 2LB)h^2 + 2(L\delta + \varepsilon + 2C)h + (3 + Lh)\delta + \vartheta,$$

302 i.e., a global error bound then directly follows from Lady Windermere's fan.

303 This result shows that the rank-adaptive BUG algorithm does not lose its ro-  
 304 bustness properties following a stochastic modification. Thus, we can now shift to  
 305 exploring the properties of the Stochastic Algorithm 3.1 as a stochastic optimization  
 306 algorithm for machine learning.

307 **5. DLRA with Stochastic Gradients.** Recall the definition of the Stochastic  
 308 Algorithm 3.1 from the previous section. In this section, we explore its properties as an  
 309 optimization algorithm for machine learning. At first, we quantify how it optimizes the  
 310 loss function  $\mathcal{L}$ . Afterward, we show its convergence properties on the task of training  
 311 individual layers of deep neural networks when using stochastic gradient descent and  
 312 momentum methods as solvers of the integration steps. Let Assumptions 4.1 and 4.2  
 313 hold throughout this section.

314 **5.1. Optimization of the Loss Function.** Let  $Y_k$  be iterates generated by  
 315 the Rank-Adaptive Stochastic Algorithm 3.1 for some starting point  $Y_0$ . We begin by  
 316 taking a look at the structure of  $P_{U_{k+1}^*} f(Y_k) P_{V_{k+1}^*}$ . Using (3.7), we obtain

$$317 \quad P_{U_{k+1}^*} = U_{k+1}^* U_{k+1}^{*\top} = U_k U_k^\top + U_k^+ U_k^{+\top} = P_{U_k} + P_{U_k^+}.$$

318 The same also holds for  $V_{k+1}^*$ . Thus, we can write

$$\begin{aligned} 319 \quad &P_{U_{k+1}^*} f(Y_k) P_{V_{k+1}^*} = \\ 320 \quad (5.1) \quad &P_{U_k} f(Y_k) P_{V_k} + P_{U_k^+} f(Y_k) P_{V_k} + P_{U_k} f(Y_k) P_{V_k^+} + P_{U_k^+} f(Y_k) P_{V_k^+}. \end{aligned}$$

321 This form will help us to prove the next theorem.

322 **THEOREM 5.1.** *If  $Y_1$  is the resulting matrix after one pass of the Rank-Adaptive*  
 323 *Stochastic Algorithm 3.1, with  $t_1 = t_0 + h$  for  $h > 0$  and some starting point  $Y_0 \in \mathcal{M}_{q_0}$ ,*  
 324 *then there exists a positive constant  $\alpha$  such that*

$$325 \quad \mathbb{E}[\mathcal{L}(Y_1)] \leq \mathcal{L}(Y_0) - h\alpha^2 + c_1 h^2 B^2 + c_2 h B \varepsilon + c_3 h^2 B^2 + B\vartheta,$$

326 where  $c_1, c_2$ , and  $c_3$  are independent of the low-rank manifold's curvature.

327 *Proof.* Consider  $\bar{Y}(t) := U_1^* S(t) V_1^{*\top}$ , where  $S(t)$  denotes the solution of the S-step  
 328 of the Stochastic Algorithm 3.1. Then

$$329 \quad \dot{\bar{Y}}(t) = U_1^* \dot{S}(t) V_1^{*\top} = U_1^* U_1^{*\top} f(\bar{Y}(t)) V_1^* V_1^{*\top}.$$

330 By (3.7), the ranges of  $U_1^*$  and  $V_1^*$  contain those of  $U_0$  and  $V_0$ . Hence, we have

$$331 \quad \bar{Y}(t_0) = U_1^* S(t_0) V_1^{*\top} = U_1^* U_1^{*\top} U_0 S_0 V_0^\top V_1^* V_1^{*\top} = U_0 S_0 V_0^\top = Y_0.$$

332 As a consequence,

$$333 \quad (5.2) \quad \bar{Y}(t_1) = Y_0 + \int_{t_0}^{t_1} P_{U_1^*} f(\bar{Y}(s)) P_{V_1^*} ds.$$

334 Using equations (5.1) and (5.2), we obtain

$$\begin{aligned} 335 \quad \frac{d}{dt} \mathcal{L}(\bar{Y}(t)) &= \langle \nabla \mathcal{L}(\bar{Y}(t)), \dot{\bar{Y}}(t) \rangle = -\langle F(\bar{Y}(t)), P_{U_1^*} f(\bar{Y}(t)) P_{V_1^*} \rangle \\ 336 &\leq -\langle F(Y_0), P_{U_1^*} f(Y_0) P_{V_1^*} \rangle + c_1 h B^2 \\ 337 \quad (5.3) &= -\langle F(Y_0), P_{U_1^*} f(Y_0) P_{V_0} + P_{U_0} f(Y_0) P_{V_1^*} - P_{U_0} f(Y_0) P_{V_0} + P_{U_0^+} f(Y_0) P_{V_0^+} \rangle \\ 338 &\quad + c_1 h B^2. \end{aligned}$$

339 The last term in the scalar product can be bounded by

$$340 \quad |\langle F(Y_0), P_{U_0^+} f(Y_0) P_{V_0^+} \rangle| \leq |\langle F(Y_0), P_{U_0^+} P(Y_0) f(Y_0) P_{V_0^+} \rangle| + c_2 B \varepsilon = c_2 B \varepsilon.$$

341 We use the structure of the K- and L-steps to bound the remaining terms in (5.3).

342 The ranges of both  $K(t_1)$  and  $K_0$  are spanned by the columns of  $U_1^*$ . Therefore,

343  $(I - P_{U_1^*})K(t_1) = (I - P_{U_1^*})K_0 = 0$  holds. Combining this equality with

$$344 \quad K(t_1) = K_0 + \int_{t_0}^{t_1} f(K(s) V_0^\top) V_0 ds$$

345 yields

$$346 \quad (5.4) \quad \int_{t_0}^{t_1} P_{U_1^*} f(K(s) V_0^\top) V_0 ds = \int_{t_0}^{t_1} f(K(s) V_0^\top) V_0 ds.$$

347 A symmetric statement also holds for the L-step yielding

$$348 \quad \|P_{U_1^*} f(Y_0) P_{V_0} - f(Y_0) P_{V_0}\| \leq \frac{1}{2} c_3 h B,$$

$$349 \quad \|P_{V_1^*} f(Y_0)^\top P_{U_0} - f(Y_0)^\top P_{U_0}\| \leq \frac{1}{2} c_3 h B.$$

350 Hence, (5.3) becomes with  $\mu := c_1 h B^2 + c_2 B \varepsilon + c_3 h B^2$

$$351 \quad \frac{d}{dt} \mathcal{L}(\bar{Y}(t)) \leq -\langle F(Y_0), f(Y_0)P_{V_0} + P_{U_0}f(Y_0) - P_{U_0}f(Y_0)P_{V_0} \rangle + \mu.$$

352 Thus it holds that

$$\begin{aligned} 353 \quad \mathbb{E} \left[ \frac{d}{dt} \mathcal{L}(\bar{Y}(t)) \right] &\leq -\langle F(Y_0), F(Y_0)P_{V_0} + P_{U_0}F(Y_0) - P_{U_0}F(Y_0)P_{V_0} \rangle + \mu \\ 354 &= -\|U_0^\top F(Y_0)\|^2 - \|F(Y_0)V_0\|^2 + \|U_0^\top F(Y_0)V_0\|^2 + \mu \\ 355 &\leq -\frac{1}{2}(\|U_0^\top F(Y_0)\|^2 + \|F(Y_0)V_0\|^2) + \mu. \end{aligned}$$

356 Now let  $\alpha^2 := \frac{1}{2}(\|U_0^\top F(Y_0)\|^2 + \|F(Y_0)V_0\|^2)$ . Then we get

$$357 \quad \mathbb{E} \left[ \frac{d}{dt} \mathcal{L}(\bar{Y}(t)) \right] \leq -\alpha^2 + \mu.$$

