

---

# Federated Dynamical Low-Rank Training with Global Loss Convergence Guarantees

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1        In this work, we propose a federated dynamical low-rank training (FeDLRT)  
2        scheme to reduce client compute and communication costs - two significant per-  
3        formance bottlenecks in horizontal federated learning. Our method builds upon  
4        dynamical low-rank splitting schemes for manifold-constrained optimization to  
5        create a global low-rank basis of network weights, which enables client training on  
6        a small coefficient matrix. A consistent global low-rank basis allows us to incorpo-  
7        rate a variance correction scheme and prove global loss descent and convergence  
8        to a stationary point. Dynamic augmentation and truncation of the low-rank bases  
9        automatically optimizes computing and communication resource utilization. We  
10       demonstrate the efficiency of FeDLRT in an array of computer vision benchmarks  
11       and show a reduction of client compute and communication costs by up to an order  
12       of magnitude with minimal impacts on global accuracy.

## 13    1 Introduction

14    Federated learning (FL) [20, 33, 23] builds a global model on a central *server* from data distributed  
15    on multiple devices, i.e., *clients*, by iteratively aggregating local models trained with the computation  
16    resource on the clients. In horizontal FL, where all clients share identical model architecture and  
17    data features, computation is often limited by (i) the communication bandwidth between clients and  
18    the server and (ii) the restricted compute and memory resources at each client. The former could be  
19    addressed by deploying various compression techniques, such as sparse randomized sketching [9],  
20    subsampling [18], or by allowing for partial [23, 26] or asynchronous [35, 4] communications. The  
21    latter could be addressed by sparse training [29, 41] and transfer learning [5].

22    Since FedAvg [23], low-rank methods have been proposed to increase communication and compute  
23    efficiency for FL in [28, 43, 21, 42, 40, 12, 18, 30]. These methods can be categorized into: 1) methods  
24    that purely reduce communication cost by communicating only the low-rank factors obtained by  
25    performing a full-size SVD (or similar factorization methods) on the weight matrix after client  
26    optimization [28, 37, 40] and 2) methods that reduce both communication and client compute costs  
27    by learning only low-rank factors on clients [21, 43, 42, 12, 18].

28    **Contribution:** This work focuses on the horizontal FL setting and addresses the challenges of  
29    communication bandwidth and client compute resources simultaneously by leveraging low-rank  
30    approximations of weight matrices that follow the dynamics of the gradient flow. The proposed  
31    method features 1) **Efficient communication** — only transmitting low-rank factors; 2) **Low client**  
32    **compute and memory footprint** — clients optimizing only a small coefficient matrix; 3) **Automatic**  
33    **server-side compression** — minimizing memory and communication requirements during training  
34    via server-side dynamical rank adjustment; 4) **Global loss convergence guarantees** — converging  
35    to a stationary point by incorporating a variance correction scheme [24]. Each of these features is

36 demonstrated on benchmark problems. To the best of the authors’ knowledge, this is the first low-rank  
 37 method possessing all these features.

## 38 2 Background and problem statement

39 **Federated optimization** typically considers *distributed* setups and with *limited communication* and  
 40 *limited client compute and memory* resources [23]. In this work, we consider a general federated  
 41 optimization problem, i.e.,

$$\min_w \mathcal{L}(w) := \frac{1}{C} \sum_{c=1}^C \mathcal{L}_c(w), \quad (1)$$

42 where  $w$  is a trainable weight,  $\mathcal{L}$  is the global loss function associated to a global dataset  
 43  $X$ , and  $\mathcal{L}_c$  is the local loss function of client  $c$  with local dataset  $X_c$  in a federated  
 44 setup with  $C$  clients. For notational simplicity, we consider that  $X = \cup_{c=1}^C X_c$  and  
 45 each  $X_c$  is of the same size. Therefore,  $\mathcal{L}$  is an average of  $\mathcal{L}_c$  with uniform weights.  
 46 The extension to handle a (non-uniform)  
 47 weighted average case is straightforward.  
 48 As the first baseline for federated optimization,  
 49 we consider FedAvg [23], see Algo-  
 50 rithm 3. Here, each client optimizes its local  
 51 loss function  $\mathcal{L}_c$  for  $s_*$  local iterations  
 52 using gradient descent,

$$w_c^{s+1} = w_c^s - \lambda \nabla_w \mathcal{L}(w_c^s), \quad (2)$$

53 with learning rate  $\lambda$ , for  $s = 0, \dots, s_* - 1$ .  
 54 The initial value for the local iteration is  
 55 the last global weight, i.e.,  $w_c^0 = w^t$ . After  
 56 local iterations, the weights are commu-  
 57 nicated to and aggregated at the server to  
 58 update the global weight following

$$w^{t+1} = \frac{1}{C} \sum_{c=1}^C w_c^{s_*}. \quad (3)$$

59 **Client-drift effect** is a common challenge  
 60 in FL, where the iterative client updates (2)  
 61 of FedAvg converge to local minima and jeopardize global training performance since the average  
 62 of the local minimizers may be far away from the global minimizer. These effects are particularly  
 63 pronounced for a large number of local iterations  $s_*$ , or high discrepancies between local loss  
 64 functions  $\mathcal{L}_c$ , as illustrated by Figure 1. Multiple methods [33, 20, 27, 14, 39] have been proposed to  
 65 mitigate this issue. However, these methods often exhibit a *speed-accuracy conflict*, where learning  
 66 rates need to be heavily reduced; thus, convergence is slow.

67 **Variance correction**<sup>1</sup> introduced in the FedLin method [24] constructs a variance correction term  
 68  $V_c = \nabla_w \mathcal{L}_c(w^t) - \frac{1}{C} \sum_{c=1}^C \nabla_w \mathcal{L}_c(w^t)$  and modifies the client update iteration to

$$w_c^{s+1} = w_c^s - \lambda (\nabla_w \mathcal{L}(w_c^s) - V_c), \quad s = 0, \dots, s_* - 1. \quad (4)$$

69 This technique leads to global convergence to the minimizer of (1) with constant learning rates [24]  
 70 for convex  $\mathcal{L}$  and else to convergence to a stationary point, at the cost of an additional communication  
 71 round for computing the variance correction.

72 **Federated neural network training** considers problem (1) with the trainable weight  $w$  being the set  
 73 of weight matrices  $\{W_i\}_i^L$  of an  $L$  layer neural network. In each iteration, the weight updates in (2)  
 74 and (4) are applied to all layers simultaneously. Therefore, w.l.o.g., we express the local loss function  
 75 as  $\mathcal{L}_c(W)$ , where  $W \in \mathbb{R}^{n \times n}$  denotes the weight matrix of an arbitrary layer.

76 **Low-rank neural network training:** An array of recent work has provided theoretical and experi-  
 77 mental evidence that layer weights of over-parameterized networks tend to be low rank [1, 2, 8, 22]  
 78 and that removing small singular values may even lead to increased model performance while dramati-  
 79 cally reducing model size [34, 32] in non-federated scenarios. This beneficial feature has spawned a

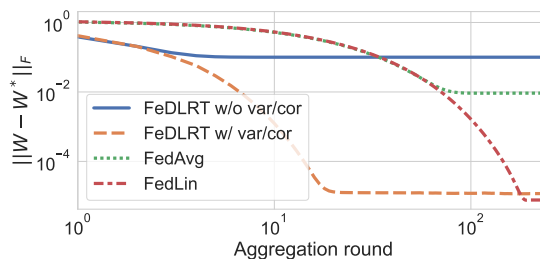


Figure 1: Federated, heterogeneous least squares regression problem, see Section 4.1, for  $C = 4$  clients,  $s_* = 100$  iterations, learning rate  $\lambda = 1e - 3$  and  $C$  rank-1 local target functions. FL methods without variance correction plateau quickly, whereas FedLin and FeDLRT with variance correction converge to  $1e - 5$ . FeDLRT converges faster than FedLin and has lower communication costs.

<sup>1</sup>Variance correction is commonly referred to as “variance reduction” [17, 24].

rich landscape of methods to compress neural networks to a low-rank factorization after training with subsequent fine-tuning [31, 6, 36, 19], train the factorized network with fixed rank [13, 38, 15], dynamically adjust the rank during training [32, 44], or use low-rank adapters for fine-tuning foundation models [11, 7, 45].

**Dynamical Low-rank Approximation of the gradient flow of neural network training.** The core contribution of this paper builds on the dynamical low-rank approximation (DLRA) method, which was initially proposed for solving matrix equations [16] and recently extended to neural network training [32, 44, 10]. Let  $\dot{W}(t) = -\nabla_W \mathcal{L}(W(t))$  denote the gradient flow for minimizing  $\mathcal{L}$ .

The DLRA method restricts the trajectory of  $W$  to  $\mathcal{M}_r$ , the manifold of  $n \times n$ , rank- $r$  matrices, by projecting  $\dot{W}$  onto a local tangent plane of  $\mathcal{M}_r$  via an orthogonal projection. This guarantees a low-rank solution when following the projected dynamics from a low-rank initial guess. Let the low-rank matrix take the form  $W_r = USV^\top \in \mathcal{M}_r$  with  $U, V \in \mathbb{R}^{n \times r}$  the orthonormal bases of  $\mathcal{M}_r$  and  $S \in \mathbb{R}^{r \times r}$  the coefficient matrix. The dynamics for each low-rank factor in DLRA are then derived in [16, Proposition 2.1] as

$$\begin{aligned}\dot{S}(t) &= -U^\top(t) \nabla_W \mathcal{L}(U(t)S(t)V(t)^\top) V(t), \\ \dot{U}(t) &= -(I - P_{U(t)}) \nabla_W \mathcal{L}(U(t)S(t)V(t)^\top) V(t) S(t)^{-1}, \\ \dot{V}(t) &= -(I - P_{V(t)}) \nabla_W \mathcal{L}(U(t)S(t)V(t)^\top) U(t) S(t)^{-\top},\end{aligned}\tag{5}$$

where  $P_U = UU^\top$  and  $P_V = VV^\top$  are the projections onto the column spaces of  $U$  and  $V$ , respectively. By using the *basis update & Galerkin* (BUG) scheme [3], (5) can be split into a basis update step for  $U$  and  $V$  and a coefficient update step for  $S$ . This splitting scheme allows for dynamic adjustment of the rank via a basis augmentation before the coefficient update step and a basis truncation after the coefficient update, as shown in [32].

### 3 FeDLRT: Federated dynamical low-rank training with variance correction

In this section, we present the core contribution of this paper, *federated dynamical low-rank training* (FeDLRT), which features a low-rank client optimization step with optional variance correction and an efficient server aggregation process that dynamically determines the optimal weight matrix rank for automatic compression.

In the context of FL, the BUG of DLRA splitting scheme is particularly interesting since it allows for learning the low-rank bases and coefficients in separate steps. This gives rise to a globally shared basis for the local client iterations, reducing communication and client compute cost of the proposed FeDLRT scheme, see Figure 2: First, the factorization is broadcast to the clients (panel 1), and the basis gradients<sup>2</sup>  $U, V$  are aggregated on the server (panel 2). Next, the basis is augmented on the server (panel 3) and broadcast. On the clients, only the augmented coefficient matrix  $S$  is updated repeatedly (panel 4) before aggregation to the server. After aggregation of the local augmented coefficient matrices, redundant basis directions are eliminated to optimize the accuracy-to-compression ratio of the model on the server.

The strategy yields the following benefits compared to “full-rank” FL schemes as FedLin [24] and low-rank schemes with local compression:

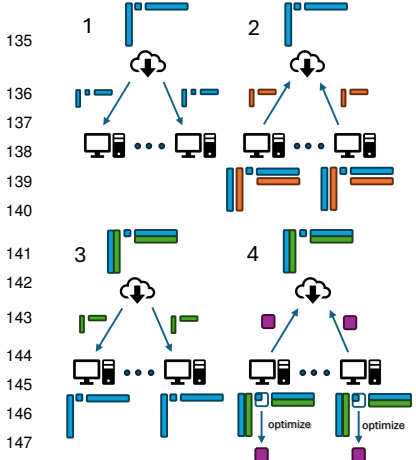
**Low client compute cost:** Server-based basis augmentation and compression enables an automatic compression without a-priori knowledge of the layer rank  $r$  and at no cost for the resource-constrained clients. The clients only evaluate gradients of low-rank factors and optimize the small matrix  $S \in \mathbb{R}^{r \times r}$ .

**Efficient communication:** Similar to FedLin, FeDLRT requires *in practice* two communication rounds – one for aggregating and distributing global gradients for basis augmentation and variance correction and one for aggregating locally updated coefficients. However, communication cost for each round is significantly reduced since only low-rank factors are communicated. We refer to Section 3.3 on communication and compute cost.

Existing federated low-rank schemes effectively generate individual and incompatible representations of  $W_r \in \mathcal{M}_r$  for each client. While the factors can still be efficiently communicated, averaging on

<sup>2</sup>and later on the coefficient gradients for variance correction

126 the server requires a reconstruction of the full weigh matrix  $W^* = \frac{1}{C} \sum_{c=1}^C U_c S_c V_c^\top$ , since the local  
 127 manifolds possibly diverge. Thus, the local rank information is lost and needs to be costly recovered  
 128 by a full  $n \times n$  SVD on the server; see Algorithm 6 for details. Since the average of low-rank matrices  
 129 is not necessarily of low rank, these schemes may lose crucial information on the manifold if client  
 130 solutions drift too far apart from each other. FeDLRT, in contrast, provides the advantage of **client-**  
 131 **wide manifold consistency**: Splitting the low-rank update and sharing bases amongst clients provides  
 132 a globally consistent manifold basis. This furthermore allows for bounding the coefficient drift, see  
 133 Theorem 1, and enables a variance correction for the federated low-rank similar to the FedLin scheme.  
 134



148 Figure 2: Communication of  
 149 FeDLRT without variance correction. 1) Broadcast global basis  
 150  $U, V$  (blue). 2) Aggregate basis  
 151 gradients  $G_{c,U}, G_{c,V}$  (orange). 3)  
 Broadcast global augmented basis  
 $\tilde{U}, \tilde{V}$  (green). 4) Aggregate client  
 coefficient update  $\tilde{S}_c^{s^*}$  (purple).

152 We denote the augmented bases by  $\tilde{U} = [U^t \mid \bar{U}]$  and  $\tilde{V} = [V^t \mid \bar{V}]$ . The orthonormalization is performed on the server, providing compute cost reduction for the client.

155 **Basis broadcasting** of  $\tilde{U}$  and  $\tilde{V}$  only requires to broadcast the new bases  $\bar{U}$  and  $\bar{V}$ , since  $U^t$  and  $V^t$   
 156 are readily available on the clients. Formally, the coefficients  $S^t$  are projected onto the augmented  
 157 basis, i.e.,  $\tilde{S} = \tilde{U}^\top U^t S^t V^t \tilde{V} \in \mathbb{R}^{2r \times 2r}$ , before broadcasting them to the clients. Exploiting the  
 158 orthonormality of the basis results in further reduction of the communication and compute cost:

159 **Lemma 1.**  $\tilde{S} = \tilde{U}^\top U^t S^t V^t \tilde{V}$  takes the form  $\tilde{S} = \begin{bmatrix} S^t & 0 \\ 0 & 0 \end{bmatrix}$ .

160 See Appendix F for the proof. With Lemma 1, only  $\bar{U}$  and  $\bar{V}$  have to be broadcast, and the augmented  
 161 bases and coefficients  $\tilde{U}, \tilde{V}$ , and  $\tilde{S}$  can be assembled on each client as needed. Furthermore, only  
 162  $S \in \mathbb{R}^{r \times r}$ , instead of  $\tilde{S} \in \mathbb{R}^{2r \times 2r}$ , needs to be communicated.

163 Below, we discuss three options for the client coefficient update step.

164 **Client coefficient update** without variance correction is implemented similarly to FedAvg (3). On  
 165 each client  $c$ , the augmented coefficient matrix  $\tilde{S}_c$  is trained for  $s_*$  iterations<sup>3</sup> with learning rate  $\lambda$ ,

$$\tilde{S}_c^{s+1} = \tilde{S}_c^s - \lambda \nabla_{\tilde{S}} \mathcal{L}_c(\tilde{U} \tilde{S}_c \tilde{V}^\top), \quad s = 0, \dots, s_* - 1, \quad \text{with} \quad \tilde{S}_c^{s=0} = \tilde{S}. \quad (7)$$

166 **Client coefficient update with variance correction** is required in certain federated scenarios, e.g.,  
 167 the case considered in Figure 1. Based on FedLin [24], we introduce a correction step for the local  
 168 coefficient update of FeDLRT. It extends the above local iteration by another communication round,

<sup>3</sup>Our analysis focuses on the case where all clients share the same number of local iterations  $s_*$ . The analysis can be extended to the case where  $s_*$  is client dependent, following a similar strategy as in [24].

### 3.1 Description of Algorithm 1 - FeDLRT

In this section, we elaborate on the details in Algorithm 1. The orthonormal factors  $U^t, V^t$  and the coefficient matrix  $S^t$  are initialized with rank  $r$  and then broadcast to the clients. Note that FeDLRT ensures that, for all  $t > 1$ ,  $U^t$  and  $V^t$  are orthonormal, and  $S^t$  is diagonal and full rank.

**Basis augmentation** of the bases  $U^t$  and  $V^t$  is performed using concatenation with the corresponding global basis gradients  $G_U = \frac{1}{C} \sum_{c=1}^C \nabla_U \mathcal{L}_c(U^t S^t V^t \top)$  and  $G_V = \frac{1}{C} \sum_{c=1}^C \nabla_V \mathcal{L}_c(U^t S^t V^t \top)$ , obtained by aggregating the local basis gradients.  $G_U$  and  $G_V$  encapsulate the gradient flow dynamics (5) projected onto the original bases, thus yielding an intuitive choice for basis augmentation. Further, this choice is consistent with the basis update step of the augmented BUG splitting scheme, see Appendix E, which ensures the robustness of the client optimizer. Subsequent orthonormalization, e.g., by a QR decomposition, yields the augmented basis, i.e.,

$$\begin{aligned} [U^t \mid \bar{U}]R &= \text{qr}([U^t \mid G_U]) \in \mathbb{R}^{n \times 2r}, \\ \text{and } [V^t \mid \bar{V}]R &= \text{qr}([V^t \mid G_V]) \in \mathbb{R}^{n \times 2r}. \end{aligned} \quad (6)$$

We denote the augmented bases by  $\tilde{U} = [U^t \mid \bar{U}]$  and  $\tilde{V} = [V^t \mid \bar{V}]$ . The orthonormalization is performed on the server, providing compute cost reduction for the client.

169 where the gradient of the augmented coefficients  $G_{\tilde{S},c} = \nabla_{\tilde{S}} \mathcal{L}_c(\tilde{U}\tilde{S}\tilde{V}^\top)$  is computed, aggregated to  
 170  $G_{\tilde{S}} = \frac{1}{C} \sum_{c=1}^C G_{\tilde{S},c}$  and subsequently broadcast. This yields a correction term  $V_c = G_{\tilde{S}} - G_{\tilde{S},c}$  for  
 171 each client  $c$  and thus the client iterations read

$$\tilde{S}_c^{s+1} = \tilde{S}_c^s - \lambda \left( \nabla_{\tilde{S}} \mathcal{L}_c(\tilde{U}\tilde{S}_c^s\tilde{V}^\top) + V_c \right), \quad s = 0, \dots, s_* - 1, \quad \text{with} \quad \tilde{S}_c^{s=0} = \tilde{S}. \quad (8)$$

172 The correction term results in a bound on the coefficient drift and leads to convergence guarantees for  
 173 FeDLRT, as detailed in Section 3.2.

174 **Client coefficient update with simplified variance correction:** Empirically, we observe that a  
 175 simplified variance correction, which only considers the correction term of the *non-augmented*  
 176 coefficients  $S^t$ , is sufficient, see Figure 6. The simplified variance correction term takes the form

$$V_c = G_{\tilde{S}} - G_{\tilde{S},c} \approx \check{V}_c := \check{G}_{\tilde{S}} - \check{G}_{\tilde{S},c} = \begin{bmatrix} \nabla_S \mathcal{L}(U^t S^t V^{t,\top}) - \nabla_S \mathcal{L}_c(U^t S^t V^{t,\top}) & 0 \\ 0 & 0 \end{bmatrix}, \quad (9)$$

177 which makes lines 10 and 12 in Algorithm 1 redundant, since  $\check{G}_{\tilde{S}}$  can be aggregated in one step with  
 178 the basis gradients  $G_U, G_V$  in line 4 and broadcast with  $\tilde{U}, \tilde{V}$  in line 6, reducing the communication  
 179 rounds to two - the same as FedLin. See Algorithm 5 for details.

180 **Coefficient averaging** is performed after (any of the above variants of) the client iterations. The server  
 181 computes the updated global coefficients by averaging the local updates, i.e.,  $\tilde{S}^* = \frac{1}{C} \sum_{c=1}^C \tilde{S}_c^{s_*}$ .  
 182 With the shared augmented bases  $\tilde{U}$  and  $\tilde{V}$ , this is equivalent to the FedAvg aggregation

$$\tilde{W}_r^* = \frac{1}{C} \sum_{c=1}^C \tilde{W}_r^{s_*} = \frac{1}{C} \sum_{c=1}^C \left( \tilde{U} \tilde{S}_c^{s_*} \tilde{V}^\top \right) = \tilde{U} \left( \frac{1}{C} \sum_{c=1}^C \tilde{S}_c^{s_*} \right) \tilde{V}^\top = \tilde{U} \tilde{S}^* \tilde{V}^\top. \quad (10)$$

183 Since the basis is fixed, the rank  $2r$  is preserved in the aggregation, which is in contrast to other  
 184 federated low-rank schemes where the aggregated weights could be full rank and, in turn, require a  
 185 full matrix SVD to determine the new rank [28, 40].

186 **Automatic compression via rank truncation** is necessary 1) to identify the optimal rank of the  
 187 weight matrix and 2) to ensure that  $S$  is full rank<sup>4</sup>. To this end, a truncated SVD of  $\tilde{S}^* \in \mathbb{R}^{2r \times 2r}$  is  
 188 performed, i.e.  $P_{r_1}, \Sigma_{r_1}, Q_{r_1}^\top = \text{svd}(\tilde{S}^*)$ , where  $P_{r_1}, Q_{r_1} \in \mathbb{R}^{2r \times r_1}$  and  $\Sigma_{r_1} = \text{diag}(\sigma_1, \dots, \sigma_{r_1})$   
 189 contains the  $r_1$  largest singular values of  $\tilde{S}^*$ . The new rank  $r_1$  can be chosen by a variety of criteria,  
 190 e.g., a singular value threshold  $\|[\sigma_{r_1}, \dots, \sigma_{2r}]\|_2 < \vartheta$ . Once a suitable rank is determined, the  
 191 factorization is updated by the projection of the bases  $U^{t+1} = \tilde{U} P_{r_1} \in \mathbb{R}^{n \times r_1}$ ,  $V^{t+1} = \tilde{V} Q_{r_1} \in$   
 192  $\mathbb{R}^{n \times r_1}$  and update of the coefficient  $S^{t+1} = \Sigma_{r_1}$ . Remarkably, Algorithm 1 is a federated low-rank  
 193 learning scheme whose solution is close to a full-rank solution, see Theorem 5.

194 FeDLRT can readily be extended to tensor-valued, e.g., convolutional, layers by applying Algorithm 1  
 195 to each basis and the core tensor in a Tucker Tensor factorization. We refer to Appendix B for details.

### 196 3.2 Analysis of FeDLRT with variance correction

197 In this section, we analyze the FeDLRT algorithm under the general assumption that  $\mathcal{L}_c$  and  $\mathcal{L}$  are  
 198  $L$ -smooth with constant  $L$ . Theorems 2 and 3 give the convergence results for FeDLRT with full  
 199 variance correction (8) in Algorithm 1. Theorem 4 and Corollary 1 provide the convergence for  
 200 FeDLRT with simplified variance correction in (9), as detailed in Algorithm 5, under additional  
 201 assumptions given therein. We note that the analysis does not require convexity of  $\mathcal{L}_c$  or  $\mathcal{L}$ .

202 **FeDLRT convergence with full variance correction.** The variance-corrected client iteration (8)  
 203 leads to the following bound the client coefficient drift.

204 **Theorem 1.** *Given augmented basis and coefficient matrices  $\tilde{U}$ ,  $\tilde{V}$ , and  $\tilde{S}$ . If the local learning rate*  
 205  *$0 < \lambda \leq \frac{1}{L s_*}$  with  $s_* \geq 1$  the number of local steps, for all clients  $c$ ,*

$$\|\tilde{S}_c^s - \tilde{S}_c\| \leq \exp(1) s_* \lambda \|\nabla_{\tilde{S}} \mathcal{L}(\tilde{U}\tilde{S}\tilde{V}^\top)\|, \quad \text{for } s = 1, \dots, s_* - 1, \quad (11)$$

206 where  $\tilde{S}_c^s$  is the variance corrected coefficient as given in (8).

<sup>4</sup>Full rank  $S$  is required to show consistency of the basis update step (6) with the robust operator splitting of [3, 32], see Appendix E.

---

**Algorithm 1: FeDLRT** (See Algorithm 2 for auxiliary function definitions)

---

**Input:** Initial orthonormal bases  $U^1, V^1 \in \mathbb{R}^{n \times r}$  and full rank  $S^1 \in \mathbb{R}^{r \times r}$ ;  
Client-server setup with clients  $c = 1, \dots, C$ ;  
var\_cor: Boolean flag to activate variance correction;  
 $\tau$ : singular value threshold for rank truncation.

```

1 for  $t = 1, \dots, T$  do
2   broadcast( $\{U^t, V^t, S^t\}$ )
3    $G_{U,c} \leftarrow \nabla_U \mathcal{L}_c(U^t S^t V^{t,\top}); G_{V,c} \leftarrow \nabla_V \mathcal{L}_c(U^t S^t V^{t,\top})$  /* On client */
4    $G_U, G_V \leftarrow \text{aggregate}(\{G_{U,c}, G_{V,c}\})$ 
5    $\bar{U} \leftarrow \text{basis\_augmentation}(U^t, G_U); \bar{V} \leftarrow \text{basis\_augmentation}(V^t, G_V)$ 
6   broadcast( $\{\bar{U}, \bar{V}\}$ )
7    $\tilde{U} \leftarrow [U^t \mid \bar{U}]; \tilde{V} \leftarrow [V^t \mid \bar{V}]$  /* Basis assembly on client */
8    $\tilde{S}^{s=0} \leftarrow \begin{bmatrix} S^t & 0 \\ 0 & 0 \end{bmatrix}$  /* Coefficient matrix assembly on client */
9   if var_cor then
10     $G_{\tilde{S},c} \leftarrow \nabla_{\tilde{S}} \mathcal{L}_c(\tilde{U} \tilde{S} \tilde{V}^\top)$  /* Augmented gradient on client */
11     $G_{\tilde{S}} \leftarrow \text{aggregate}(\{G_{\tilde{S},c}\})$ 
12    broadcast( $\{G_{\tilde{S}}\}$ )
13    coefficient_update_var_cor( $c, G_{\tilde{S}} - G_{\tilde{S},c}$ ) /* On client */
14  else
15    coefficient_update( $c$ ) /* On client */
16   $\tilde{S}^* \leftarrow \text{aggregate}(\{\tilde{S}_c^*\})$ 
17   $P_{r_1}, \Sigma_{r_1}, Q_{r_1} \leftarrow \text{svd}(\tilde{S}^*)$  with threshold  $\vartheta$  /* Compression step */
18   $U^{t+1} \leftarrow \tilde{U} P_{r_1}; V^{t+1} \leftarrow \tilde{V} Q_{r_1}; S^{t+1} \leftarrow \Sigma_{r_1}$  /* Basis and coefficient update */

```

---

Table 1: Comparison of the computational footprint of FeDLRT with FedAvg, FedLin and several low-rank FL methods. The FeDLRT variants are the only low-rank schemes with linearly scaling (in  $n$ ) memory, compute, and communication costs with automatic compression and variance correction.

Method	Client compute	Client memory	Server compute	Server memory	Com. Cost	Com. Rounds	var/cor.	rank adaptive
FedAVG [23]	$\mathcal{O}(s_s b n^2)$	$\mathcal{O}(2n^2)$	$\mathcal{O}(n^2)$	$\mathcal{O}(2n^2)$	$\mathcal{O}(2n^2)$	1	✗	✗
FedLin [24]	$\mathcal{O}(s_s b n^2)$	$\mathcal{O}(2n^2)$	$\mathcal{O}(n^2)$	$\mathcal{O}(2n^2)$	$\mathcal{O}(4n^2)$	2	✓	✗
FeDLRT w/o var/cor	$\mathcal{O}(s_s b(4nr + 4r^2))$	$\mathcal{O}(4(nr + 2r^2))$	$\mathcal{O}(2nr + (8 + 4n)r^2 + 8r^3)$	$\mathcal{O}(2nr + 4r^2)$	$\mathcal{O}(6nr + 6r^2)$	2	✗	✓
FeDLRT simpl. var/cor	$\mathcal{O}(s_s b(4nr + 4r^2) + r^2)$	$\mathcal{O}(4(nr + 2r^2))$	$\mathcal{O}(2nr + (8 + 4n)r^2 + 8r^3)$	$\mathcal{O}(2nr + 4r^2)$	$\mathcal{O}(6nr + 8r^2)$	2	✓	✓
FeDLRT full var/cor	$\mathcal{O}(s_s b(4nr + 4r^2) + 4r^2)$	$\mathcal{O}(4(nr + 2r^2))$	$\mathcal{O}(2nr + (8 + 4n)r^2 + 8r^3)$	$\mathcal{O}(2nr + 4r^2)$	$\mathcal{O}(6nr + 10r^2)$	3	✓	✓
FeDLR [28]	$\mathcal{O}(s_s b n^2 + n^3)$	$\mathcal{O}(2n^2)$	$\mathcal{O}(n^2 + n^3)$	$\mathcal{O}(4nr)$	$\mathcal{O}(4nr)$	1	✗	✓
Riemannian FL [40]	$\mathcal{O}(2n^2 r + 4nr^2 + 2nr)$	$\mathcal{O}(2n^2)$	$\mathcal{O}(2nr + n^2 r)$	$\mathcal{O}(4nr)$	$\mathcal{O}(4nr)$	1	✗	✓

207 The critical ingredient for the proof, provided in Appendix G.1, is the globally shared augmented  
208 bases. Theorem 1 bounds the drift of the low-rank representations of the local weight, which gives  
209 rise to the following global loss descent guarantee.