358 Integrating this equation and using Fubini's theorem (which applies since the inte-  
359 grand  $\frac{d}{dt} \mathcal{L}(\bar{Y}(t))$  is bounded) yields

$$\begin{aligned} 360 \quad \mathbb{E}[\mathcal{L}(\bar{Y}(t_1))] &= \mathbb{E}[\mathcal{L}(\bar{Y}(t_0))] + \mathbb{E} \left[ \int_{t_0}^{t_1} \frac{d}{dt} \mathcal{L}(\bar{Y}(t)) dt \right] \\ 361 \quad (5.5) \quad &= \mathcal{L}(Y_0) + \int_{t_0}^{t_1} \mathbb{E} \left[ \frac{d}{dt} \mathcal{L}(\bar{Y}(t)) \right] dt \leq \mathcal{L}(Y_0) - h\alpha^2 + h\mu. \end{aligned}$$

362 Since  $S(t_1) = S_1^*$ , we can write  $\bar{Y}(t_1) = U_1^* S_1^* V_1^{*\top}$ . Consequently, by (3.8),

$$363 \quad (5.6) \quad \|Y_1 - \bar{Y}(t_1)\| \leq \vartheta.$$

364 By Taylor, there exists a  $\tau \in [0, 1]$ , such that

$$365 \quad \mathcal{L}(Y_1) = \mathcal{L}(\bar{Y}(t_1)) - \langle F(\tau Y_1 + (1-\tau)\bar{Y}(t_1)), Y_1 - \bar{Y}(t_1) \rangle.$$

366 Applying the Cauchy-Schwarz inequality and (5.6) to the above we obtain

$$367 \quad (5.7) \quad \mathcal{L}(Y_1) \leq \mathcal{L}(\bar{Y}(t_1)) + B\vartheta.$$

368 Putting (5.5) and (5.7) together and taking the expected value yields

$$369 \quad \mathbb{E}[\mathcal{L}(Y_1)] \leq \mathbb{E}[\mathcal{L}(\bar{Y}(t_1))] + B\vartheta \leq \mathcal{L}(Y_0) - h\alpha^2 + h\mu + B\vartheta. \quad \square$$

370 This theorem expresses an upper bound on the expected value of the loss function  
371  $\mathcal{L}$  after one step of the Stochastic Algorithm 3.1 in terms of the free parameters  $h$  and  
372  $\vartheta$ . By construction, it holds that  $\alpha \leq B$ , so this bound might actually be larger than  
373  $\mathcal{L}(Y_0)$ . This is somewhat unsatisfactory in the context of minimization. We will see  
374 later, however, that practical modifications of this algorithm are nonetheless capable  
375 of assuring descent, given an appropriate choice of the step size.

376 **5.2. Stochastic Gradient Descent.** Until now, we have assumed that the  
377 K-, L- and S-steps of Algorithm 3.1 are solved exactly. In practice, however, they  
378 are solved using various discrete-time methods, such as (stochastic) gradient descent  
379 and Adam [21, 3]. In this section, we will investigate the convergence properties of

380 the Stochastic Algorithm 3.1, where the continuous time variable  $t$  is replaced by  
 381 a discrete one ( $t = 0, 1, \dots$ ) and the integration steps are replaced by the gradient  
 382 descent algorithm, generally expressed as

$$383 \quad Z_{t+1} = Z_t - h\nabla\mathcal{G}(Z_t)$$

384 for a step size  $h > 0$ , a smooth function  $\mathcal{G}$ , and iterates  $Z_t$ . More precisely, the  
 385 integration in the K-step is replaced by

$$386 \quad (5.8) \quad K_{t+1} = U_t S_t + hf(U_t S_t V_t^\top) V_t,$$

387 in the L-step by

$$388 \quad L_{t+1} = V_t S_t^\top + hf(U_t S_t V_t^\top)^\top U_t,$$

389 and in the S-step by

$$390 \quad (5.9) \quad S_{t+1}^* = U_{t+1}^{*\top} U_t S_t V_t^\top V_{t+1}^* + h U_{t+1}^{*\top} f(P_{U_{t+1}^*} U_t S_t V_t^\top P_{V_{t+1}^*}) V_{t+1}^*.$$

391 Once again, a single realization of  $\xi$  is used to compute  $f$  for all these steps within  
 392 one pass of the `for`-loop of the algorithm. In this section, we will refer to this mod-  
 393 ification of the Stochastic Algorithm 3.1 as the *Stochastic Gradient Descent (SGD)*  
 394 Algorithm 3.1. This modification can also be seen as solving the differential equations  
 395 in the K-, L-, and S-steps with the explicit Euler method.

396 The logic in this section generally follows from adapting the treatment of SGD  
 397 found in Chapter 4 of [3] to the DLRA setting.

398 LEMMA 5.2. *For any  $Y, \bar{Y}$  it holds that*

$$399 \quad (5.10) \quad \mathcal{L}(Y) \leq \mathcal{L}(\bar{Y}) - \langle F(\bar{Y}), Y - \bar{Y} \rangle + \frac{L}{2} \|Y - \bar{Y}\|^2.$$

*Proof.*

$$\begin{aligned} 400 \quad \mathcal{L}(Y) &= \mathcal{L}(\bar{Y}) + \int_0^1 \frac{d}{dt} \mathcal{L}(\bar{Y} + t(Y - \bar{Y})) dt = \mathcal{L}(\bar{Y}) - \int_0^1 \langle F(\bar{Y} + t(Y - \bar{Y})), Y - \bar{Y} \rangle dt \\ 401 \quad &= \mathcal{L}(\bar{Y}) - \langle F(\bar{Y}), Y - \bar{Y} \rangle - \int_0^1 \langle F(\bar{Y} + t(Y - \bar{Y})) - F(\bar{Y}), Y - \bar{Y} \rangle dt \\ 402 \quad &\leq \mathcal{L}(\bar{Y}) - \langle F(\bar{Y}), Y - \bar{Y} \rangle + \int_0^1 L t \|Y - \bar{Y}\|^2 dt \\ 403 \quad &= \mathcal{L}(\bar{Y}) - \langle F(\bar{Y}), Y - \bar{Y} \rangle + \frac{L}{2} \|Y - \bar{Y}\|^2, \end{aligned}$$

404 where in the first line we use the chain rule and  $F = -\nabla\mathcal{L}$  and in the third line the  
 405 Cauchy-Schwarz inequality.  $\square$

406 Now, let  $Y_0 = U_0 S_0 V_0^\top \in \mathcal{M}_{q_0}$  be some fixed initial value for the SGD Algo-  
 407 rithm 3.1. Let  $Y_k$  be iterates generated by this algorithm.

408 THEOREM 5.3. *If  $Y_1$  is the resulting matrix after one pass of the Rank-Adaptive*  
 409 *SGD Algorithm 3.1, then it holds that*

$$410 \quad \mathbb{E}[\mathcal{L}(Y_1)] \leq \mathcal{L}(Y_0) - h \|P(Y_0) F(Y_0)\|^2 + \frac{1}{2} h^2 L B^2 + h B \varepsilon + B \vartheta.$$

411 *Proof.* Consider  $Y_1^* = U_1^* S_1^* V_1^{*\top}$ . By the definition of the S-step (5.9), we can  
 412 write

$$413 \quad S_1^* = U_1^{*\top} U_0 S_0 V_0^\top V_1^* + h U_1^{*\top} f(P_{U_1^*} U_0 S_0 V_0^\top P_{V_1^*}) V_1^*.$$

414 Using (3.7), we obtain

$$415 \quad U_1^* U_1^{*\top} U_0 S_0 V_0^\top V_1^* V_1^{*\top} = P_{U_1^*} U_0 S_0 V_0^\top P_{V_1^*} = U_0 S_0 V_0^\top = Y_0.$$

416 Combining the above equations yields

$$417 \quad Y_1^* = Y_0 + h P_{U_1^*} f(Y_0) P_{V_1^*}.$$

418 Using Lemma 5.2 and equation (5.1) we obtain

$$\begin{aligned} 419 \quad \mathcal{L}(Y_1^*) - \mathcal{L}(Y_0) &\leq -h \langle F(Y_0), P_{U_1^*} f(Y_0) P_{V_1^*} \rangle + \frac{h^2 L}{2} \|P_{U_1^*} f(Y_0) P_{V_1^*}\|^2 \\ 420 \quad &= -h \langle F(Y_0), P_{U_0} f(Y_0) P_{V_0} + P_{U_0^+} f(Y_0) P_{V_0} + P_{U_0} f(Y_0) P_{V_0^+} \rangle \\ 421 \quad (5.11) \quad &\quad -h \langle F(Y_0), P_{U_0^+} f(Y_0) P_{V_0^+} \rangle + \frac{h^2 L}{2} \|P_{U_1^*} f(Y_0) P_{V_1^*}\|^2. \end{aligned}$$

422 By Assumption 4.2 we have

$$423 \quad (5.12) \quad \frac{h^2 L}{2} \|P_{U_1^*} f(Y_0) P_{V_1^*}\|^2 \leq \frac{1}{2} h^2 L B^2.$$

424 Furthermore, by applying the logic preceding (5.4) to the new K-step (5.8), we obtain

$$425 \quad P_{U_1^*} f(Y_0) P_{V_0} = f(Y_0) P_{V_0}.$$

426 Combining this with a symmetric argument on the L-step yields

$$\begin{aligned} 427 \quad P_{U_0} f(Y_0) P_{V_0} + P_{U_0^+} f(Y_0) P_{V_0} + P_{U_0} f(Y_0) P_{V_0^+} \\ 428 \quad = P_{U_1^*} f(Y_0) P_{V_0} + P_{U_0} f(Y_0) P_{V_1^*} - P_{U_0} f(Y_0) P_{V_0} \\ 429 \quad (5.13) \quad = f(Y_0) P_{V_0} + P_{U_0} f(Y_0) - P_{U_0} f(Y_0) P_{V_0} = P(Y_0) f(Y_0). \end{aligned}$$

430 The first term in (5.11),  $P_{U_0^+} f(Y_0) P_{V_0^+}$ , can be bounded using Assumption 4.1. Since  
 431 we know that  $f(Y_0) = m(Y_0) + r(Y_0)$  with  $m(Y_0) \in \mathcal{T}_{Y_0} \mathcal{M}_{q_0}$ , the definition of  $U_0^+$  and  
 432  $V_0^+$  and the Assumption 4.1 allow us to write

$$433 \quad (5.14) \quad \|P_{U_0^+} f(Y_0) P_{V_0^+}\| = \|P_{U_0^+} r(Y_0) P_{V_0^+}\| \leq \varepsilon.$$

434 Applying (5.12), (5.13), and (5.14) onto (5.11) and using the Cauchy-Schwarz  
 435 inequality as well as Assumption 4.2 yields

$$436 \quad \mathcal{L}(Y_1^*) - \mathcal{L}(Y_0) \leq \langle F(Y_0), P(Y_0) f(Y_0) \rangle + \frac{1}{2} h^2 L B^2 + h B \varepsilon.$$

437 Taking the expected value, this becomes

$$438 \quad \mathbb{E}[\mathcal{L}(Y_1^*)] - \mathcal{L}(Y_0) \leq -h \|P(Y_0) F(Y_0)\|^2 + \frac{1}{2} h^2 L B^2 + h B \varepsilon.$$

439 Furthermore, by (3.8), we know that

$$440 \quad \|Y_1 - Y_1^*\| \leq \vartheta.$$

441 Thus, using Taylor, for some  $\tau \in [0, 1]$ , we get

$$442 \quad \mathcal{L}(Y_1) = \mathcal{L}(Y_1^*) - \langle F(\tau Y_1 + (1 - \tau)Y_1^*), Y_1 - Y_1^* \rangle \leq \mathcal{L}(Y_1^*) + B\vartheta.$$

443 Putting everything together and taking the expected value yields

$$444 \quad \mathbb{E}[\mathcal{L}(Y_1)] \leq \mathbb{E}[\mathcal{L}(Y_1^*)] + B\vartheta \leq \mathcal{L}(Y_0) - h\|P(Y_0)F(Y_0)\|^2 + \frac{1}{2}h^2LB^2 + hB\varepsilon + B\vartheta. \quad \square$$

445 Whenever we write  $f(Y) = f(Y, \xi)$ , the randomness is hidden in the term  $\xi$ .  
 446 Thus, when reading  $\mathbb{E}[f(Y_0)]$ , one should understand  $\mathbb{E}_\xi[f(Y_0, \xi)]$ . The matrix  $Y_1$   
 447 and its factors  $U_1$ ,  $V_1$ , and  $S_1$  are stochastic in  $\xi$ , since they have been generated  
 448 using  $f(Y_0, \xi)$ . Therefore,  $\mathbb{E}[\mathcal{L}(Y_1)] = \mathbb{E}_\xi[\mathcal{L}(Y_1(\xi))]$  holds. In general, when consid-  
 449 ering the evolution of  $Y_t$  up to the  $k$ -th step, if we denote with  $\xi_1, \dots, \xi_k$  the i.i.d.  
 450 realizations of  $\xi$  made in each pass of the algorithm, we can say that  $Y_t$ ,  $U_t$ ,  $V_t$ ,  
 451 and  $S_t$  are stochastic in  $\xi_1, \dots, \xi_t$ , so the simplified notation  $\mathbb{E}[f(Y_t)]$  corresponds to  
 452  $\mathbb{E}_{\xi_1, \dots, \xi_{t+1}}[f(Y_t(\xi_1, \dots, \xi_t), \xi_{t+1})]$ .

453 Thus, we can rewrite the result of Theorem 5.3 in a more general manner:

$$454 \quad (5.15) \quad \mathbb{E}_{\xi_t}[\mathcal{L}(Y_t(\xi_t))] \leq \mathcal{L}(Y_{t-1}) - h\|P(Y_{t-1})F(Y_{t-1})\|^2 + \frac{1}{2}h^2LB^2 + hB\varepsilon + B\vartheta.$$

455 Now, we can make statements about the behavior of the algorithm as  $t \rightarrow \infty$ .

456 **THEOREM 5.4.** *Let  $\mathcal{L}$  be non-negative. Let  $Y_1, \dots, Y_k$  be iterates generated by the*  
 457 *Rank-Adaptive SGD Algorithm 3.1 over  $k$  steps. Then it holds that*

$$458 \quad \frac{1}{k} \sum_{t=1}^k \mathbb{E}[\|P(Y_{t-1})F(Y_{t-1})\|^2] \leq$$

$$459 \quad (5.16) \quad \frac{\mathcal{L}(Y_0)}{kh} + \frac{1}{2}hLB^2 + B\varepsilon + \frac{1}{h}B\vartheta \xrightarrow{k \rightarrow \infty} \frac{1}{2}hLB^2 + B\varepsilon + \frac{1}{h}B\vartheta,$$

460 where the expected value is taken over all  $\xi_t$ .

461 *Proof.* By taking the expected value over all  $\xi_t$  in (5.15), we get

$$462 \quad \mathbb{E}[\mathcal{L}(Y_t)] - \mathbb{E}[\mathcal{L}(Y_{t-1})] \leq -h\mathbb{E}[\|P(Y_{t-1})F(Y_{t-1})\|^2] + \frac{1}{2}h^2LB^2 + hB\varepsilon + B\vartheta.$$

463 Using  $\mathcal{L} \geq 0$ , we can now conclude

$$464 \quad -\mathcal{L}(Y_0) \leq \mathbb{E}[\mathcal{L}(Y_k)] - \mathcal{L}(Y_0)$$

$$465 \quad \leq -h \sum_{t=1}^k \mathbb{E}[\|P(Y_{t-1})F(Y_{t-1})\|^2] + k \left( \frac{1}{2}h^2LB^2 + hB\varepsilon + B\vartheta \right).$$

466 Rearranging the terms, we obtain

$$467 \quad \sum_{t=1}^k \mathbb{E}[\|P(Y_{t-1})F(Y_{t-1})\|^2] \leq \frac{1}{h} \left( \mathcal{L}(Y_0) + k \left( \frac{1}{2}h^2LB^2 + hB\varepsilon + B\vartheta \right) \right)$$

$$468 \quad = \frac{\mathcal{L}(Y_0)}{h} + k \left( \frac{1}{2}hLB^2 + B\varepsilon + \frac{1}{h}B\vartheta \right).$$

469 Dividing by  $k$  and taking the limit yields the desired result. □

470 Theorem 5.4 states that the running average of the expected squared norms of  
 471 the projected gradients of  $\mathcal{L}$  does not surpass  $\frac{1}{2}hLB^2 + B\varepsilon + \frac{1}{h}B\vartheta$ . If this bound could  
 472 be made arbitrarily small, we could claim convergence towards a stationary point. In  
 473 the setting of the usual stochastic gradient descent algorithm such as in Section 4 of  
 474 [3], this bound is linear in  $h$ . This can be used to obtain

$$475 \quad \liminf_{t \rightarrow \infty} \mathbb{E}[\|P(Y_t)F(Y_t)\|^2] = 0,$$

476 by using a variable and shrinking step size  $h = h_t$ .