210 **Theorem 2.** Let  $U^t S^t V^{t,\top}$  and  $U^{t+1} S^{t+1} V^{t+1,\top}$  be the factorization before and after iteration  $t$   
211 of Algorithm 1 with variance correction and singular value truncation threshold  $\vartheta$ . Let the local  
212 learning rate be  $0 < \lambda \leq \frac{1}{12Ls_*}$ , then the global loss descent is bounded by

$$\mathcal{L}(U^{t+1} S^{t+1} V^{t+1,\top}) - \mathcal{L}(U^t S^t V^{t,\top}) \leq -s_* \lambda (1 - 12s_* \lambda L) \|\nabla_{\tilde{S}} \mathcal{L}(\tilde{U} \tilde{S} \tilde{V}^\top)\|^2 + L\vartheta. \quad (12)$$

213 The proof is provided in Appendix G.2. Theorem 2 paves the way for the following result on  
214 convergence to a global stationary point.

215 **Theorem 3.** Algorithm 1 guarantees that, for learning rate  $\lambda \leq \frac{1}{12Ls_*}$  and final iteration  $T$ ,

$$\min_{t=1, \dots, T} \|\nabla_{\tilde{S}} \mathcal{L}(U^t S^t V^{t,\top})\|^2 \leq \frac{48L}{T} (\mathcal{L}(U^1 S^1 V^{1,\top}) - \mathcal{L}(U^{T+1} S^{T+1} V^{T+1,\top})) + 48L^2 \vartheta. \quad (13)$$

216 The proof is given in Appendix G.3. In particular, this theorem implies convergence of Algorithm 1  
217 for  $T \rightarrow \infty$  up to a  $\vartheta$ -distance to a global stationary point. This is consistent with the numerical

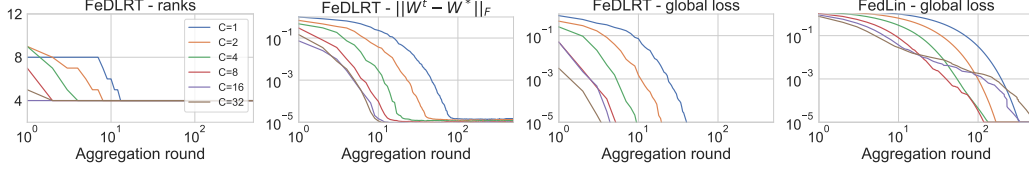


Figure 3: Comparison between FeDLRT with simplified variance correction and FedLin in the homogeneous linear least squares regression test. Each line represents the median result of 20 random initialization with  $C$  clients. The plots from left to right show the rank evolution, the distance to the global optimizer, the global loss values by FeDLRT, and the global loss values by FedLin. The results show that FeDLRT converges faster in this low-rank test case by identifying (and never underestimating) the target rank  $r = 4$  early in the training.

218 results in Figure 1, where FedLin converges to the global minimizer (the only stationary point) while  
 219 FeDLRT with variance correction stops at a point with slightly higher loss value due to a nonzero  $\vartheta$ .  
 220 In the case that the FL problem has a low-rank solution, the truncation error bounded by  $\vartheta$  vanishes,  
 221 and convergence to a stationary point is guaranteed, see, e.g., Figure 3.

222 **FeDLRT convergence with simplified variance correction.** FeDLRT with simplified variance  
 223 correction is detailed in Algorithm 5 with the variance correction term given in (9), which makes  
 224 variance correction more communication and computation efficient but comes at a cost of the  
 225 following additional assumption for convergence analysis.

226 **Assumption 1.** *There exists  $\delta \ll 1$  such that, at each client coefficient update,*

$$\|\nabla_{\tilde{S}} \mathcal{G}(\tilde{U} \tilde{S}_c^s \tilde{V}^\top)\| - \|\nabla_S \mathcal{G}(\tilde{U} \tilde{S}_c^s \tilde{V}^\top)\| < \delta \|\nabla_{\tilde{S}} \mathcal{L}(\tilde{U} \tilde{S} \tilde{V}^\top)\|, \quad (14)$$

227 for functions  $\mathcal{G} = \mathcal{L}$  and  $\mathcal{G} = \mathcal{L}_c$ ,  $c = 1, \dots, C$ .

228 This assumption can be interpreted as that most of dynamics in the gradient flow are captured in  
 229 the coefficient update for the original rank- $r$  matrix  $S$ , and the basis augmentation provides little  
 230 information. This scenario occurs when FeDLRT identifies the optimal rank, which could happen  
 231 early for simpler problems as shown in Figure 3, or when FeDLRT approaches a stationary point.

232 **Theorem 4.** *Under Assumption 1, let  $C := s_* \lambda (1 - \delta^2 - 12s_* \lambda L + \delta^2 s_* \lambda)$ . If the local learning  
 233 rate  $0 < \lambda \leq \frac{1}{12Ls_*}$ , Algorithm 5 leads to the global loss descent*

$$\mathcal{L}(U^{t+1} S^{t+1} V^{t+1, \top}) - \mathcal{L}(U^t S^t V^{t, \top}) \leq -C \|\nabla_{\tilde{S}} \mathcal{L}(\tilde{W}_r)\|^2 + L\vartheta.$$

234 The proof is provided in Appendix H.1. When  $\delta$  is small, this bound is slightly weaker than the one  
 235 in Theorem 2, which leads to the following corollary.

236 **Corollary 1.** *Assume that Assumption 1 holds. Algorithm 5 guarantees that, for the local learning  
 237 rate  $0 < \lambda \leq \frac{1}{s_*(12L + \delta^2)}$ ,*

$$\min_{t=1, \dots, T} \|\nabla_{\tilde{S}} \mathcal{L}(U^t S^t V^{t, \top})\|^2 \leq \frac{96L}{T} (\mathcal{L}(U^1 S^1 V^{1, \top}) - \mathcal{L}(U^{T+1} S^{T+1} V^{T+1, \top})) + 96L^2 \vartheta.$$

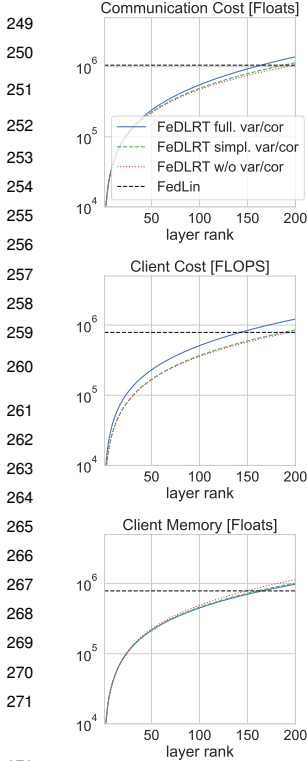
238 The proof is analogous to the one for Theorem 3, see Appendix H.2.

### 239 3.3 Compute and communication cost

240 The proposed FeDLRT methods significantly reduce server and client memory footprint, the required  
 241 communication bandwidth, as well as the client compute cost compared to various baselines, see  
 242 Table 1. We remark that the complete federated learning process is performed on the low-rank factors,  
 243 and the full matrix  $\tilde{W}_r$  is never required, as, e.g., in [28, 40] and FeDLRT is the only low-rank  
 244 method with adaptive compression incorporating variance correction, whose server compute cost  
 245 scales linearly with the layer dimension since the SVD for rank truncation only needs to be computed  
 246 on the augmented coefficient matrix of size  $2r \times 2r$ .

## 247 4 Numerical evaluation

### 248 4.1 Distributed linear least squares regression



272 Figure 4: Scaling of  
 273 communication cost (top)  
 274 compute cost at a single  
 275 client (middle), and  
 276 client memory footprint  
 277 (bottom) for  $s_* = 1$   
 278 client iteration and a single  
 279 data-point for  $W \in$   
 280  $\mathbb{R}^{n \times n}$  with  $n = 512$ . In  
 281 practice we have  $r \ll n$ ,  
 282 see Section 4.

## 283 4.2 ResNet18 on CIFAR10

284 We demonstrate the performance of FeDLRT for training the exemplary ResNet18 model on CIFAR10,  
 285 where we apply FeDLRT to train its fully connected head. The truncation tolerance is set to  
 286  $\vartheta = \tau \|\tilde{S}^*\|$  with  $\tau = 0.01$ . The test case setup is summarized in Table 2. The training data is equally  
 287 partitioned across clients; see Appendix C.2 for the data-preprocessing details. A local iteration  
 288 of Algorithm 1 at client  $c$  describes one mini-batch update on the client training data set  $X_c$  for  
 289 a given batch size,  $s_*$  is the maximum number of local iterations, and  $T$  denotes the number of  
 290 aggregation rounds. We display the statistics for 10 random initializations; each warm-started with  
 291 5 iterations with one client. We set  $s_* = 240/C$  so that in each training run, the global network  
 292 iterates through the same amount of data. This setup favors low client counts, and, as expected, the  
 293 validation accuracy drops as  $C$  grows for FedAvg and FeDLRT without variance correction, see  
 294 Figure 6 (upper row). We note that FeDLRT ties or outperforms FedAvg in terms of final validation  
 295 accuracy. Using full variance correction (second row) increases the validation accuracy of FeDLRT  
 296 by up to 12% in this test case, matching the accuracy of FedLin and enabling FL with 93% accuracy  
 297 for 32 clients. For  $C = 8$  clients, the communication cost saving of the compressed layers is up to  
 298 90%. The computationally more efficient simplified variance correction, using Algorithm 5, (third  
 299 row), yields similar validation accuracy, notably at higher compression ratio and communication cost  
 300 reduction. Similar results are obtained for AlexNet, VGG16 on CIFAR10, and ViT on CIFAR100,

<sup>5</sup>We chose to display the median trajectory to point out its convergence and monotonicity. The test case also converges in the mean.

**Homogeneous test.** We first consider a (convex) FL problem (1) for linear least squares regression with local loss  $\mathcal{L}_c(W) = \frac{1}{2|X_c|} \sum_{(x,y) \in X_c} \|p(x)^\top W p(y) - f(x,y)\|_2^2$ , where  $W \in \mathbb{R}^{n \times n}$  and  $p : [-1, 1] \rightarrow \mathbb{R}^n$  is the Legendre polynomial basis of degree  $n - 1$ . The target function  $f$  is manufactured as  $f(x,y) = p(x)^\top W_r p(y)$ , where  $\text{rank}(W_r) = r$ . We consider problems with  $n = 20$ ,  $r = 4$ , and randomly generated  $W_r$ , with 10,000 data points uniformly sampled on  $[-1, 1]^2$  and uniformly distributed among clients. We compare FeDLRT with variance correction and FedLin with  $s_* = 20$  local iterations and  $\lambda = 1e - 3$  learning rate on  $C = 1, 2, 4, 8, 16, 32$  clients. This setting satisfies the step-size restriction given in Theorem 2. In FeDLRT, the singular value truncation threshold  $\vartheta = \tau \|\tilde{S}^*\|$  with  $\tau = 0.1$  was used.

Figure 3 reports the dynamically updated ranks, errors, and loss values with respect to the aggregation rounds. The reported data are the medians over 20 randomly generated initial weights<sup>5</sup> The results indicate that FeDLRT is able to identify the correct rank within a few aggregation rounds and, furthermore, never underestimates it – which would have increased the loss value significantly. FeDLRT converges to the minimizer  $W^* = W_r$  up to a  $1e - 5$  error and converges faster with more clients. On this problem, FeDLRT shows up to 10x faster convergence than FedLin. We attribute this behavior to the fact that, by identifying a suitable low-rank manifold early in the training, FeDLRT significantly reduces the degrees of freedom in the FL problem.

**Heterogeneous test.** Inspired by [24], we consider a variation of the linear least squares regression with  $\mathcal{L}_c(W) = \frac{1}{2|X_c|} \sum_{(x,y) \in X_c} \|p(x)^\top W p(y) - f_c(x,y)\|_2^2$ , where the target function  $f_c$  is different for each client, and the 10,000 training data points are available to all clients. The local target functions  $f_c$  cause each client to optimize a different local problem. We choose problem size  $n = 10$  with  $C = 4$  clients and use learning rate  $\lambda = 1e - 3$  with  $s_* = 100$  local epochs. As seen in Figure 1, FeDLRT with variance correction converges (to single precision accuracy) to the minimizer  $W^*$  of (1) much faster than FedLin, whereas FeDLRT without correction quickly plateaus, similar to FedAvg.



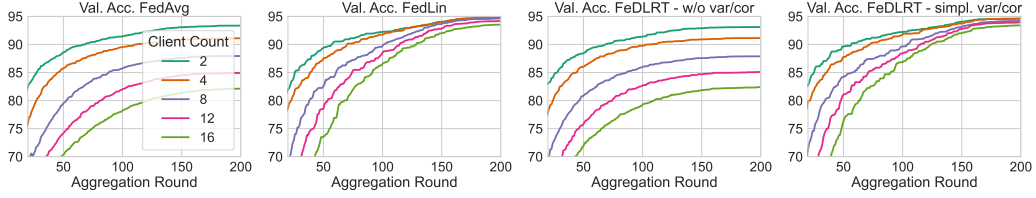


Figure 5: ResNet18 CIFAR10. We compare the convergence behavior of the median result of 10 initializations displaying the best validation accuracy until the current epoch for FedAvg (top left), FedLin (top right), FeDLRT w/o var/cor (bottom left) and FeDLRT w/ simplified var/cor (bottom right). We observe 1) the low-rank methods (bottom) closely follows the convergence dynamics of their full rank counterpart (top), and 2) variance correction starts to improve the convergence behavior during later stages of the training, where the non-corrected methods level off.

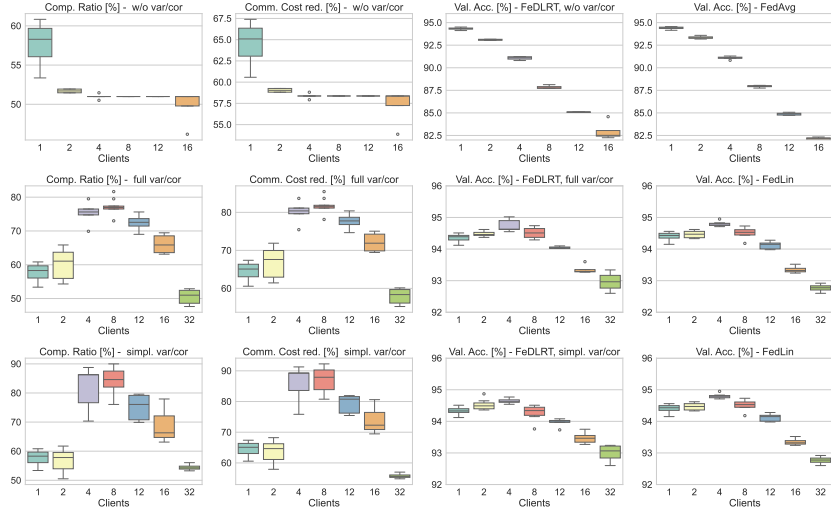


Figure 6: Comparisons for training ResNet18 on CIFAR10 benchmark. Top row compares FeDLRT without variance correction to FedAvg, middle and bottom rows compare FeDLRT with full and simplified variance correction to FedLin, respectively. In each row, the left two panels show the model compression ratio and the communication cost reduction from FeDLRT, and the right two panels show the validation accuracy for FeDLRT and the full-rank counterparts. In each plot, the results are reported for  $C = 1, \dots, 16$  or 32 clients with  $240/C$  local iterations. FeDLRT matches the accuracy of FedAvg and FedLin well, while substantially reducing the server and client memory and communication costs. Variance correction leads to an up to 12% increase in validation accuracy for large  $C$ , mitigating the client drift problem. The simplified variance correction (bottom row) gives comparable results to full version (middle row) at a lower communication and computation cost.