477 Two terms are preventing us from taking this approach here. Firstly,  $B\varepsilon$  is con-  
 478 stant in  $h$ , and reducing the step size would not affect it. This term stems from  
 479 the bound (5.14) we use on terms of the form  $\langle F(Y_t), P_{U_t^+} f(Y_t) P_{V_t^+} \rangle$ , which we can-  
 480 not easily integrate because both the projections  $P_{U_t^+}$  and the gradients  $f(Y_t)$  are  
 481 stochastic in  $\xi_{t+1}$ .

482 The other even more problematic term in (5.16) is  $\frac{1}{h}B\vartheta$ , as it is inversely pro-  
 483 portional to  $h$ . This term appears because the truncation step is independent of the  
 484 learning rate  $h$  and the stochastic gradient  $f$ . All the steps preceding the truncation  
 485 move the objective towards a (stochastic) decrease, while the truncation can seem-  
 486 ingly throw it off in any direction at step distance  $\vartheta$ . If the learning rate  $h$  is made  
 487 smaller, this truncation displacement becomes increasingly dominant in the progress  
 488 of the algorithm.

489 We generally cannot expect  $\vartheta$  to be smaller than  $h$ , see, e.g., the values in Section 5  
 490 of [21]. Letting  $\vartheta \rightarrow 0$  would defeat the purpose of rank reduction since, over many  
 491 algorithm passes, this might yield matrices of high or even full rank. Notice that at  
 492 every pass of Rank-Adaptive Algorithm 3.1, the rank is initially increased by up to  
 493 two times with respect to the rank in the previous pass, and only those dimensions  
 494 that have singular values below the threshold are later removed; if the threshold goes  
 495 to 0, fewer dimensions are removed at each pass.

496 Interestingly, rank-adaptivity seems to be less impactful in the later stages of  
 497 training. The experiments in [21] indicate that the ranks of the parameter matrices  
 498 become close to constant after sufficiently many training steps. Thus, a sensible  
 499 solution would be to perform rank-adaptivity only at the beginning of model training  
 500 and afterward continue in a non-adaptive manner with a shrinking  $h_t$ . This being  
 501 said, such an approach makes it challenging to obtain a theoretical result akin to the  
 502 ones above since the proofs depend on the property that the ranges of  $U_t$  and  $V_t$  are  
 503 contained in those of  $U_{t+1}^*$  and  $V_{t+1}^*$ . This property does not generally hold in the  
 504 non-adaptive setting.

505 We avoid this issue by exploring a method called *S-fine-tuning*. This method has  
 506 been implemented in the source code [22] released alongside [21], which consists of  
 507 dropping the K- and L-steps and only performing S-steps. S-fine-tuning is performed  
 508 after the model has been trained for several epochs (see the `train_and_finetune`  
 509 function in the `DLRT-Net/optimizer_KLS/train_experiments.py` scrip of the source  
 510 code). It seems to rely on the assumption that, in the later stages of model training,  
 511 the computational costs related to calculating the K- and L-steps outweigh the gain  
 512 in precision obtained from updating the  $U$  and  $V$  matrices.

513 Performing S-fine-tuning can be seen as assuming that the ranges of the solutions  
 514  $K_t$  of the K-steps remain constant. In other words, we get

$$515 \quad (I - P_U)K_t = 0$$

516 for all  $t$ . An analogous interpretation also holds for the L-step. Just as in (5.13),  
 517 applying this to the definition of the K- and L-steps yields

$$518 \quad P(Y_t)f(Y_t) = P_U f(Y_t)P_V.$$

519 We can now state our assumption.

520 **ASSUMPTION 5.5.** *There exists an index  $t_0$  such that for all  $t \geq t_0$*

- 521 1.  $U_t = U_{t-1}$  and  $V_t = V_{t-1}$ ,
- 522 2.  $P(Y_t)f(Y_t) = P_{U_t}f(Y_t)P_{V_t}$ .

523 *Remark 5.6.* When using this assumption, we drop the indices in our notation  
 524 of the  $U$  and  $V$  matrices. This assumption implies that, after the index  $t_0$ , we not  
 525 only use the non-adaptive method, which makes the term  $\frac{1}{h}B\vartheta$  disappear, but also  
 526 no longer perform the K- and L-steps. It follows immediately that  $U_t^+ = V_t^+ = 0$  for  
 527 all  $t \geq t_0$ , so the term  $B\varepsilon$  must also disappear by (5.1) and the discussion above.

528 **THEOREM 5.7.** *In the setting of Theorem 5.3, let Assumption 5.5 hold. Then*

$$529 \quad (5.17) \quad \mathbb{E}_{\xi_t}[\mathcal{L}(Y_t)] \leq \mathcal{L}(Y_{t-1}) - h\|P(Y_{t-1})F(Y_{t-1})\|^2 + \frac{1}{2}h^2LB^2$$

530 for all  $t \geq t_0$  for  $t_0$  from Assumption 5.5.

531 *Proof.* For  $t \geq t_0$ , by the definition of the S-step (5.9),

$$532 \quad S_t = U^\top U S_{t-1} V^\top V + hU^\top f(P_U U S_{t-1} V^\top P_V) V = S_{t-1} + hU^\top f(Y_{t-1}) V,$$

533 so

$$534 \quad Y_t = U S_t V^\top = Y_{t-1} + hP_U f(Y_{t-1}) P_V = Y_{t-1} + hP(Y_{t-1}) f(Y_{t-1}).$$

535 Using Lemma 5.2 and Assumption 4.2, we can write

$$536 \quad \begin{aligned} \mathcal{L}(Y_t) - \mathcal{L}(Y_{t-1}) &\leq -h\langle F(Y_{t-1}), P(Y_{t-1})f(Y_{t-1}) \rangle + \frac{h^2L}{2}\|P(Y_{t-1})f(Y_{t-1})\|^2 \\ 537 &\leq -h\langle F(Y_{t-1}), P(Y_{t-1})f(Y_{t-1}) \rangle + \frac{1}{2}h^2LB^2. \end{aligned}$$

538 Taking the expected value with respect to  $\xi_t$  yields the desired result.  $\square$

539 *Remark 5.8.* In particular, (5.17) shows that, for a small enough  $h$ , the algorithm  
 540 assures a decrease of the loss function in expectation.

541 The proof of Theorem 5.4 where (5.15) is replaced by (5.17) immediately yields,  
 542 up to a shift in the indices such that  $t_0 = 0$ ,

$$543 \quad (5.18) \quad \frac{1}{k} \sum_{t=1}^k \mathbb{E}[\|P(Y_{t-1})F(Y_{t-1})\|^2] \leq \frac{\mathcal{L}(Y_0)}{kh} + \frac{1}{2}hLB^2 \xrightarrow{k \rightarrow \infty} \frac{1}{2}hLB^2.$$

544 We can see that Assumption 5.5 has removed both problematic terms in equation  
 545 (5.16). It is clear that the right-hand side of (5.18) goes to 0 if we let  $h \rightarrow 0$ . To  
 546 achieve this, let us now choose a variable step size  $h = h_t$  that satisfies

$$547 \quad (5.19) \quad \sum_{t=0}^{\infty} h_t = \infty \quad \text{and} \quad \sum_{t=0}^{\infty} h_t^2 < \infty.$$

548 THEOREM 5.9. *In the setting of Theorem 5.4, let Assumption 5.5 hold and a*  
 549 *variable step size  $h_t$  satisfy (5.19). Then one has*

550 (5.20) 
$$\sum_{t=1}^{\infty} h_{t-1} \mathbb{E}[\|P(Y_{t-1})F(Y_{t-1})\|^2] < \infty.$$

551 *Proof.* Without loss of generality let  $t_0 = 0$  in Assumption 5.5, otherwise set the  
 552 index after which this holds to 0.