301 see Appendix C, where we observe that FeDLRT closely matches the full-rank accuracy of FedLin.  
 302 Lastly, we remark that variance correction is beneficial for convergence behavior in neural network  
 303 training, as shown in Figure 5.

304 **In conclusion**, we have presented FeDLRT, an efficient low-rank FL scheme with convergence  
 305 guarantees and automatic compression, and demonstrated its capabilities in several test cases.

306 **Limitations and future work:** We remark that the underlying assumption for this work is that  
 307 the target model can be expressed sufficiently well via a low-rank representation. Although the  
 308 communication cost in terms of transferred parameters is significantly reduced compared to existing  
 309 method, FeDLRT still requires two communication handshakes for one aggregation round, just like  
 310 its full-rank counterpart FedLin. Therefore, the method needs to be refined for scenarios where the  
 311 clients have different communication latencies or for completely asynchronous scenarios. Potential  
 312 future research directions include performing large-scale tests with thousands of clients, extending the  
 313 algorithm to accommodate partial client participation or asynchronous communication, and analyzing  
 314 the convergence properties in these scenarios.

315 **References**

- 316 [1] S. Arora, N. Cohen, W. Hu, and Y. Luo. Implicit regularization in deep matrix factorization.  
317 *Advances in Neural Information Processing Systems*, 32, 2019.
- 318 [2] B. Bah, H. Rauhut, U. Terstiege, and M. Westdickenberg. Learning deep linear neural networks:  
319 Riemannian gradient flows and convergence to global minimizers. *Inf. Inference*, 11(1):307–353,  
320 2022.
- 321 [3] G. Ceruti, J. Kusch, and C. Lubich. A rank-adaptive robust integrator for dynamical low-rank  
322 approximation. *BIT Numerical Mathematics*, 2022.
- 323 [4] Y. Chen, Y. Ning, M. Slawski, and H. Rangwala. Asynchronous online federated learning for  
324 edge devices with non-iid data. In *2020 IEEE International Conference on Big Data (Big Data)*,  
325 pages 15–24. IEEE, 2020.
- 326 [5] Y. Chen, X. Qin, J. Wang, C. Yu, and W. Gao. Fedhealth: A federated transfer learning  
327 framework for wearable healthcare. *IEEE Intelligent Systems*, 35(4):83–93, 2020.
- 328 [6] E. L. Denton, W. Zaremba, J. Bruna, Y. LeCun, and R. Fergus. Exploiting linear structure within  
329 convolutional networks for efficient evaluation. *Advances in neural information processing  
330 systems*, 27, 2014.
- 331 [7] T. Dettmers, A. Pagnoni, A. Holtzman, and L. Zettlemoyer. Qlora: Efficient finetuning of  
332 quantized llms. *arXiv preprint arXiv:2305.14314*, 2023.
- 333 [8] T. Galanti, Z. S. Siegel, A. Gupte, and T. Poggio. Sgd and weight decay provably induce a  
334 low-rank bias in neural networks. 2022.
- 335 [9] F. Haddadpour, B. Karimi, P. Li, and X. Li. Fedsketch: Communication-efficient and private  
336 federated learning via sketching, 2020.
- 337 [10] A. Hnatiuk, J. Kusch, L. Kusch, N. R. Gauger, and A. Walther. Stochastic aspects of dynamical  
338 low-rank approximation in the context of machine learning. *Optimization Online*, 2024.
- 339 [11] E. J. Hu, Y. Shen, P. Wallis, Z. Allen-Zhu, Y. Li, S. Wang, L. Wang, and W. Chen. Lora:  
340 Low-rank adaptation of large language models, 2021.
- 341 [12] N. Hyeon-Woo, M. Ye-Bin, and T.-H. Oh. Fedpara: Low-rank hadamard product for  
342 communication-efficient federated learning. In *International Conference on Learning Represen-  
343 tations*, 2022.
- 344 [13] M. Jaderberg, A. Vedaldi, and A. Zisserman. Speeding up convolutional neural networks with  
345 low rank expansions. In *Proceedings of the British Machine Vision Conference. BMVA Press*,  
346 2014.
- 347 [14] S. P. Karimireddy, S. Kale, M. Mohri, S. Reddi, S. Stich, and A. T. Suresh. Scaffold: Stochastic  
348 controlled averaging for federated learning. In *International Conference on Machine Learning*,  
349 pages 5132–5143. PMLR, 2020.
- 350 [15] M. Khodak, N. Tenenholtz, L. Mackey, and N. Fusi. Initialization and regularization of  
351 factorized neural layers. In *International Conference on Learning Representations*, 2021.
- 352 [16] O. Koch and C. Lubich. Dynamical low-rank approximation. *SIAM Journal on Matrix Analysis  
353 and Applications*, 29(2):434–454, 2007.
- 354 [17] J. Konečný, H. B. McMahan, D. Ramage, and P. Richtárik. Federated optimization: Distributed  
355 machine learning for on-device intelligence, 2016.
- 356 [18] J. Konečný, H. B. McMahan, F. X. Yu, P. Richtárik, A. T. Suresh, and D. Bacon. Federated  
357 learning: Strategies for improving communication efficiency, 2017.
- 358 [19] V. Lebedev, Y. Ganin, M. Rakhuba, I. Oseledets, and V. Lempitsky. Speeding-up convolutional  
359 neural networks using fine-tuned CP-decomposition. In *International Conference on Learning  
360 Representations*, 2015.

- 361 [20] T. Li, A. K. Sahu, A. Talwalkar, and V. Smith. Federated learning: Challenges, methods, and  
362 future directions. *IEEE Signal Processing Magazine*, 37(3):50–60, 2020.
- 363 [21] X.-Y. Liu, R. Zhu, D. Zha, J. Gao, S. Zhong, and M. Qiu. Differentially private low-rank  
364 adaptation of large language model using federated learning, 2023.
- 365 [22] C. H. Martin and M. W. Mahoney. Implicit self-regularization in deep neural networks: Evidence  
366 from random matrix theory and implications for learning. *arXiv preprint arXiv:1810.01075*,  
367 2018.
- 368 [23] H. B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. y Arcas. Communication-  
369 efficient learning of deep networks from decentralized data, 2016.
- 370 [24] A. Mitra, R. Jaafar, G. J. Pappas, and H. Hassani. Linear convergence in federated learning:  
371 Tackling client heterogeneity and sparse gradients. In M. Ranzato, A. Beygelzimer, Y. Dauphin,  
372 P. Liang, and J. W. Vaughan, editors, *Advances in Neural Information Processing Systems*,  
373 volume 34, pages 14606–14619. Curran Associates, Inc., 2021.
- 374 [25] A. Mitra, R. Jaafar, G. J. Pappas, and H. Hassani. Linear convergence in federated learning:  
375 Tackling client heterogeneity and sparse gradients, 2021.
- 376 [26] T. Nishio and R. Yonetani. Client selection for federated learning with heterogeneous resources  
377 in mobile edge. In *ICC 2019-2019 IEEE international conference on communications (ICC)*,  
378 pages 1–7. IEEE, 2019.
- 379 [27] R. Pathak and M. J. Wainwright. Fedsplit: An algorithmic framework for fast federated  
380 optimization. *arXiv*.
- 381 [28] Z. Qiao, X. Yu, J. Zhang, and K. B. Letaief. Communication-efficient federated learning with  
382 dual-side low-rank compression, 2021.
- 383 [29] X. Qiu, J. Fernandez-Marques, P. P. Gusmao, Y. Gao, T. Parcollet, and N. D. Lane. Ze-  
384 rofl: Efficient on-device training for federated learning with local sparsity. *arXiv preprint*  
385 *arXiv:2208.02507*, 2022.
- 386 [30] A. Reisizadeh, A. Mokhtari, H. Hassani, A. Jadbabaie, and R. Pedarsani. Fedpaq: A  
387 communication-efficient federated learning method with periodic averaging and quantization,  
388 2020.
- 389 [31] T. N. Sainath, B. Kingsbury, V. Sindhvani, E. Arisoy, and B. Ramabhadran. Low-rank matrix  
390 factorization for deep neural network training with high-dimensional output targets. In *2013*  
391 *IEEE International Conference on Acoustics, Speech and Signal Processing*, pages 6655–6659,  
392 2013.
- 393 [32] S. Schotthöfer, E. Zangrando, J. Kusch, G. Ceruti, and F. Tudisco. Low-rank lottery tickets:  
394 finding efficient low-rank neural networks via matrix differential equations. In *Advances in*  
395 *Neural Information Processing Systems*, 2022.
- 396 [33] O. Shamir, N. Srebro, and T. Zhang. Communication-efficient distributed optimization using  
397 an approximate newton-type method. In *International conference on machine learning*, pages  
398 1000–1008. PMLR, 2014.
- 399 [34] P. Sharma, J. T. Ash, and D. Misra. The truth is in there: Improving reasoning in language models  
400 with layer-selective rank reduction. In *International Conference on Learning Representations*  
401 *(ICLR)*, 2024.
- 402 [35] M. R. Sprague, A. Jalalirad, M. Scavuzzo, C. Capota, M. Neun, L. Do, and M. Kopp. Asyn-  
403 chronous federated learning for geospatial applications. In *Joint European Conference on*  
404 *Machine Learning and Knowledge Discovery in Databases*, pages 21–28. Springer, 2018.
- 405 [36] A. Tjandra, S. Sakti, and S. Nakamura. Compressing recurrent neural network with tensor train.  
406 In *2017 International Joint Conference on Neural Networks (IJCNN)*, pages 4451–4458. IEEE,  
407 2017.

- 408 [37] T. Vogels, S. P. Karimireddy, and M. Jaggi. Powersgd: Practical low-rank gradient compression  
409 for distributed optimization, 2020.
- 410 [38] H. Wang, S. Agarwal, and D. Papailiopoulos. Pufferfish: Communication-efficient models at no  
411 extra cost. *Proceedings of Machine Learning and Systems*, 3:365–386, 2021.
- 412 [39] J. Wang, Q. Liu, H. Liang, G. Joshi, and H. V. Poor. Tackling the objective inconsistency  
413 problem in heterogeneous federated optimization. *Advances in Neural Information Processing*  
414 *Systems*, 33, 2020.
- 415 [40] Y. Xue and V. Lau. Riemannian low-rank model compression for federated learning with  
416 over-the-air aggregation. *IEEE Transactions on Signal Processing*, 71:2172–2187, 2023.
- 417 [41] K. Yang, T. Jiang, Y. Shi, and Z. Ding. Federated learning via over-the-air computation. *IEEE*  
418 *transactions on wireless communications*, 19(3):2022–2035, 2020.
- 419 [42] D. Yao, W. Pan, M. J. O’Neill, Y. Dai, Y. Wan, H. Jin, and L. Sun. Fedhm: Efficient federated  
420 learning for heterogeneous models via low-rank factorization, 2022.
- 421 [43] L. Yi, H. Yu, G. Wang, X. Liu, and X. Li. pfdlora: Model-heterogeneous personalized federated  
422 learning with lora tuning, 2024.
- 423 [44] E. Zangrando, S. Schotthöfer, G. Ceruti, J. Kusch, and F. Tudisco. Rank-adaptive spectral  
424 pruning of convolutional layers during training, 2023.
- 425 [45] J. Zhao, Z. Zhang, B. Chen, Z. Wang, A. Anandkumar, and Y. Tian. Galore: Memory-efficient  
426 llm training by gradient low-rank projection, 2024.
- 427 [46] H. Zhu, B. Chen, and C. Yang. Understanding why vit trains badly on small datasets: An  
428 intuitive perspective, 2023.

429 **A Additional algorithms**

430 In the following, we list a set of algorithms that are used in the paper as a contribution or as a  
 431 baseline method. In particular, Algorithm 2 contains auxiliary function definitions for Algorithm 1  
 432 and Algorithm 5. Algorithm 3 is the standard FedAvg method as presented in [23]. Algorithm 4 is  
 433 the FedLin Algorithm [24], i.e. the extension of Algorithm 4 with variance correction. Algorithm 5  
 434 represents the FeDLRT method with simplified variance correction, as analyzed in Theorem 4 and  
 435 Corollary 1 with the additional Assumption 1.

---

**Algorithm 2:** Auxiliary functions

---

```

1 def broadcast( $\{M_i\}_i$ : list of matrices):
2   | Send  $M_i$  from server to all clients  $\forall i$ 
3 def aggregate( $\{M_{c,i}\}_i$ : list of matrices):
4   | Send  $M_{c,i}$  from client to server  $\forall c, i$ 
5   |  $M_i \leftarrow \frac{1}{C} \sum_{c=1}^C M_c \quad \forall i$ 
6   | return  $\{M_i\}_i$ ;
7 def coefficient_update_var_cor( $c$ : client,  $V_c$ : correction term):
8   | for  $s = 0, \dots, s_* - 1$  do                                     /* On client */
9   |   |  $\tilde{S}_c^{s+1} \leftarrow \tilde{S}_c^s - \lambda \left( \nabla_{\tilde{S}} \mathcal{L}_c(\tilde{U}_c \tilde{S}_c^s \tilde{V}_c^\top) + V_c \right)$ 
10 def coefficient_update( $c$ : client):
11  | for  $s = 0, \dots, s_* - 1$  do                                     /* On client */
12  |   |  $\tilde{S}_c^{s+1} \leftarrow \tilde{S}_c^s - \lambda \nabla_{\tilde{S}} \mathcal{L}_c(\tilde{U}_c \tilde{S}_c^s \tilde{V}_c^\top)$ 
13 def basis_augmentation( $B$ : old basis,  $G_B$ : basis dynamics):
14  |  $[B \mid \bar{B}] \leftarrow \text{qr}([B \mid G_B])$                              /* On server */
15  | return  $\bar{B}$ 

```

---



---

**Algorithm 3:** FedAvg [23]. (See Algorithm 2 for auxiliary function definitions)

---

**Input:** Initial values for weight matrix  $W$   
 Client-server setup with clients  $c = 1, \dots, C$ .

```

1 for  $t = 1, \dots, T$  do
2   | broadcast( $\{W^t\}$ )
3   |  $W_c^{s=0} \leftarrow W^t$ 
4   | for  $s = 0, \dots, s_* - 1$  do
5   |   |  $W_c^{s+1} \leftarrow W_c^s - \lambda \nabla_W \mathcal{L}_c(W_c^s)$            /* Gradient descent on client */
6   |   |  $W^{t+1} \leftarrow \text{aggregate}(\{W_c^{s_*}\})$                  /* Aggregation on server */