553 Using (5.17),

554 
$$\mathbb{E}[\mathcal{L}(Y_t)] - \mathbb{E}[\mathcal{L}(Y_{t-1})] \leq -h_{t-1} \mathbb{E}[\|P(Y_{t-1})F(Y_{t-1})\|^2] + \frac{1}{2} h_{t-1}^2 LB^2.$$

555 Similarly to the proof of Theorem 5.4, we obtain

556 
$$-\mathcal{L}(Y_0) \leq \mathbb{E}[\mathcal{L}(Y_k)] - \mathcal{L}(Y_0) \leq -\frac{1}{2} \sum_{t=1}^k h_{t-1} \mathbb{E}[\|P(Y_{t-1})F(Y_{t-1})\|^2] + \frac{LB^2}{2} \sum_{t=1}^k h_{t-1}^2.$$

557 Rearranging the terms,

558 
$$\sum_{t=1}^k h_{t-1} \mathbb{E}[\|P(Y_{t-1})F(Y_{t-1})\|^2] \leq 2\mathcal{L}(Y_0) + LB^2 \sum_{t=1}^k h_{t-1}^2.$$

559 Taking the limit  $k \rightarrow \infty$  and using  $\sum_{t=0}^{\infty} h_t^2 < \infty$  yields the result. □

560 COROLLARY 5.10. *In the setting of Theorem 5.9, it holds that*

561 (5.21) 
$$\liminf_{t \rightarrow \infty} \mathbb{E}[\|P(Y_t)F(Y_t)\|^2] = 0.$$

562 *Proof.* The statement follows directly from (5.20) and (5.19). □

563 This shows that in the setting of Theorem 5.9, the SGD Algorithm 3.1 yields  
 564 a sequence of iterates, such that a subsequence comes arbitrarily close to satisfying  
 565 the necessary condition for optimality (3.9) in expectation. If Assumption 5.5.2 is  
 566 dropped, the same result holds, albeit for the weaker necessary condition of optimal-  
 567 ity (3.10).

568 This convergence result indicates that the SGD Algorithm 3.1 is a valid optimiza-  
 569 tion algorithm for training individual layers of DNNs. The discussion in this section  
 570 also suggests that the best way to apply it would be to first train the model in a  
 571 rank-adaptive manner and, once the ranks of the  $S$  matrices have stabilized, finish  
 572 the training using S-fine-tuning.

573 **5.3. Momentum Methods.** Although stochastic gradient descent has an im-  
 574 portant place among optimization methods for machine learning, in practice, it has  
 575 been largely outperformed by momentum methods [3, 25, 10]. These methods do  
 576 not simply use the gradient of the current step but rather the accumulated gradient  
 577 information from all previous steps [3, 25]. Some common examples are the heavy  
 578 ball and Nesterov methods. In this section, we will use a momentum method as a  
 579 solver of the differential equations that constitute the Stochastic Algorithm 3.1 and  
 580 will investigate the convergence properties of the resulting algorithm.

581 For  $f$  as in the previous sections, consider the two-step algorithm

582 (5.22) (SUM) : 
$$\begin{cases} X_t = \mu X_{t-1} + h_t f(Y_t), \\ Y_{t+1} = Y_t + \lambda h_t f(Y_t) + (1 - \lambda + \lambda \mu) X_t \end{cases}$$

583 with parameters  $\mu \in [0, 1)$ ,  $\lambda \in [0, \frac{1}{1-\mu}]$  and step sizes  $h_t$ , that generates a sequence  
 584 of iterates  $Y_t$  from starting values  $Y_0$  and  $X_0 := 0$ . This is the *Stochastic Unified*  
 585 *Momentum (SUM)* algorithm proposed in [18], which generalizes the stochastic heavy  
 586 ball ( $\lambda = 0$ ) and Nesterov ( $\lambda = 1$ ) methods. When applied as an integrator of the S-  
 587 step of the Stochastic Algorithm 3.1 once again replacing the continuous time variable  
 588  $t$  by a discrete one, it becomes

$$589 \quad (5.23) \quad \begin{aligned} X_t &= \mu X_{t-1} + h_t U_{t+1}^{*\top} f(P_{U_{t+1}^*} U_t S_t V_t^\top P_{V_{t+1}^*}) V_{t+1}^*, \\ S_{t+1}^* &= U_{t+1}^{*\top} U_t S_t V_t^\top V_{t+1}^* + \lambda h_t U_{t+1}^{*\top} f(P_{U_{t+1}^*} U_t S_t V_t^\top P_{V_{t+1}^*}) V_{t+1}^* \\ &\quad + (1 - \lambda + \lambda \mu) X_t. \end{aligned}$$

590 This algorithm is not immediately applicable in the rank-adaptive setting since the  
 591 dimension of  $X_{t-1}$  and that of the stochastic gradient

$$592 \quad U_{t+1}^{*\top} f(P_{U_{t+1}^*} U_t S_t V_t^\top P_{V_{t+1}^*}) V_{t+1}^*$$

593 are not necessarily the same at any given step.

594 In practice, this issue can be circumvented by using heuristics. For example, in the  
 595 source code of [21], when applying the Adam algorithm, which also uses momentum  
 596 and thus suffers from the same issue [10], in the rank-adaptive setting, the dimensions  
 597 of  $X_t$  and  $S_t$  are kept constant and set to the largest possible value, while only  
 598 *submatrices* of  $X_t$  and  $S_t$  with appropriate dynamical dimensions are being used and  
 599 updated.

600 Such approaches do not correspond to the SUM algorithm (5.22) that we want to  
 601 study. Luckily, the experimental findings in [21] suggest that even when using these  
 602 heuristics, the adaptive rank of the matrices  $Y_t$  stabilizes during training. It is thus  
 603 reasonable to once again use Assumption 5.5.1. Just like in the previous section, we  
 604 can generally assume that  $t_0 = 0$ , for the number of steps  $t_0$  after which the basis is  
 605 kept fixed since we are only interested in the behavior as  $t \rightarrow \infty$ .

606 Under this assumption, we can rewrite (5.23) as

$$607 \quad \begin{aligned} X_t &= \mu X_{t-1} + h_t U^\top f(Y_t) V, \\ 608 \quad S_{t+1} &= S_t + \lambda h_t U^\top f(Y_t) V + (1 - \lambda + \lambda \mu) X_t. \end{aligned}$$

609 Since there is no truncation step, multiplying by  $U$  from the left and by  $V^\top$  from the  
 610 right yields, with a new definition of  $X_t$ ,

$$611 \quad (5.24) \quad \begin{aligned} X_t &= \mu X_{t-1} + h_t P_U f(Y_t) P_V, \\ Y_{t+1} &= Y_t + \lambda h_t P_U f(Y_t) P_V + (1 - \lambda + \lambda \mu) X_t. \end{aligned}$$

612 Consider the modification of the non-adaptive Algorithm 3.1, where, as under As-  
 613 sumption 5.5.1, only the S-step is performed and is further replaced by (5.24). In this  
 614 section, we will refer to this modification as the *Stochastic Unified Momentum (SUM)*  
 615 Algorithm 3.1. We will now investigate its convergence properties.

616 Let Assumptions 4.2 and 5.5.1 hold throughout this section. Let  $\{Y_t\}_{t \geq 0}$  be a  
 617 sequence of iterates generated by the SUM Algorithm 3.1 for some  $Y_0 \in \mathcal{M}_q$  and  
 618  $\{X_t\}_{t \geq 0}$ ,  $X_0 = 0$ , be its corresponding sequence from (5.24).

619 The logic in this section generally follows from applying the treatment of the SUM  
 620 algorithm from [18] to the DLRA setting.

621 First, we need to state a few technical lemmas.

622 LEMMA 5.11. Let  $\{a_t\}_{t \geq 1}$ ,  $\{b_t\}_{t \geq 1}$ , and  $\{\tilde{a}_t\}_{t \geq 1}$  be non-negative real sequences  
 623 such that  $\sum_{t=1}^{\infty} a_t = \infty$ ,  $\sum_{t=1}^{\infty} a_t b_t^2 < \infty$ ,  $\lim_{t \rightarrow \infty} \frac{a_t}{\tilde{a}_t} = 1$ , and  $|b_{t+1} - b_t| \leq C \tilde{a}_t$  for a  
 624 positive constant  $C$ . Then  $\lim_{t \rightarrow \infty} b_t = 0$ .

625 *Proof.* This lemma is proven as Corollary 3.1 in [18].  $\square$

626 LEMMA 5.12. Let  $\{a_t\}_{t \geq 0}$ ,  $\{b_t\}_{t \geq 0}$ , and  $\{c_t\}_{t \geq 0}$  be real sequences where  $\{b_t\}_{t \geq 0}$   
 627 is non-negative. Further let  $a_{t+1} \leq a_t - b_t + c_t$  and  $\sum_{t=0}^{\infty} c_t$  converge. Then either  
 628  $\lim_{t \rightarrow \infty} a_t = -\infty$ , or  $a_t$  converges and  $\sum_{t=0}^{\infty} b_t < \infty$ .

629 *Proof.* This lemma is proven as Lemma 1 in [1].  $\square$

630 LEMMA 5.13. Let the step sizes  $h_t$  satisfy (5.19). Then it holds that

- 631 1.  $\sum_{t=1}^{\infty} \left( \sum_{s=1}^t \mu^{t-s} h_s^2 \right) < \infty$ ,  
 632 2.  $\sum_{k=1}^{\infty} \left( \sum_{t=1}^{k-1} \mu^{k-t} \mathbb{E}[\|X_t\|^2] \right) < \infty$ ,  
 633 3.  $\sum_{t=1}^{\infty} \mathbb{E}[\|Y_{t+1} - Y_t\|^2] < \infty$ .

634 *Proof.* We can rewrite  $X_t$  from (5.24) as

$$635 \quad (5.25) \quad X_t = \sum_{s=1}^t \mu^{t-s} h_s P_U f(Y_s) P_V.$$

636 Thus, using Jensen's inequality in the second line,

$$637 \quad \mathbb{E}[\|X_t\|^2] \leq \mathbb{E} \left[ \left( \sum_{s=1}^t \mu^{t-s} h_s \|P_U f(Y_s) P_V\| \right)^2 \right]$$

$$638 \quad \leq \mathbb{E} \left[ \left( \sum_{s=1}^t \mu^{t-s} \right) \sum_{s=1}^t \mu^{t-s} h_s^2 \|P_U f(Y_s) P_V\|^2 \right]$$

$$639 \quad \leq \frac{B^2}{1-\mu} \sum_{s=1}^t \mu^{t-s} h_s^2.$$

640 Since  $\sum_{t=1}^{\infty} h_t^2$  and  $\sum_{t=1}^{\infty} \mu^t$  converge absolutely, their Cauchy product

$$641 \quad \sum_{t=1}^{\infty} \left( \sum_{s=1}^t \mu^{t-s} h_s^2 \right)$$

642 also converges. Thus

$$643 \quad (5.26) \quad \sum_{t=1}^{\infty} \mathbb{E}[\|X_t\|^2] < \infty \quad \text{and} \quad \sum_{k=1}^{\infty} \left( \sum_{t=1}^{k-1} \mu^{k-t} \mathbb{E}[\|X_t\|^2] \right) < \infty,$$

644 where we notice that the sum in the second inequality is once again a Cauchy product  
 645 of two absolutely convergent series. Using the inequality  $(a+b)^2 \leq 2(a^2+b^2)$ , we can  
 646 also write

$$647 \quad \|Y_{t+1} - Y_t\|^2 = \|\lambda h_t P_U f(Y_t) P_V + (1-\lambda + \lambda\mu) X_t\|^2$$

$$648 \quad (5.27) \quad \leq 2\lambda^2 h_t^2 \|P_U f(Y_t) P_V\|^2 + 2(1-\lambda + \lambda\mu)^2 \|X_t\|^2$$

$$649 \quad \leq 2\lambda^2 h_t^2 B^2 + 2(1-\lambda + \lambda\mu)^2 \|X_t\|^2.$$

650 Thus,  $\sum_{t=1}^{\infty} h_t^2 < \infty$  and (5.26) yield

$$651 \quad \sum_{t=1}^{\infty} \mathbb{E}[\|Y_{t+1} - Y_t\|^2] < \infty. \quad \square$$

652 LEMMA 5.14. *It holds that*

$$653 \quad -\mathbb{E}[\langle F(Y_t), X_t \rangle] \leq -\sum_{s=1}^t \mu^{t-s} h_s \mathbb{E}[\|P_U F(Y_s) P_V\|^2]$$

$$654 \quad + 2L \sum_{s=1}^{t-1} \mu^{t-s} \mathbb{E}[\|X_s\|^2] + L\lambda^2 B^2 \sum_{s=1}^{t-1} \mu^{t-s} h_s^2.$$

655 *Proof.* By the definition of  $X_k$ , we have

$$656 \quad -\mathbb{E}[\langle F(Y_t), X_t \rangle] = -\mu \mathbb{E}[\langle F(Y_t), X_{t-1} \rangle] - h_t \mathbb{E}[\|P_U F(Y_t) P_V\|^2].$$

657 Using (5.27), the inequality  $ab \leq \frac{1}{2}(a^2 + b^2)$ , and the Cauchy-Schwarz inequality, we  
658 obtain

$$659 \quad \langle F(Y_{t-1}) - F(Y_t), X_{t-1} \rangle \leq \|F(Y_{t-1}) - F(Y_t)\| \|X_{t-1}\| \leq L \|Y_{t-1} - Y_t\| \|X_{t-1}\|$$

$$660 \quad \leq \frac{L}{2} \|Y_{t-1} - Y_t\|^2 + \frac{L}{2} \|X_{t-1}\|^2 \leq L\lambda^2 h_{t-1}^2 B^2 + \left( L(1 - \lambda + \lambda\mu)^2 + \frac{L}{2} \right) \|X_{t-1}\|^2$$

$$661 \quad \leq L\lambda^2 h_{t-1}^2 B^2 + 2L \|X_{t-1}\|^2,$$

662 where in the last line we use the fact that  $\mu - 1 < 0$  and thus  $1 - \lambda + \lambda\mu \leq 1$ .

663 Combining these findings yields

$$664 \quad -\mathbb{E}[\langle F(Y_t), X_t \rangle] = -\mu \mathbb{E}[\langle F(Y_t) + F(Y_{t-1}) - F(Y_{t-1}), X_{t-1} \rangle] - h_t \mathbb{E}[\|P_U F(Y_t) P_V\|^2]$$

$$665 \quad \leq -\mu \mathbb{E}[\langle F(Y_{t-1}), X_{t-1} \rangle] - h_t \mathbb{E}[\|P_U F(Y_t) P_V\|^2] + \mu L\lambda^2 h_{t-1}^2 B^2 + 2\mu L \mathbb{E}[\|X_{t-1}\|^2]$$

$$666 \quad \leq -\sum_{s=1}^t \mu^{t-s} h_s \mathbb{E}[\|P_U F(Y_s) P_V\|^2] + 2L \sum_{s=1}^{t-1} \mu^{t-s} \mathbb{E}[\|X_s\|^2] + L\lambda^2 B^2 \sum_{s=1}^{t-1} \mu^{t-s} h_s^2. \quad \square$$

667 We can now show the central result of this section.

668 THEOREM 5.15. *Let  $\mathcal{L}$  be non-negative and the step sizes  $h_t$  be such that (5.19)*  
669 *as well as  $\lim_{t \rightarrow \infty} \frac{h_{t-1}}{h_t} = 1$  hold. Then*

$$670 \quad (5.28) \quad \lim_{t \rightarrow \infty} \mathbb{E}[\|P_U F(Y_t) P_V\|] = 0.$$

671 *Proof.* Putting Lemmas 5.2 and 5.14 together, we can write

$$672 \quad \mathbb{E}[\mathcal{L}(Y_{t+1})] - \mathbb{E}[\mathcal{L}(Y_t)]$$

$$673 \quad \leq -\mathbb{E}[\langle F(Y_t), \lambda h_t P_U f(Y_t) P_V + (1 - \lambda + \lambda\mu) X_t \rangle] + \frac{L}{2} \mathbb{E}[\|Y_{t+1} - Y_t\|^2]$$

$$674 \quad (5.29) \quad \leq -\lambda h_t \mathbb{E}[\|P_U F(Y_t) P_V\|^2] + \frac{L}{2} \mathbb{E}[\|Y_{t+1} - Y_t\|^2]$$

$$675 \quad + (1 - \lambda + \lambda\mu) \left( 2L \sum_{s=1}^{t-1} \mu^{t-s} \mathbb{E}[\|X_s\|^2] + L\lambda^2 B^2 \sum_{s=1}^{t-1} \mu^{t-s} h_s^2 - \right.$$

$$676 \quad \left. \sum_{s=1}^t \mu^{t-s} h_s \mathbb{E}[\|P_U F(Y_s) P_V\|^2] \right).$$