```

---



---

**Algorithm 4:** FedLin [24]. (See Algorithm 2 for auxiliary function definitions)

---

**Input:** Initial values for weight matrix  $W$   
 Client-server setup with clients  $c = 1, \dots, C$ .

```

1 for  $t = 1, \dots, T$  do
2   | broadcast( $\{W^t\}$ )
3   |  $G_{W,c} \leftarrow \nabla_W \mathcal{L}_c(W^t)$                                /* Gradient computation on client */
4   |  $G_W \leftarrow \text{aggregate}(\{G_{W,c}\})$                          /* Aggregation on server */
5   | broadcast( $\{G_W\}$ )
6   |  $W_c^{s=0} \leftarrow W^t$ 
7   |  $V_c \leftarrow G_W - G_{W,c}$                                      /* Correction term computation on client */
8   | for  $s = 0, \dots, s_* - 1$  do
9   |   |  $W_c^{s+1} \leftarrow W_c^s - \lambda \nabla_W \mathcal{L}_c(W_c^s) + V_c$  /* Corrected iteration on client */
10  |   |  $W^{t+1} \leftarrow \text{aggregate}(\{W_c^{s_*}\})$                  /* Aggregation on server */

```

---

---

**Algorithm 5:** FeDLRT with simplified variance correction. (See Algorithm 2 for auxiliary function definitions)

---

**Input:** Initial orthonormal bases  $U^1, V^1 \in \mathbb{R}^{n \times r}$  and full rank  $S^1 \in \mathbb{R}^{r \times r}$ ;

Client-server setup with clients  $c = 1, \dots, C$ ;

$\tau$ : singular value threshold for rank truncation.

```

1 for  $t = 1, \dots, T$  do
2   broadcast ( $\{U^t, V^t, S^t\}$ )
3    $G_{U,c} \leftarrow \nabla_U \mathcal{L}_c(U^t S^t V^{t,\top})$  /* On client */
4    $G_{V,c} \leftarrow \nabla_V \mathcal{L}_c(U^t S^t V^{t,\top})$  /* On client */
5    $G_{S,c} \leftarrow \nabla_S \mathcal{L}_c(U^t S^t V^{t,\top})$  /* On client */
6    $G_U, G_V, G_S \leftarrow \text{aggregate}(\{G_{U,c}, G_{V,c}, G_{S,c}\})$ 
7    $\bar{U} \leftarrow \text{basis\_augmentation}(U^t, G_U), \bar{V} \leftarrow \text{basis\_augmentation}(V^t, G_V)$ 
8   broadcast ( $\{\bar{U}, \bar{V}, G_S\}$ )
9    $\tilde{U} \leftarrow [U^t \mid \bar{U}], \tilde{V} \leftarrow [V^t \mid \bar{V}]$  /* Basis assembly on client */
10   $\tilde{S}^{s=0} \leftarrow \begin{bmatrix} S^t & 0 \\ 0 & 0 \end{bmatrix}$  /* Coefficient matrix assembly on client */
11   $\check{G}_{\tilde{S},c} \leftarrow \begin{bmatrix} G_{S,c} & 0 \\ 0 & 0 \end{bmatrix}$  /* Client coeff. gradient approximation on client */
12   $\check{G}_{\tilde{S}} \leftarrow \begin{bmatrix} G_S & 0 \\ 0 & 0 \end{bmatrix}$  /* Global coeff. gradient approximation on client */
13   $\text{coefficient\_update\_var\_cor}(c, \check{G}_{\tilde{S}} - \check{G}_{\tilde{S},c})$  /* On client */
14   $\tilde{S}^* \leftarrow \text{aggregate}(\{\tilde{S}_c^{s*}\})$ 
15   $P_{r_1}, \Sigma_{r_1}, Q_{r_1} \leftarrow \text{svd}(\tilde{S}^*)$  with threshold  $\vartheta$  /* Compression step */
16   $U^{t+1} \leftarrow \tilde{U} P_{r_1}$ , and  $V^{t+1} \leftarrow \tilde{V} Q_{r_1}$  /* Basis projection */
17   $S^{t+1} \leftarrow \Sigma_{r_1}$ 

```

---

## 436 B Extension to convolutions and tensor-valued weights

437 FeDLRT can readily be extended to tensor-valued neural network layers, e.g. convolutional layers,  
438 following [44], where, e.g., a 2D convolution kernel is interpreted as an order-4 tensor and factorized  
439 by using the Tucker decomposition. To this end, the Tucker bases  $U_i \in \mathbb{R}^{n_i \times r_i}$  for  $i = 1, 2, 3, 4$   
440 replace the  $U$  and  $V$  bases in the matrix case, and the Tucker core tensor  $C \in \mathbb{R}^{r_1 \times r_2 \times r_3 \times r_4}$  replaces  
441 the coefficient matrix  $S$ , to which the variance correction is applied. The analysis holds for the Tucker  
442 Tensor case, since Tucker Tensors have a manifold structure. In the analysis, one needs to consider  
443 the gradient projected upon all bases  $U_i$  instead of  $U$  and  $V$ . The compression step is performed with  
444 an truncated Tucker decomposition of the core tensor  $C$ , instead of an SVD of  $S$ . For intuition, one  
445 can also refer to the matrix case as the order-2 Tucker Tensor case. Remark that the bases  $U_i$  are all  
446 updated simultaneously, thus the adaption to the tensor case does not require more communication  
447 rounds.

## 448 C Additional numerical evaluation

### 449 C.1 Compute resources

450 The convex test cases are computed on a single Nvidia RTX 4090 GPU. The computer vision bench-  
451 marks use a set of Nvidia Tesla V100-SXM2-16GB and Tesla P100-PCIE-16GB. For prototyping, a  
452 Nvidia GTX1080ti is used.

### 453 C.2 Data augmentation

454 We use standard data augmentation techniques for the proposed test cases. That is, for CIFAR10,  
455 we augment the training data set by a random horizontal flip of the image, followed by a normal-

---

**Algorithm 6:** Naive implementation of FeDLRT. (See Algorithm 2 for auxiliary function definitions)

---

**Input:** Initial orthonormal bases  $U^1, V^1 \in \mathbb{R}^{n \times r}$  and full rank  $S^1 \in \mathbb{R}^{r \times r}$ ;  
 Client-server setup with clients  $c = 1, \dots, C$ ;  
 $\tau$ : singular value threshold for rank truncation.

```

1 for  $t = 1, \dots, T$  do
2   broadcast  $(\{U^t, V^t, S^t\})$ 
3    $U_c^{s=0}, V_c^{s=0}, S_c^{s=0} \leftarrow U^t, V^t, S^t$ 
4   for  $s = 0, \dots, s_* - 1$  do /* On client */
5      $G_{U,c} \leftarrow \nabla_U \mathcal{L}_c(U_c^s S_c^s V_c^{s,\top})$ 
6      $G_{V,c} \leftarrow \nabla_V \mathcal{L}_c(U_c^s S_c^s V_c^{s,\top})$ 
7      $\tilde{U}_{c,-} \leftarrow \text{qr}([U_c^s \mid G_{U,c}])$ 
8      $\tilde{V}_{c,-} \leftarrow \text{qr}([V_c^s \mid G_{V,c}])$ 
9      $\tilde{S}_c = \tilde{U}_{c,-}^\top U_c^s S_c^s V_c^{s,\top} \tilde{V}_{c,-}$ 
10     $\tilde{S}_c^* \leftarrow \tilde{S}_c - \lambda \nabla_{\tilde{S}} \mathcal{L}_c(\tilde{U}_{c,-} \tilde{S}_c \tilde{V}_{c,-}^\top)$ 
11     $\tilde{S}^* \leftarrow \text{aggregate}(\{\tilde{S}_c^*\})$ 
12     $P_{r_1}, \Sigma_{r_1}, Q_{r_1} \leftarrow \text{svd}(\tilde{S}^*)$  with threshold  $\vartheta$  /* Compression step */
13     $U^{t+1} \leftarrow \tilde{U} P_{r_1}$ , and  $V^{t+1} \leftarrow \tilde{V} Q_{r_1}$  /* Basis projection */
14     $S^{t+1} \leftarrow \Sigma_{r_1}$ 

```

---

456 ization using mean  $[0.4914, 0.4822, 0.4465]$  and std. dev.  $[0.2470, 0.2435, 0.2616]$ . The test data  
 457 set is only normalized. The same augmentation is performed for CIFAR100, where with mean  
 458  $[0.5071, 0.4867, 0.4408]$  and std. dev.  $[0.2673, 0.2564, 0.2762]$ .

### 459 C.3 Additional computer vision results

460 **AlexNet on CIFAR10:** We train AlexNet on CIFAR10, where the fully connected head of the  
 461 network is replaced by a low-rank counterpart. A federated neural network setup with  $C$  clients  
 462 trains on  $CTs_*$  random batches of the dataset, that is the number of seen training data batches scales  
 463 with the client count. Figure 7 displays the validation accuracy of FeDLRT with variance correction  
 464 compared to FedLin, where one can see that the performance of FeDLRT mirrors the performance of  
 465 FedLin with more degrees of freedom. The measured validation accuracy peaks at  $C = 4$  clients in  
 466 both cases, where the higher number of seen training data-points offsets the negative effects of more  
 467 clients on the validation performance. All reported runs are within close distance of the non-federated,  
 468 full-rank baseline accuracy of 85.6%. Communication cost savings of the fully connected layers  
 469 amount between 96% and 97%<sup>6</sup> We observe, similarly to the results in Section 4.1, that the maximum  
 470 achieved communication cost savings, which depend on the layer ranks scales with the number of  
 471 clients  $C = 4$ , indicating that the decay rate of the singular values of the averaged coefficient matrix  
 472  $\tilde{S}^*$  depends on  $C$ .

473 **VGG16 on CIFAR10:** We train AlexNet on CIFAR10, where the fully connected head of the network  
 474 is replaced by a low-rank counterpart. A federated neural network setup with  $240/C$  local iterations  
 475 for  $C$  clients. Figure 8 displays the validation accuracy of FeDLRT with variance correction compared  
 476 to FedLin, where one can see that the performance of FeDLRT mirrors the performance of FedLin  
 477 with more degrees of freedom. All reported runs are within close distance of the non-federated,  
 478 full-rank baseline accuracy of 85.6%. Communication cost savings of the fully connected layers  
 479 amount between 96% and 97%<sup>7</sup> We observe, similarly results as in the ResNet18 test case.

480 **VGG16 on CIFAR10 with low-rank convolutions:** Mirroring the compute setup of the VGG16  
 481 test-case above, we now rewrite all convolutional layers of VGG16 as order 4 tensors in low-rank

<sup>6</sup>For clarity of exposition we consider only the fully connected layers. Taking into account the non low-rank convolution layers, the communication cost savings reduces to 87.5% to 87.3%.

<sup>7</sup>For clarity of exposition we consider only the fully connected layers. Taking into account the non low-rank convolution layers, the communication cost savings reduces to 87.5% to 87.3%.

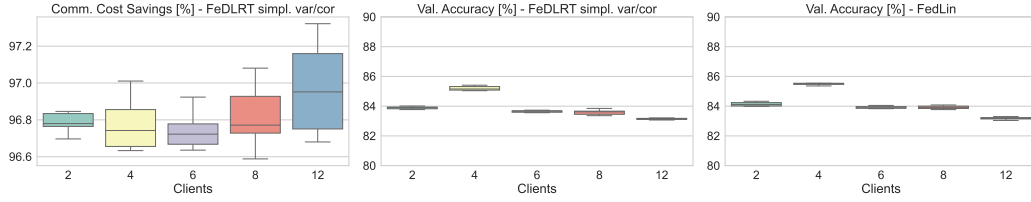


Figure 7: AlexNet CIFAR10 benchmark with fixed number of local iterations. (Left Panel) shows the savings in communication cost of simplified variance corrected FeDLRT vs FedLin. (Mid and right panel) compares the validation accuracy of FeDLRT and FedLin, where we see that FeDLRT behaves similarly to FedLin and achieves accuracy levels near the non-federated baseline value of 85.6%.

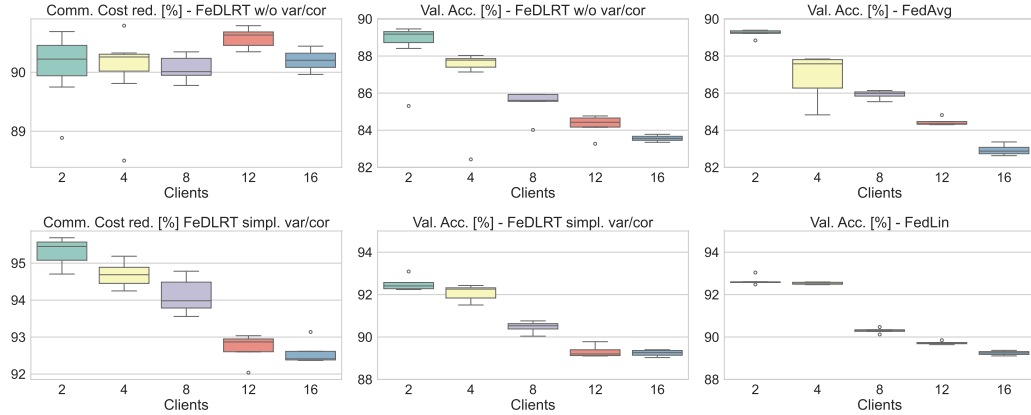


Figure 8: VGG16 CIFAR10 benchmark with  $240/C$  local iterations for  $C$  clients with simplified (lower row) and without (upper row) variance correction. (Left panel) show the savings in communication cost corresponding to FedLin at final time. (Mid and right panel top row) compares the validation accuracy of FeDLRT and FedAvg, where we see that FeDLRT behaves similarly to FedAvg, where higher  $C$  correlates with a drop in accuracy. FeDLRT with variance correction mitigates this issue and achieves similar performance as FedLin, close to the non-federated baseline accuracy is 93.15%.

482 Tucker format, as described in appendix B. The full-connected head of the network is treated with the  
 483 matrix low-rank method. The corresponding training results can be seen in Figure 9, and correspond  
 484 well with the previous results for VGG16. The reduction of communication cost is slightly higher,  
 485 due to the compression of the convolutions.

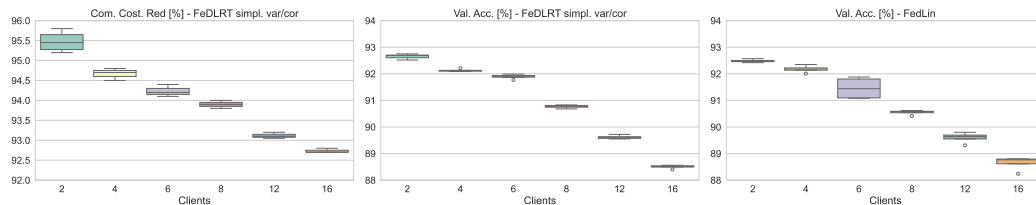


Figure 9: VGG16 CIFAR10, low-rank convolutional layers and low-rank fully connected layers. We report the communication cost savings and the validation accuracy of VGG16 with FeDLRT applied to training convolution and classifier layers. 2D convolutions are interpreted as an order-4 tensor and factorized in the Tucker format. The statistics over five random network initializations are reported using the training hyperparameters of Table 2 of the main manuscript. The results are similar to Fig. 7 in the main manuscript, where only the classifier is compressed. Remark that here the classifier contains most of the network parameters.