677 Consider the sequences

$$678 \quad a_t := \mathbb{E}[\mathcal{L}(Y_t)],$$

$$679 \quad b_t := \lambda h_t \mathbb{E}[\|P_U F(Y_t) P_V\|^2] + (1 - \lambda + \lambda\mu) \sum_{s=1}^t \mu^{t-s} h_s \mathbb{E}[\|P_U F(Y_s) P_V\|^2],$$

$$680 \quad c_t := \frac{L}{2} \mathbb{E}[\|Y_{t+1} - Y_t\|^2] + (1 - \lambda + \lambda\mu) \left( 2L \sum_{s=1}^{t-1} \mu^{t-s} \mathbb{E}[\|X_s\|^2] + L\lambda^2 B^2 \sum_{s=1}^{t-1} \mu^{t-s} h_s^2 \right).$$

681 The  $b_t$  are non-negative by definition of  $\lambda$  and  $\mu$ . By (5.29) we know that  $a_{t+1} \leq$   
 682  $a_t - b_t + c_t$ . Also, from Lemma 5.13 it follows that  $\sum_{t=0}^{\infty} c_t$  converges, so Lemma 5.12  
 683 applies to these sequences. By the non-negativity of  $\mathcal{L}$ , this yields  $a_t \rightarrow a_* < \infty$  and

$$684 \quad (5.30) \quad \sum_{t=1}^{\infty} \left( \lambda h_t \mathbb{E}[\|P_U F(Y_t) P_V\|^2] + (1 - \lambda + \lambda\mu) \sum_{s=1}^t \mu^{t-s} h_s \mathbb{E}[\|P_U F(Y_s) P_V\|^2] \right)$$

$$685 \quad = \sum_{t=1}^{\infty} b_t < \infty.$$

686 In particular, since for any  $k$  it holds that

$$687 \quad \sum_{t=1}^k \sum_{s=1}^t \mu^{t-s} h_s \mathbb{E}[\|P_U F(Y_s) P_V\|^2] = \sum_{s=1}^k \left( h_s \mathbb{E}[\|P_U F(Y_s) P_V\|^2] \sum_{t=s}^k \mu^{t-s} \right)$$

$$688 \quad \geq \sum_{s=1}^k h_s \mathbb{E}[\|P_U F(Y_s) P_V\|^2],$$

689 we can conclude with (5.30) that

$$690 \quad (5.31) \quad \sum_{t=1}^{\infty} h_t \mathbb{E}[\|P_U F(Y_t) P_V\|^2] \leq \sum_{t=1}^{\infty} \sum_{s=1}^t \mu^{t-s} h_s \mathbb{E}[\|P_U F(Y_s) P_V\|^2] < \infty.$$

691 Now, using (5.25), we write

$$692 \quad \|Y_{t+1} - Y_t\| \leq \lambda h_t \|P_U f(Y_t) P_V\| + (1 - \lambda + \lambda\mu) \|X_t\|$$

$$693 \quad \leq \lambda h_t \|P_U f(Y_t) P_V\| + (1 - \lambda + \lambda\mu) \sum_{s=1}^t \mu^{t-s} h_s \|P_U f(Y_s) P_V\|.$$

694 From this inequality as well as the definitions of  $\lambda$  and  $\mu$  it follows that

$$695 \quad (5.32) \quad \mathbb{E}[\|Y_{t+1} - Y_t\|] \leq \frac{2B}{1-\mu} \left( \frac{h_t}{2} + \frac{(1-\mu)}{2} \sum_{s=1}^t \mu^{t-s} h_s \right) =: \frac{2B}{1-\mu} \frac{h_t + \tilde{h}_t}{2}$$

696 for  $\tilde{h}_t = (1-\mu) \sum_{s=1}^t \mu^{t-s} h_s$ . To bring these values together, consider now a new set  
 697 of sequences

$$698 \quad a_t := h_t,$$

$$699 \quad \tilde{a}_t := \frac{h_t + \tilde{h}_t}{2},$$

$$700 \quad b_t := \mathbb{E}[\|P_U F(Y_t) P_V\|].$$

701 Using the Stolz theorem, we can compute the limit

$$\begin{aligned}
702 \quad \lim_{t \rightarrow \infty} \frac{h_t}{\tilde{h}_t} &= \frac{1}{1 - \mu} \lim_{t \rightarrow \infty} \frac{h_t/\mu^t}{h_t/\mu^t + h_{t-1}/\mu^{t-1} + \dots + h_1/\mu} \\
703 \quad &= \frac{1}{1 - \mu} \lim_{t \rightarrow \infty} \frac{h_{t+1}/\mu^{t+1} - h_t/\mu^t}{h_{t+1}/\mu^{t+1}} = \frac{1}{1 - \mu} \lim_{t \rightarrow \infty} \left(1 - \mu \frac{h_t}{h_{t+1}}\right) = 1,
\end{aligned}$$

704 from which one gets

$$705 \quad \lim_{t \rightarrow \infty} \frac{a_t}{\tilde{a}_t} = 1.$$

706 Also, by using Jensen's inequality and (5.31), we can see

$$707 \quad \sum_{t=1}^{\infty} a_t b_t^2 = \sum_{t=1}^{\infty} h_t (\mathbb{E}[\|P_U F(Y_t) P_V\|])^2 \leq \sum_{t=1}^{\infty} h_t \mathbb{E}[\|P_U F(Y_t) P_V\|^2] < \infty.$$

708 Lastly, using the inverse triangle inequality and (5.32), we obtain

$$\begin{aligned}
709 \quad |b_{t+1} - b_t| &= |\mathbb{E}[\|P_U F(Y_{t+1}) P_V\|] - \mathbb{E}[\|P_U F(Y_t) P_V\|]| \\
710 \quad &\leq \mathbb{E}[\|P_U (F(Y_{t+1}) - F(Y_t)) P_V\|] \leq L \mathbb{E}[\|Y_{t+1} - Y_t\|] \\
711 \quad &\leq \frac{2LB}{1 - \mu} \frac{h_t + \tilde{h}_t}{2} = \frac{2LB}{1 - \mu} \tilde{a}_t.
\end{aligned}$$

712 Lemma 5.11 is thus applicable and yields

$$713 \quad \lim_{t \rightarrow \infty} \mathbb{E}[\|P_U F(Y_t) P_V\|] = \lim_{t \rightarrow \infty} b_t = 0. \quad \square$$

714 This result shows that, under Assumption 5.5.1, the SUM Algorithm 3.1 yields a  
715 sequence of iterates that comes arbitrarily close to satisfying the necessary condition  
716 for optimality (3.10) in expectation. If, like in the previous section, Assumption 5.5.2  
717 is also imposed, then this exact statement also holds for the stronger optimality  
718 condition (3.9).

719 Notice that the result in Theorem 5.15 even improves upon the previous one  
720 from Corollary 5.10. By setting  $\lambda = \mu = 0$  in (5.24), we obtain the stochastic  
721 gradient descent algorithm, so (5.28) is valid even for the SGD Algorithm 3.1 under  
722 Assumption 5.5.1. Although (5.28) does not induce a result for squared norms like  
723 (5.21), it shows so-called last-iterate convergence as opposed to the convergence of a  
724 subsequence or in average. Furthermore, since the equation (5.31) is identical to the  
725 statement of Theorem 5.9, the exact result of Corollary 5.10 is effectively obtained at  
726 that moment in the above proof.

727 This finding indicates that the SUM Algorithm 3.1 is also a valid optimization  
728 algorithm for training individual layers of DNNs. The fact that it relies entirely on  
729 Assumption 5.5, however, means that it cannot be used at the beginning of training. In  
730 practice, one should use an algorithm such as the Rank-Adaptive SGD Algorithm 3.1  
731 until the ranks of the parameter matrices become constant and then switch to the  
732 SUM Algorithm 3.1. The fact that momentum methods perform better when training  
733 DNNs in practice [3] indicates that this might be a better approach than only using  
734 the SGD Algorithm 3.1.

735 **6. Conclusion and Further Research.** This work aimed to investigate the  
 736 properties of DLRA (more precisely, of the Rank-Adaptive BUG integrator [7]) as an  
 737 optimization algorithm for machine learning, an approach proposed in [21]. This in-  
 738 volved using stochastic gradients instead of deterministic ones and common machine  
 739 learning methods as solvers of the differential equations that constitute the integra-  
 740 tor. We showed that these modifications not only had little impact on the method’s  
 741 ability to perform low-rank approximation but also yielded algorithms that exhibit  
 742 convergence on the task of training individual layers of deep neural networks.