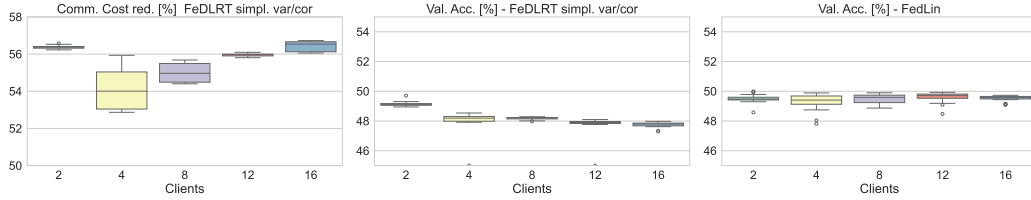


Figure 10: ViT CIFAR100 benchmark. (Left Panel) shows the savings in communication cost of variance corrected FeDLRT vs FedLin. (Mid and right panel) compares the validation accuracy of FeDLRT and FedLin, where we see that FeDLRT behaves similarly to FedLin and achieves accuracy levels near the non-federated baseline value of 50%, which is similar to literature results [46].

Table 2: Experimental setup object detection benchmarks. All test cases use a cosine annealing learning rate scheduler.

	Alexnet/Cifar10	ResNet18/Cifar10	VGG16/Cifar10	ViT/Cifar100
Batch size	128	128	128	256
Start Learningrate	$1e-2$	$1e-3$	$1e-2$	$3e-4$
End Learningrate	$1e-5$	$5e-4$	$5e-4$	$1e-5$
Aggregation Rounds	200	200	200	200
Local Iterations	100	$240/C$	$240/C$	$240/C$
Truncation tolerance $\tau$	0.01	0.01	0.01	0.01
Momentum	0.0	0.9	0.1	n.a.
Weight Decay	$1e-4$	$1e-3$	$1e-4$	$1e-2$
Optimizer	SGD	SGD	SGD	Adam w/ std pytorch parameters

486 **Vision Transformer on CIFAR100:** We consider a small vision transformer for CIFAR100, with  
 487 6 attention layers with 2 heads each followed by a ResNet block and a drop-out layer, all with  
 488 weight matrices of dimension  $512 \times 512$ . The tokenizer takes patches of size 8 with embedding  
 489 dimension 512. Training hyperparameters are given in Table 2. Remark that we do not aim for SOTA  
 490 performance, since transformer architectures are notoriously difficult to compress with low-rank  
 491 approaches, but rather compare the performance of FedLin to FeDLRT for a given compute budget.  
 492 We use  $s_* = 240/C$  local iterations for  $C$  clients. Observe in Figure 10 that FeDLRT achieves  
 493 similar performance as ViT with over 55% communication cost savings on average.

## 494 D Notation overview for the numerical analysis

495 We establish a set of notations to simplify the notation in the proofs

- 496 •  $\mathcal{L}_c(W)$  denotes the local loss function based on dataset  $X_c$  at client  $c$ .
- 497 •  $\mathcal{L}(W) = \frac{1}{C} \sum_{c=1}^C \mathcal{L}_c(W)$  is the global loss function.
- 498 •  $F_c(W) = -\nabla_W \mathcal{L}_c(W)$  is the negate of local loss gradient.
- 499 •  $F(W) = \frac{1}{C} \sum_{c=1}^C F_c(W)$  is the negate of global loss gradient.
- 500 •  $\mathcal{M}_r = \{W \in \mathbb{R}^{n \times n} : \text{rank}(W) = r\}$  is a manifold of rank  $r$  matrices.
- 501 •  $W_r = USV^\top \in \mathcal{M}_r$  is a rank- $r$  approximation of a matrix  $W$ .
- 502 •  $\mathcal{T}_{W_r} \mathcal{M}_r$  is the tangent space of  $\mathcal{M}_r$  at  $W_r$ .
- 503 •  $P(W_r)$  is the orthogonal projection onto  $\mathcal{T}_{W_r} \mathcal{M}_r$ .
- 504 •  $P_U = UU^\top$  is the orthogonal projection onto the range of orthonormal  $U \in \mathbb{R}^{n \times r}$ .
- 505 •  $P_V = VV^\top$  is the orthogonal projection onto the range of orthonormal  $V \in \mathbb{R}^{n \times r}$ .
- 506 • When applied to vectors,  $\|\cdot\|$  denotes the Euclidean norm ( $\ell_2$ -norm). When applied to matrices,  $\|\cdot\|$
- 507 denotes the Frobenius norm.

## 508 E Efficient basis gradient dynamics for basis augmentation

509 We first consider the basis update & Galerkin splitting scheme of (5). The splitting performs a  
510 reparametrization of the form  $K(t) = U(t)S(t)$  and  $L(t) = V(t)S(t)^\top$ . The basis update then reads

$$\begin{aligned} \dot{K} &= -\nabla_K \mathcal{L}(K(t)V_0^\top) \in \mathbb{R}^{n \times r}, & K(0) &= U_0 S_0, \\ \dot{L} &= -\nabla_L \mathcal{L}(U_0 L(t)^\top) \in \mathbb{R}^{n \times r}, & L(0) &= V_0 S_0^\top. \end{aligned} \quad (15)$$

511 Given the solution  $K(t_1)$  and  $L(t_1)$  at time  $t_1$ , the bases  $U_0$  and  $V_0$  are augmented by the orthonor-  
512 malization of the new directions  $K(t_1)$  and  $L(t_1)$ , i.e.

$$\begin{aligned} \tilde{U}R &= \text{qr}([U_0 \mid K(t_1)]) \in \mathbb{R}^{n \times 2r}, \\ \text{and } \tilde{V}R &= \text{qr}([V_0 \mid L(t_1)]) \in \mathbb{R}^{n \times 2r}, \end{aligned} \quad (16)$$

513 where  $R$  is the right factor of the respective QR decomposition and can be discarded. The initial  
514 condition of the coefficient update is  $\tilde{S}(t_0)$  projected onto the new bases, i.e.,

$$\dot{\tilde{S}} = -\nabla_S \mathcal{L}(\tilde{U}\tilde{S}(t)\tilde{V}^\top), \quad \tilde{S}(0) = \tilde{U}^\top U_0 \tilde{S}(0) V_0^\top \tilde{V}. \quad (17)$$

515 After the integration of the coefficient dynamics above, the redundant basis functions are typically  
516 truncated via an SVD of  $\tilde{S}$  ensuring that  $\tilde{S}$  is always full rank. In its continuous form above, the  
517 splitting yields a robust integrator for the projected gradient flow, without manifold dependent  
518 step-size restrictions:

519 **Theorem 5.** ([32]) *Assume  $\mathcal{L}$  is  $L$ -smooth with constant  $L$ , and locally bounded by  $B$ . Let  $W_r(t)$   
520 be the low-rank continuous time solution of (15) and (17) and let  $W(t)$  be the full rank solution at  
521  $t = 0$ . Assume the  $K, L$ , and  $S$  equations are integrated exactly from time  $t = 0$  to  $\Delta t$ . Assume that  
522 for any  $Y \in \mathcal{M}_r$  sufficiently close to  $W_r(t)$  the gradient  $F(Y)$  is  $\epsilon$  close to  $\mathcal{M}_r$ . Then*

$$\|W(\Delta t) - W_r(\Delta t)\| \leq d_1 \epsilon + d_2 \Delta t + d_3 \frac{\vartheta}{\Delta t},$$

523 where  $d_1, d_2, d_3$  depend only on  $L$  and  $B$ .

524 The theorem guarantees, that the low-rank representation does not imply any step-size restrictions on  
525 the optimization scheme. This is in stark contrast to a naive alternating descent optimization of the  
526 low-rank factors  $U, S, V$ .

527 To build an discretized numerical optimizer in a resource constrained federated scenario from the  
528 above continuous splitting equations, we avoid the reparametrization, which implies a 200% memory  
529 cost increase on the client side, since three versions of the low-rank layer need to be tracked.

530 **Lemma 2.** *Let  $USV \in \mathcal{M}_r$  be a low rank factorization that follows the projected gradient (5) flow  
531 using the splitting scheme (15) with  $K = US$  and  $V = VS^\top$ . Further, assume that equations for the*

532  $K$  and  $L$  factors are solved by an explicit Euler time integration with learning rate  $\lambda$ , i.e.

$$\begin{aligned} K(t_1) &= K(0) - \lambda \nabla_K \mathcal{L}(K(0)V_0^\top), & K(0) &= U_0 S_0, \\ L(t_1) &= L(0) - \lambda \nabla_L \mathcal{L}(U_0 L(0)^\top), & L(0) &= V_0 S_0^\top. \end{aligned} \quad (18)$$

533 Then, the basis augmentation (16) can be expressed as

$$\begin{aligned} \tilde{U}R &= \text{qr}([U_0 \mid -\nabla_U \mathcal{L}(U_0 S_0 V_0^\top)]) \in \mathbb{R}^{n \times 2r}, \\ \text{and } \tilde{V}R &= \text{qr}([V_0 \mid -\nabla_V \mathcal{L}(U_0 S_0 V_0^\top)]) \in \mathbb{R}^{n \times 2r}. \end{aligned} \quad (19)$$

534 and maintains the structure of the basis update and Galerkin operator split.

535 *Proof.* We consider the proof for the  $K$  equation and the  $U$  basis; the proof for  $L$  and  $V$  follows  
536 analogously.

537 Considering (16), we obtain with the explicit Euler discretization (18),

$$\begin{aligned} \text{span}([U_0 \mid K(t_1)]) &= \text{span}([U_0 \mid U_0 - \lambda \nabla_K \mathcal{L}(K(0)V_0^\top)]) \\ &= \text{span}([U_0 \mid -\lambda \nabla_K \mathcal{L}(K(0)V_0^\top)]) \\ &= \text{span}([U_0 \mid -\nabla_K \mathcal{L}(K(0)V_0^\top)]). \end{aligned} \quad (20)$$

538 Next, consider the continuous time dynamics of  $\dot{K}$ , where we omit explicit time dependence on  
539  $U, S, V$  and  $K$  for the sake of brevity, i.e.,

$$\begin{aligned} \dot{K} &= (\dot{U}S) \\ &= \dot{U}S + U\dot{S} \\ &\stackrel{(5)}{=} -(I - UU^\top) \nabla_W \mathcal{L}(USV^\top) V S^{-1} S - UU^\top \nabla_W \mathcal{L}(USV^\top) V \\ &= -(I - P_U) \nabla_W \mathcal{L}(USV^\top) V - P_U \nabla_W \mathcal{L}(USV^\top) V \\ &= (P_U - I) \nabla_W \mathcal{L}(USV^\top) V - P_U \nabla_W \mathcal{L}(USV^\top) V \\ &= -\nabla_W \mathcal{L}(USV^\top) V \end{aligned} \quad (21)$$

540 Further, using the chain rule, we observe

$$\nabla_U \mathcal{L}(USV^\top) = \nabla_W \mathcal{L}(USV^\top) \nabla_U (USV^\top) = \nabla_W \mathcal{L}(USV^\top) V S^\top$$

541 Thus,  $-\nabla_U \mathcal{L}(USV^\top) S^{-\top} = -\nabla_W \mathcal{L}(USV^\top) V = \dot{K}$ . Full rankness of  $S$  and (21) yield that  
542  $\text{span}(-\nabla_U \mathcal{L}(USV^\top)) = \text{span}(\dot{K})$ . Together with (20) this yields the proof.  $\square$

543 Lemma 2 adopts a more general result for Tucker tensors in an unpublished manuscript and simplifies  
544 the analysis for the matrix case considered here.

## 545 **F Efficient basis and coefficient communication**

546 Note that we have by orthogonality of the bases  $\tilde{U} = [U, \bar{U}]$  with  $\bar{U} \in \mathbb{R}^{n \times r}$  and  $\bar{U}^\top U = 0$  and  
547  $\tilde{V} = [V, \bar{V}]$  with  $\bar{V} \in \mathbb{R}^{n \times r}$  and  $\bar{V}^\top V = 0$ .

548 *Proof.* (Lemma 1) The basis augmented basis  $[U, G_U]$  before orthonormalization already contains  
549 the orthonormal vectors given by the columns of  $U$ . A QR decomposition therefor only rearranges  
550 the columns of  $G_U$  such that  $\tilde{U} = [U, \bar{U}]$  with  $\bar{U} \in \mathbb{R}^{n \times r}$  and  $\bar{U}^\top U = 0$ . The analogous result holds  
551 for  $\tilde{V} = [V, \bar{V}]$ . The projection onto the augmented basis therefore reads

$$\tilde{U}^\top U = \begin{bmatrix} U^\top U \\ \bar{U}^\top U \end{bmatrix} = \begin{bmatrix} I \\ 0 \end{bmatrix} \quad \text{and} \quad \tilde{V}^\top V = \begin{bmatrix} V^\top V \\ \bar{V}^\top V \end{bmatrix} = \begin{bmatrix} I \\ 0 \end{bmatrix}. \quad (22)$$

552 Consequently, the augmented coefficient matrix takes the form

$$\tilde{S} = \tilde{U}^\top U S V^\top \tilde{V} = \begin{bmatrix} S & 0 \\ 0 & 0 \end{bmatrix}. \quad (23)$$

553  $\square$

554 **G Analysis for FeDLRT with full variance correction**

555 In this section we establish bounds on the coefficient drift of the FeDLRT method with full variance  
 556 correction. We use the established coefficient drift bound to derive a loss-descend guarantee. The  
 557 strategy of our analysis follows the one of FedLin [24]. We first state an auxiliary lemma.

558 **Lemma 3.** *Let  $U \in \mathbb{R}^{n \times r}$  and  $V \in \mathbb{R}^{n \times r}$  be orthonormal matrices. Let  $F$  be an  $L$ -continuous  
 559 function. Then, for  $S_1, S_2 \in \mathbb{R}^{r \times r}$ ,*

$$\|P_U (F(US_1V^\top) - F(US_2V^\top)) P_V\| \leq L \|S_1 - S_2\| \quad (24)$$

560 and

$$\|U (F(US_1V^\top) - F(US_2V^\top)) V^\top\| \leq L \|S_1 - S_2\|, \quad (25)$$

561 where  $P_U$  and  $P_V$  are orthogonal projections defined in Appendix D.

562 *Proof.* For the first statement, consider

$$\begin{aligned} & \|P_U (F(US_1V^\top) - F(US_2V^\top)) P_V\| \\ &= \|UU^\top (F(US_1V^\top) - F(US_2V^\top)) VV^\top\| \\ &\stackrel{(I)}{\leq} \|U\| \|U^\top\| \|F(US_1V^\top) - F(US_2V^\top)\| \|V\| \|V^\top\| \\ &\stackrel{(II)}{=} \|F(US_1V^\top) - F(US_2V^\top)\| \\ &\stackrel{(III)}{\leq} L \|US_1V^\top - US_2V^\top\| = L \|U(S_1 - S_2)V^\top\| \\ &\stackrel{(I)}{\leq} L \|U\| \|S_1 - S_2\| \|V^\top\| \\ &\stackrel{(II)}{=} L \|S_1 - S_2\|, \end{aligned}$$

563 where we have used in (I) the operator norm inequality of the Frobenius norm, in (II) orthonormality  
 564 of  $U$ ,  $V$ , and in (III)  $L$ -continuity of  $F$ . The second statement is proven analogously.  $\square$

565 **G.1 Coefficient drift bound for FeDLRT with full variance correction**

566 We consider the FeDLRT method with variance correction, see Algorithm 1. Key difference to the  
 567 FeDLRT method without variance correction is the modified coefficient update, incorporating global  
 568 gradient information of the augmented coefficient matrix  $\tilde{S}$  and local, stale gradient information  
 569 of the augmented coefficient matrix  $\tilde{S}_c$ . The variance corrected local coefficient update (8) can be  
 570 expressed in terms of the projected Riemannian gradient as

$$\tilde{S}_c^{s+1} = \tilde{S}_c^s + \lambda \tilde{U}^\top \left( F_c(\tilde{W}_{r,c}^s) - F_c(\tilde{W}_r) + F(\tilde{W}_r) \right) \tilde{V}, \quad (26)$$

571 where  $\tilde{U}^\top F_c(\tilde{W}_{r,c}^s) \tilde{V} = \nabla_{\tilde{S}_c} \mathcal{L}_c(\tilde{U} \tilde{S}_c^s \tilde{V})$ ,  $\tilde{U}^\top F_c(\tilde{W}_{r,c}) \tilde{V} = \nabla_{\tilde{S}_c} \mathcal{L}_c(\tilde{U} \tilde{S}_c^{s=0} \tilde{V})$  and  
 572  $\tilde{U}^\top F_c(\tilde{W}_{r,c}) \tilde{V} = \nabla_{\tilde{S}_c} \mathcal{L}(\tilde{U} \tilde{S}_c^s \tilde{V})$ . Recall that  $\tilde{S} = \tilde{S}_c$  for  $s = 0$ .

573 We provide proof for Theorem 1 to bound the drift term  $\|\tilde{S}_c^s - \tilde{S}_c\|$ . We restate this theorem to the  
 574 Riemannian notation and restate it below.

575 **Theorem 6.** *(Restatement of Theorem 1) Given augmented basis and coefficient matrices  $\tilde{U}$ ,  $\tilde{V}$ , and  
 576  $\tilde{S}$ , and  $\tilde{W}_r = \tilde{U} \tilde{S} \tilde{V}^\top$ . If the local learning rate  $0 < \lambda \leq \frac{1}{L s_*}$  with  $s_* \geq 1$  the number of local steps,  
 577 for all clients  $c$ ,*

$$\|\tilde{S}_c^s - \tilde{S}_c\| \leq \exp(1) s_* \lambda \left\| \tilde{U}^\top F(\tilde{W}_r) \tilde{V} \right\|, \quad \text{for } s = 1, \dots, s_* - 1, \quad (27)$$

578 where  $\tilde{S}_c^s$  is the variance corrected coefficient as given in (8).

579 *Proof.* From the adjusted coefficient update in (26), we get

$$\begin{aligned}
\left\| \tilde{S}_c^{s+1} - \tilde{S}_c \right\| &= \left\| \tilde{S}_c^s - \tilde{S}_c + \lambda \tilde{U}^\top \left( F_c(\tilde{W}_{r,c}^s) - F_c(\tilde{W}_r) + F(\tilde{W}_r) \right) \tilde{V} \right\| \\
&\leq \left\| \tilde{S}_c^s - \tilde{S}_c \right\| + \lambda \left\| \tilde{U}^\top \left( F_c(\tilde{W}_{r,c}^s) - F_c(\tilde{W}_r) \right) \tilde{V} \right\| + \lambda \left\| \tilde{U}^\top F(\tilde{W}_r) \tilde{V} \right\| \\
&\stackrel{(I)}{\leq} \left\| \tilde{S}_c^s - \tilde{S}_c \right\| + \lambda L \left\| \tilde{S}_c^s - \tilde{S} \right\| + \lambda \left\| \tilde{U}^\top F(\tilde{W}_r) \tilde{V} \right\| \\
&\leq (1 + \lambda L) \left\| \tilde{S}_c^s - \tilde{S} \right\| + \lambda \left\| \tilde{U}^\top F(\tilde{W}_r) \tilde{V} \right\| \\
&\leq \left( 1 + \frac{1}{s_*} \right) \left\| \tilde{S}_c^s - \tilde{S} \right\| + \lambda \left\| \tilde{U}^\top F(\tilde{W}_r) \tilde{V} \right\|.
\end{aligned}$$

580 We use in (I) Lemma 3 Recursively plugging in the above inequality yields for  $a = (1 + \frac{1}{s_*})$

$$\begin{aligned}
\left\| \tilde{S}_c^{s+1} - \tilde{S}_c \right\| &\leq a^{s+1} \left\| \tilde{S}_c^{s=0} - \tilde{S} \right\| + \left( \sum_{j=0}^s a^j \right) \lambda \left\| \tilde{U}^\top F(\tilde{W}_r) \tilde{V} \right\| \\
&= \left( \sum_{j=0}^s a^j \right) \lambda \left\| \tilde{U}^\top F(\tilde{W}_r) \tilde{V} \right\| \\
&= \frac{a^{s+1} - 1}{a - 1} \lambda \left\| \tilde{U}^\top F(\tilde{W}_r) \tilde{V} \right\| \\
&\leq \left( 1 + \frac{1}{s_*} \right)^{s+1} s_* \lambda \left\| \tilde{U}^\top F(\tilde{W}_r) \tilde{V} \right\| \\
&\leq \left( 1 + \frac{1}{s_*} \right)^{s_*} s_* \lambda \left\| \tilde{U}^\top F(\tilde{W}_r) \tilde{V} \right\| \\
&\leq \exp(1) s_* \lambda \left\| \tilde{U}^\top F(\tilde{W}_r) \tilde{V} \right\|.
\end{aligned}$$

581

□

## 582 G.2 Global loss descend for FeDLRT with full variance correction

583 We first state a few auxiliary lemmas, which provide common inequalities that will be used in the  
584 following analysis.

585 **Lemma 4.** ([10, Lemma 5.2]) For any two matrices  $Y_1, Y_2 \in \mathbb{R}^{n \times n}$  and an  $L$ -smooth  $\mathcal{L}$  with constant  
586  $L$  it holds

$$\mathcal{L}(Y_1) - \mathcal{L}(Y_2) \leq -\langle Y_1 - Y_2, F(Y_2) \rangle + \frac{L}{2} \|Y_1 - Y_2\|^2, \quad (28)$$

587 where  $F(Y) = -\nabla_Y \mathcal{L}(Y)$ .

588 **Lemma 5.** ([25, Lemma 5]) For two vectors  $x_1, x_2 \in \mathbb{R}^d$  it holds for  $\gamma > 0$

$$\|x_1 + x_2\|^2 \leq (1 + \gamma) \|x_1\|^2 + \left( 1 + \frac{1}{\gamma} \right) \|x_2\|^2. \quad (29)$$

589 **Lemma 6.** ([25, Lemma 6]) For  $C$  vectors  $x_1, \dots, x_C \in \mathbb{R}^d$  the application of Jensen's inequality  
590 yields

$$\left\| \sum_{c=1}^C x_c \right\|^2 \leq C \sum_{c=1}^C \|x_c\|^2. \quad (30)$$

591 First, we consider the loss function value at the augmentation step.

592 **Lemma 7.** We have  $\mathcal{L}(\tilde{W}_r) = \mathcal{L}(W_r^t)$  for the loss before and after basis augmentation.

593 *Proof.* Due to Lemma 1,  $\tilde{S} = \begin{bmatrix} S^t & 0 \\ 0 & 0 \end{bmatrix}$ , thus  $\tilde{W}_r = \tilde{U}\tilde{S}\tilde{V}^\top = USV^\top = W^t$ .  $\square$

594 We next bound the loss descent between the augmentation step and the truncation step - having  
595 performed the aggregation of the client updates.

596 **Theorem 7.** Let  $\tilde{W}_r = \tilde{U}\tilde{S}\tilde{V}^\top$  be the augmented factorization at global iteration  $t$  and let  $\tilde{W}_r^* =$   
597  $\tilde{U}\tilde{S}^*\tilde{V}^\top$  be the aggregated solution after client iterations, i.e.,  $\tilde{S}^* = \frac{1}{C} \sum_{c=1}^C \tilde{S}_c^{s_*}$ . Then the variance  
598 corrected coefficient update (26) yields the guarantee

$$\begin{aligned} \mathcal{L}(\tilde{W}_r^*) - \mathcal{L}(\tilde{W}_r) &\leq -(s_*\lambda)(1 - (s_*\lambda)L) \left\| \tilde{U}^\top F(\tilde{W}_r) \tilde{V} \right\|^2 \\ &\quad + \left( \frac{L\lambda}{C} \sum_{c=1}^C \sum_{s=0}^{s_*-1} \left\| \tilde{S}_c^s - \tilde{S} \right\| \right) \left\| \tilde{U}^\top F(\tilde{W}_r) \tilde{V} \right\| \\ &\quad + \frac{L^3\lambda^2 s_*}{C} \sum_{c=1}^C \sum_{s=0}^{s_*-1} \left\| \tilde{S}_c^s - \tilde{S}_c \right\|^2. \end{aligned} \quad (31)$$

599 *Proof.* From (8),  $P_{\tilde{U}} = \tilde{U}\tilde{U}^\top$ ,  $P_{\tilde{V}} = \tilde{V}\tilde{V}^\top$ , and the fact that  $\tilde{W}_{r,c}^{s=0} = \tilde{W}_r$  for all  $c = 1, \dots, C$ ,

$$\begin{aligned} \tilde{W}_{r,c}^{s_*} &= \tilde{U}\tilde{S}_c^{s_*}\tilde{V}^\top = \tilde{U}\tilde{S}_c^{s=0}\tilde{V}^\top + \tilde{U}\tilde{U}^\top \sum_{s=0}^{s_*-1} \lambda \left( F_c(\tilde{W}_{r,c}^s) - F_c(\tilde{W}_r) + F(\tilde{W}_r) \right) \tilde{V}\tilde{V}^\top \\ &= \tilde{W}_r - \lambda \sum_{s=0}^{s_*-1} P_{\tilde{U}} F_c(\tilde{W}_{r,c}^s) P_{\tilde{V}} - \lambda P_{\tilde{U}} \left( F(\tilde{W}_r) - F_c(\tilde{W}_r) \right) P_{\tilde{V}}. \end{aligned}$$

600 Averaging across clients leads to

$$\begin{aligned} \tilde{W}_r^* &= \frac{1}{C} \sum_{c=1}^C \tilde{W}_{r,c}^{s_*} = \tilde{W}_r - \frac{\lambda}{C} \sum_{c=1}^C \sum_{s=0}^{s_*-1} P_{\tilde{U}} F_c(\tilde{W}_{r,c}^s) P_{\tilde{V}} - \frac{\lambda}{C} \sum_{c=1}^C P_{\tilde{U}} \left( F(\tilde{W}_r) - F_c(\tilde{W}_r) \right) P_{\tilde{V}} \\ &= \tilde{W}_r - \frac{\lambda}{C} \sum_{c=1}^C \sum_{s=0}^{s_*-1} P_{\tilde{U}} F_c(\tilde{W}_{r,c}^s) P_{\tilde{V}}, \end{aligned} \quad (32)$$

601 where we have used the definition of the global and local gradient at  $\tilde{W}_r$ , i.e.,  $\frac{1}{C} \sum_{c=1}^C F_c(\tilde{W}_r) =$   
602  $F(\tilde{W}_r)$ . Based on  $L$ -continuity of  $F$  and  $F_c$ , (32), and Lemma 4, we obtain further

$$\begin{aligned} \mathcal{L}(\tilde{W}_r^*) - \mathcal{L}(\tilde{W}_r) &\leq \left\langle \tilde{W}_r^* - \tilde{W}_r, F(\tilde{W}_r) \right\rangle + \frac{L}{2} \left\| \tilde{W}_r^* - \tilde{W}_r \right\|^2 \\ &= - \left\langle \frac{\lambda}{C} \sum_{c=1}^C \sum_{s=0}^{s_*-1} P_{\tilde{U}} F_c(\tilde{W}_{r,c}^s) P_{\tilde{V}}, F(\tilde{W}_r) \right\rangle + \frac{L}{2} \left\| \frac{\lambda}{C} \sum_{c=1}^C \sum_{s=0}^{s_*-1} P_{\tilde{U}} F_c(\tilde{W}_{r,c}^s) P_{\tilde{V}} \right\|^2. \end{aligned} \quad (33)$$

603 Next, we bound each of the two right-hand-side terms separately. We first express the first term as

$$\begin{aligned}
& - \left\langle \frac{\lambda}{C} \sum_{c=1}^C \sum_{s=0}^{s_*-1} P_{\tilde{U}} F_c(\tilde{W}_{r,c}^s) P_{\tilde{V}}, F(\tilde{W}_r) \right\rangle \\
& = - \left\langle \frac{\lambda}{C} \sum_{c=1}^C \sum_{s=0}^{s_*-1} P_{\tilde{U}} \left( F_c(\tilde{W}_{r,c}^s) - F_c(\tilde{W}_r) \right) P_{\tilde{V}} + P_{\tilde{U}} \left( \frac{\lambda}{C} \sum_{c=1}^C \sum_{s=0}^{s_*-1} F_c(\tilde{W}_r) \right) P_{\tilde{V}}, F(\tilde{W}_r) \right\rangle \\
& = - \left\langle \frac{\lambda}{C} \sum_{c=1}^C \sum_{s=0}^{s_*-1} P_{\tilde{U}} \left( F_c(\tilde{W}_{r,c}^s) - F_c(\tilde{W}_r) \right) P_{\tilde{V}} + P_{\tilde{U}} \frac{s_* \lambda}{C} \sum_{c=1}^C F_c(\tilde{W}_r) P_{\tilde{V}}, F(\tilde{W}_r) \right\rangle \\
& = - \left\langle P_{\tilde{U}} \left( \frac{\lambda}{C} \sum_{c=1}^C \sum_{s=0}^{s_*-1} F_c(\tilde{W}_{r,c}^s) - F_c(\tilde{W}_r) \right) P_{\tilde{V}} + P_{\tilde{U}} s_* \lambda F(\tilde{W}_r) P_{\tilde{V}}, F(\tilde{W}_r) \right\rangle \\
& = - \left\langle \tilde{U}^\top \left( \frac{\lambda}{C} \sum_{c=1}^C \sum_{s=0}^{s_*-1} F_c(\tilde{W}_{r,c}^s) - F_c(\tilde{W}_r) \right) \tilde{V}, \tilde{U}^\top F(\tilde{W}_r) \tilde{V}^\top \right\rangle - s_* \lambda \left\langle \tilde{U}^\top F(\tilde{W}_r) \tilde{V}, \tilde{U}^\top F(\tilde{W}_r) \tilde{V} \right\rangle \\
& = - \left\langle \frac{\lambda}{C} \sum_{c=1}^C \sum_{s=0}^{s_*-1} \tilde{U}^\top \left( F_c(\tilde{W}_{r,c}^s) - F_c(\tilde{W}_r) \right) \tilde{V}, \tilde{U}^\top F(\tilde{W}_r) \tilde{V} \right\rangle - s_* \lambda \left\| \tilde{U}^\top F(\tilde{W}_r) \tilde{V} \right\|^2,
\end{aligned}$$