743 Future research on this topic could tackle some of the weaker points of this work.  
 744 For example, we are always fixing all layers of the network apart from the one that  
 745 we are currently training. It might be possible to frame our approach as a coordi-  
 746 nate descent method and obtain convergence results for the whole network instead  
 747 of just one layer. Also, it appears pertinent to look for a deeper understanding of  
 748 Assumption 5.5 and why or why not it is a sensible one, since all of our convergence  
 749 results rely heavily on it. Since it is essentially an assumption about the S-step of  
 750 the algorithm, understanding it better could offer insight into the weak optimality  
 751 condition (3.10) in the context of the manifold  $\mathcal{M}_q$ , since it derives its structure also  
 752 from the S-step. In general, the need for the condition (3.10) and Assumption 5.5  
 753 stems from an over-reliance on the S-step in our proofs.

754 Our findings can be interpreted as a theoretical validation of the advantageous  
 755 convergence behavior reported in [21]. They show that DLRA is generally capable of  
 756 finding optimal low-rank layers of neural networks during training. We can safely say  
 757 that this makes it a very promising network pruning technique that should be studied  
 758 further and will save considerable computational resources for practitioners.

## REFERENCES

- 759
- 760 [1] D. P. BERTSEKAS AND J. N. TSITSIKLIS, *Gradient convergence in gradient methods with errors*,  
 761 SIAM journal on optimization, 10 (2000), pp. 627–642.
- 762 [2] D. BLALOCK, J. J. GONZALEZ ORTIZ, J. FRANKLE, AND J. GUTTAG, *What is the state of neural*  
 763 *network pruning?*, Proceedings of machine learning and systems, 2 (2020), pp. 129–146.
- 764 [3] L. BOTTOU, F. E. CURTIS, AND J. NOCEDAL, *Optimization methods for large-scale machine*  
 765 *learning*, SIAM review, 60 (2018), pp. 223–311.
- 766 [4] G. CERUTI, L. EINKEMMER, J. KUSCH, AND C. LUBICH, *A robust second-order low-rank BUG*  
 767 *integrator based on the midpoint rule*, arXiv preprint arXiv:2402.08607, (2024).
- 768 [5] G. CERUTI, J. KUSCH, AND C. LUBICH, *A rank-adaptive robust integrator for dynamical low-*  
 769 *rank approximation*, BIT, 62 (2022), pp. 1149–1174.
- 770 [6] G. CERUTI, J. KUSCH, AND C. LUBICH, *A parallel rank-adaptive integrator for dynamical low-*  
 771 *rank approximation*, (2023), <https://arxiv.org/abs/2304.05660>.
- 772 [7] G. CERUTI AND C. LUBICH, *An unconventional robust integrator for dynamical low-rank ap-*  
 773 *proximation*, BIT, 62 (2022), pp. 23–44.
- 774 [8] Y. CHENG, F. X. YU, R. S. FERIS, S. KUMAR, A. CHOUDHARY, AND S.-F. CHANG, *An explo-*  
 775 *ration of parameter redundancy in deep networks with circulant projections*, in 2015 IEEE  
 776 International Conference on Computer Vision (ICCV), IEEE, 2015, pp. 2857–2865.
- 777 [9] Y. CHENG, F. X. YU, R. S. FERIS, S. KUMAR, A. CHOUDHARY, AND S.-F. CHANG, *An explo-*  
 778 *ration of parameter redundancy in deep networks with circulant projections*, in Proceedings  
 779 of the IEEE international conference on computer vision, 2015, pp. 2857–2865.
- 780 [10] A. DÉFOSSEZ, L. BOTTOU, F. BACH, AND N. USUNIER, *A simple convergence proof of Adam*  
 781 *and Adagrad*, arXiv.org, (2022).
- 782 [11] J. FRANKLE AND M. CARBIN, *The lottery ticket hypothesis: Finding sparse, trainable neural*  
 783 *networks*, arXiv preprint arXiv:1803.03635, (2018).
- 784 [12] M. KHODAK, N. TENENHOLTZ, L. MACKEY, AND N. FUSI, *Initialization and regularization of*  
 785 *factorized neural layers*, arXiv preprint arXiv:2105.01029, (2021).
- 786 [13] E. KIERI, C. LUBICH, AND H. WALACH, *Discretized dynamical low-rank approximation in the*  
 787 *presence of small singular values*, SIAM journal on numerical analysis, 54 (2016), pp. 1020–

- 788 1038.  
789 [14] E. KIERI AND B. VANDEREYCKEN, *Projection methods for dynamical low-rank approximation of*  
790 *high-dimensional problems*, Computational Methods in Applied Mathematics, 19 (2019),  
791 pp. 73–92.  
792 [15] O. KOCH AND C. LUBICH, *Dynamical low-rank approximation*, SIAM Journal on Matrix Analy-  
793 sis and Applications, 29 (2007), pp. 434–454, <https://doi.org/10.1137/050639703>.  
794 [16] J. KUSCH, *Second-order robust parallel integrators for dynamical low-rank approximation*, arXiv  
795 preprint arXiv:2403.02834, (2024).  
796 [17] J. M. LEE, *Introduction to Smooth Manifolds by John M. Lee*, Graduate Texts in Mathematics,  
797 218, New York, NY, 2003.  
798 [18] J. LIU, D. XU, Y. LU, J. KONG, AND D. P. MANDIC, *Last-iterate convergence analysis of*  
799 *stochastic momentum methods for neural networks*, Neurocomputing (Amsterdam), 527  
800 (2023), pp. 27–35.  
801 [19] C. LUBICH AND I. V. OSELEDETS, *A projector-splitting integrator for dynamical low-rank ap-*  
802 *proximation*, BIT Numerical Mathematics, 54 (2014), pp. 171–188.  
803 [20] H. SATO, *Riemannian Optimization and Its Applications by Hiroyuki Sato*, SpringerBriefs in  
804 Control, Automation and Robotics, Cham, 1st ed. 2021. ed., 2021.  
805 [21] S. SCHOTTHÖFER, E. ZANGRANDO, J. KUSCH, G. CERUTI, AND F. TUDISCO, *Low-rank lot-*  
806 *tery tickets: finding efficient low-rank neural networks via matrix differential equa-*  
807 *tions*, in Advances in Neural Information Processing Systems, S. Koyejo, S. Mohamed,  
808 A. Agarwal, D. Belgrave, K. Cho, and A. Oh, eds., vol. 35, Curran Associates,  
809 Inc., 2022, pp. 20051–20063, [https://proceedings.neurips.cc/paper\\_files/paper/2022/file/](https://proceedings.neurips.cc/paper_files/paper/2022/file/7e98b00eeafcdab0c5661fb9355be3a-Paper-Conference.pdf)  
810 [7e98b00eeafcdab0c5661fb9355be3a-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2022/file/7e98b00eeafcdab0c5661fb9355be3a-Paper-Conference.pdf).  
811 [22] S. SCHOTTHÖFER, E. ZANGRANDO, J. KUSCH, G. CERUTI, AND F. TUDISCO, *Dlrt-net*. [https:](https://github.com/COMPiLELab/DLRT-Net)  
812 [//github.com/COMPiLELab/DLRT-Net](https://github.com/COMPiLELab/DLRT-Net), 2023.  
813 [23] P. VILLALOBOS, J. SEVILLA, T. BESIROGLU, L. HEIM, A. HO, AND M. HOBBAHN, *Machine*  
814 *learning model sizes and the parameter gap*, arXiv.org, (2022).  
815 [24] H. WANG, S. AGARWAL, AND D. PAPALIOPOULOS, *Pufferfish: Communication-efficient models*  
816 *at no extra cost*, Proceedings of Machine Learning and Systems, 3 (2021), pp. 365–386.  
817 [25] S. J. WRIGHT AND B. RECHT, *Optimization for Data Analysis*, Cambridge University Press,  
818 2022.  
819 [26] E. ZANGRANDO, S. SCHOTTHÖFER, G. CERUTI, J. KUSCH, AND F. TUDISCO, *Rank-adaptive*  
820 *spectral pruning of convolutional layers during training*, arXiv preprint arXiv:2305.19059,  
821 (2023).