604 where the definitions of  $P_{\tilde{U}}$  and  $P_{\tilde{V}}$  are used. Following this, the first term then can be bounded by

$$\begin{aligned}
& - \left\langle \frac{\lambda}{C} \sum_{c=1}^C \sum_{s=0}^{s_*-1} P_{\tilde{U}} F_c(\tilde{W}_{r,c}^s) P_{\tilde{V}}, F(\tilde{W}_r) \right\rangle \\
& \leq \frac{\lambda}{C} \sum_{c=1}^C \sum_{s=0}^{s_*-1} \left\| \tilde{U}^\top \left( F_c(\tilde{W}_{r,c}^s) - F_c(\tilde{W}_r) \right) \tilde{V} \right\| \left\| \tilde{U}^\top F(\tilde{W}_r) \tilde{V} \right\| - s_* \lambda \left\| \tilde{U}^\top F(\tilde{W}_r) \tilde{V} \right\|^2 \\
& \leq \frac{L\lambda}{C} \sum_{c=1}^C \sum_{s=0}^{s_*-1} \left\| \tilde{S}_c^s - \tilde{S} \right\| \left\| \tilde{U}^\top F(\tilde{W}_r) \tilde{V} \right\| - s_* \lambda \left\| \tilde{U}^\top F(\tilde{W}_r) \tilde{V} \right\|^2,
\end{aligned}$$

605 where Lemma 3 is invoked in the last inequality. Following a similar approach, we express the second  
606 term as

$$\frac{L}{2} \left\| \frac{\lambda}{C} \sum_{c=1}^C \sum_{s=0}^{s_*-1} P_{\tilde{U}} F_c(\tilde{W}_{r,c}^s) P_{\tilde{V}} \right\|^2 = \frac{L}{2} \left\| \frac{\lambda}{C} \sum_{c=1}^C \sum_{s=0}^{s_*-1} P_{\tilde{U}} \left( F_c(\tilde{W}_{r,c}^s) - F_c(\tilde{W}_r) \right) P_{\tilde{V}} + s_* \lambda P_{\tilde{U}} F(\tilde{W}_r) P_{\tilde{V}} \right\|^2,$$

607 which can be bounded by

$$\begin{aligned}
& \frac{L}{2} \left\| \frac{\lambda}{C} \sum_{c=1}^C \sum_{s=0}^{s_*-1} P_{\tilde{U}} F_c(\tilde{W}_{r,c}^s) P_{\tilde{V}} \right\|^2 \\
& \stackrel{(I)}{\leq} L \left\| \frac{\lambda}{C} \sum_{c=1}^C \sum_{s=0}^{s_*-1} P_{\tilde{U}} \left( F_c(\tilde{W}_{r,c}^s) - F_c(\tilde{W}_r) \right) P_{\tilde{V}} \right\|^2 + (s_* \lambda)^2 L \left\| P_{\tilde{U}} F(\tilde{W}_r) P_{\tilde{V}} \right\|^2 \\
& \stackrel{(II)}{\leq} \frac{L}{C} \sum_{c=1}^C \lambda^2 s_*^2 \sum_{s=0}^{s_*-1} \left\| P_{\tilde{U}} \left( F_c(\tilde{W}_{r,c}^s) - F_c(\tilde{W}_r) \right) P_{\tilde{V}} \right\|^2 + (s_* \lambda)^2 L \left\| P_{\tilde{U}} F(\tilde{W}_r) P_{\tilde{V}} \right\|^2 \\
& \stackrel{(III)}{\leq} \frac{L^3 \lambda^2 s_*^2}{C} \sum_{c=1}^C \sum_{s=0}^{s_*-1} \left\| \tilde{S}_c^s - \tilde{S} \right\|^2 + (s_* \lambda)^2 L \left\| P_{\tilde{U}} F(\tilde{W}_r) P_{\tilde{V}} \right\|^2 \\
& \stackrel{(IV)}{\leq} \frac{L^3 \lambda^2 s_*^2}{C} \sum_{c=1}^C \sum_{s=0}^{s_*-1} \left\| \tilde{S}_c^s - \tilde{S} \right\|^2 + (s_* \lambda)^2 L \left\| \tilde{U}^\top F(\tilde{W}_r) \tilde{V} \right\|^2,
\end{aligned}$$

608 where Lemma 5 with  $\gamma = 1$  is used in (I), Jensen's inequality is used in (II), Lemma 3 is used in  
609 in (III), and (IV) follows from the Operator norm inequality of the Frobenius norm in combination  
610 with orthonormality of  $U$  and  $V^\top$ .

611 Plugging these two bounds into (33) gives

$$\begin{aligned}
\mathcal{L}(\widetilde{W}_r^*) - \mathcal{L}(\widetilde{W}_r) &\leq - \left\langle \frac{\lambda}{C} \sum_{c=1}^C \sum_{s=0}^{s_*-1} P_{\widetilde{U}} F_c(\widetilde{W}_{r,c}^s) P_{\widetilde{V}}, F(\widetilde{W}_r) \right\rangle + \frac{L}{2} \left\| \frac{\lambda}{C} \sum_{c=1}^C \sum_{s=0}^{s_*-1} P_{\widetilde{U}} F_c(\widetilde{W}_{r,c}^s) P_{\widetilde{V}} \right\|^2 \\
&\leq \frac{L\lambda}{C} \sum_{c=1}^C \sum_{s=0}^{s_*-1} \left\| \widetilde{S}_c^s - \widetilde{S} \right\| \left\| \widetilde{U}^\top F(\widetilde{W}_r) \widetilde{V} \right\| - s_* \lambda \left\| \widetilde{U}^\top F(\widetilde{W}_r) \widetilde{V} \right\|^2 \\
&\quad + \frac{L^3 \lambda^2 s_*}{C} \sum_{c=1}^C \sum_{s=0}^{s_*-1} \left\| \widetilde{S}_c^s - \widetilde{S}_c \right\|^2 + (s_* \lambda)^2 L \left\| \widetilde{U}^\top F(\widetilde{W}_r) \widetilde{V} \right\|^2 \\
&= - (s_* \lambda) (1 - (s_* \lambda) L) \left\| \widetilde{U}^\top F(\widetilde{W}_r) \widetilde{V} \right\|^2 \\
&\quad + \left( \frac{L\lambda}{C} \sum_{c=1}^C \sum_{s=0}^{s_*-1} \left\| \widetilde{S}_c^s - \widetilde{S} \right\| \right) \left\| \widetilde{U}^\top F(\widetilde{W}_r) \widetilde{V} \right\| \\
&\quad + \frac{L^3 \lambda^2 s_*}{C} \sum_{c=1}^C \sum_{s=0}^{s_*-1} \left\| \widetilde{S}_c^s - \widetilde{S}_c \right\|^2,
\end{aligned}$$

612 which concludes the proof.  $\square$

613 With this result, we next bound the loss descent between the augmentation and coefficient aggregation  
614 step in the following theorem.

615 **Theorem 8.** *Under the same assumptions as in Theorem 7. Let the local learning rate be  $0 < \lambda \leq$   
616  $\frac{1}{12Ls_*}$  with number of local iterations  $s_* \geq 1$ . Then,*

$$\mathcal{L}(\widetilde{W}_r^*) - \mathcal{L}(\widetilde{W}_r) \leq -s_* \lambda (1 - 12s_* \lambda L) \left\| \widetilde{U}^\top F(\widetilde{W}_r) \widetilde{V} \right\|^2. \quad (34)$$

617 *Proof.* Applying the drift bound given in Theorem 1 to the loss descent bound given by Theorem 7  
618 in (31) leads to

$$\begin{aligned}
&- (s_* \lambda) (1 - (s_* \lambda) L) \left\| \widetilde{U}^\top F(\widetilde{W}_r) \widetilde{V} \right\|^2 \\
&+ \left( \frac{L\lambda}{C} \sum_{c=1}^C \sum_{s=0}^{s_*-1} \left( \exp(1) s_* \lambda \left\| \widetilde{U}^\top F(\widetilde{W}_r) \widetilde{V} \right\| \right) \right) \left\| \widetilde{U}^\top F(\widetilde{W}_r) \widetilde{V} \right\| \\
&+ \frac{L^3 \lambda^2 s_*}{C} \sum_{c=1}^C \sum_{s=0}^{s_*-1} \left( \exp(1) s_* \lambda \left\| \widetilde{U}^\top F(\widetilde{W}_r) \widetilde{V} \right\| \right)^2 \\
&= - (s_* \lambda) (1 - (s_* \lambda) L) \left\| \widetilde{U}^\top F(\widetilde{W}_r) \widetilde{V} \right\|^2 + L \lambda^2 s_*^2 \exp(1) \left\| \widetilde{U}^\top F(\widetilde{W}_r) \widetilde{V} \right\|^2 \\
&\quad + L^3 \lambda^4 s_*^4 \exp(2) \left\| \widetilde{U}^\top F(\widetilde{W}_r) \widetilde{V} \right\|^2 \\
&= - (s_* \lambda) (1 - (s_* \lambda) L - (s_* \lambda) L \exp(1) - (s_* \lambda)^3 L^2 \exp(2)) \left\| \widetilde{U}^\top F(\widetilde{W}_r) \widetilde{V} \right\|^2 \\
&\leq - (s_* \lambda) (1 - (s_* \lambda) L (1 + \exp(1) + \exp(2))) \left\| \widetilde{U}^\top F(\widetilde{W}_r) \widetilde{V} \right\|^2 \\
&\leq - (s_* \lambda) (1 - 12(s_* \lambda) L) \left\| \widetilde{U}^\top F(\widetilde{W}_r) \widetilde{V} \right\|^2,
\end{aligned}$$

619 where we have used that  $(s_* \lambda) L \leq 1$  and that  $1 + \exp(1) + \exp(2) \approx 11.107 \leq 12$ .  $\square$

620 We are now prepared to prove Theorem 2, which we restate in terms of Riemannian gradients as  
621 below.

622 **Theorem 9.** *(Restatement of Theorem 2) Let  $U^t S^t V^{t,\top}$  and  $U^{t+1} S^{t+1} V^{t+1,\top}$  be the factorization  
623 before and after iteration  $t$  of Algorithm 1 with variance correction and singular value truncation*



624 threshold  $\vartheta$ . Let  $\mathcal{L}_c$  and  $\mathcal{L}$  be  $L$ -smooth with constant  $L$ , and let the local learning rate be  $0 \leq \lambda \leq$   
625  $\frac{1}{12Ls_*}$ . Then the global loss descent is bounded by

$$\mathcal{L}(U^{t+1}S^{t+1}V^{t+1,\top}) - \mathcal{L}(U^tS^tV^{t,\top}) \leq -(s_*\lambda)(1 - 12(s_*\lambda)L) \left\| \tilde{U}^\top F(\tilde{W}_r) \tilde{V} \right\|^2 + L\vartheta. \quad (35)$$

626 *Proof.* Consider  $\mathcal{L}(W_r^{t+1})$  and  $\mathcal{L}(\tilde{W}_r^*)$ , i.e., the loss values before and after the truncation step. By  
627 the mean value theorem, we obtain for some  $h \in [0, 1]$

$$\begin{aligned} \mathcal{L}(W_r^{t+1}) &= \mathcal{L}(\tilde{W}_r^*) + \left\langle -F(hW_r^{t+1} + (1-h)\tilde{W}_r^*), W_r^{t+1} - \tilde{W}_r^* \right\rangle \\ &\leq \mathcal{L}(\tilde{W}_r^*) + \left\| F(hW_r^{t+1} + (1-h)\tilde{W}_r^*) \right\| \left\| W_r^{t+1} - \tilde{W}_r^* \right\| \\ &\leq \mathcal{L}(\tilde{W}_r^*) + L\vartheta \end{aligned} \quad (36)$$

628 where  $L$ -smoothness and the fact that  $\vartheta \geq \left\| W_r^{t+1} - \tilde{W}_r^* \right\|$  are used in (II), where the latter follows  
629 from the singular value truncation threshold. Combining the above arguments with Lemma 7 and  
630 Theorem 8 yields

$$\begin{aligned} \mathcal{L}(W_r^{t+1}) - \mathcal{L}(W_r^t) &= (\mathcal{L}(W_r^{t+1}) - \mathcal{L}(\tilde{W}_r^*)) + (\mathcal{L}(\tilde{W}_r^*) - \mathcal{L}(\tilde{W}_r)) + (\mathcal{L}(\tilde{W}_r) - \mathcal{L}(W_r^t)) \\ &\leq L\vartheta - (s_*\lambda)(1 - 12(s_*\lambda)L) \left\| \tilde{U}^\top F(\tilde{W}_r) \tilde{V} \right\|^2, \end{aligned}$$

631 which concludes the proof.  $\square$

### 632 G.3 Global convergence of FeDLRT with full variance correction

633 **Theorem 10.** (Restatement of Theorem 3) Assume that  $\mathcal{L}$  is  $L$ -smooth with constant  $L$  for all  
634  $c = 1, \dots, C$ . Let  $\tilde{U}^t \tilde{S}^t \tilde{V}^{t,\top}$  be the augmented representation at iteration  $t$ . Then Algorithm 1  
635 guarantees for the learning rate  $\lambda \leq \frac{1}{12Ls_*}$  and final iteration  $T$

$$\min_{t=1, \dots, T} \left\| \nabla_{\tilde{S}} \mathcal{L}(U^t S^t V^{t,\top}) \right\|^2 \leq \frac{48L}{T} (\mathcal{L}(W_r^{t=1}) - \mathcal{L}(W_r^{t=T+1})) + 48L^2\vartheta. \quad (37)$$

636 *Proof.* Consider Theorem 2,

$$\mathcal{L}(W_r^{t+1}) - \mathcal{L}(W_r^t) \leq L\vartheta - (s_*\lambda)(1 - 12(s_*\lambda)L) \left\| \nabla_{\tilde{S}} \mathcal{L}(U^t S^t V^{t,\top}) \right\|^2, \quad (38)$$

637 and assume that  $\lambda s_* = \frac{1}{24L}$ , i.e.  $\lambda = \frac{1}{24Ls_*} \leq \frac{1}{Ls_*}$ , which obeys the learning rate requirement of  
638 Theorem 2. Plugging this learning rate into (38) gives

$$\left\| \nabla_{\tilde{S}} \mathcal{L}(U^t S^t V^{t,\top}) \right\|^2 \leq 48L (\mathcal{L}(W_r^t) - \mathcal{L}(W_r^{t+1}) + L\vartheta).$$

639 Averaging from  $t = 1$  to  $t = T$  yields

$$\begin{aligned} \min_{t=1, \dots, T} \left\| \nabla_{\tilde{S}} \mathcal{L}(U^t S^t V^{t,\top}) \right\|^2 &\leq \frac{1}{T} \sum_{t=1}^T \left\| \nabla_{\tilde{S}} \mathcal{L}(U^t S^t V^{t,\top}) \right\|^2 \\ &\leq \frac{48L}{T} (\mathcal{L}(W_r^{t=1}) - \mathcal{L}(W_r^{t=T+1})) + 48L^2\vartheta, \end{aligned}$$

640 which concludes the proof.  $\square$

641 We remark that for a general loss function, it is possible that a point with small gradient magnitude  
642 can be far from the stationary points. However, assuming that the loss function is locally strongly  
643 convex in a neighborhood of a stationary point, then the gradient magnitude can be used to bound  
644 the distance to this stationary point in the neighborhood. For further reference, we point to [?, Eq.  
645 (4.12)] for the estimate.

## 646 H Analysis for FeDLRT with simplified variance correction

647 We consider the FeDLRT method with simplified variance correction, see Algorithm 5. Key difference  
 648 to the standard FeDLRT with full variance correction, see Algorithm 1 is the modified coefficient  
 649 update, incorporating global gradient information of the non-augmented coefficient matrix  $S$  for the  
 650 variance correction term, that is

$$\check{V}_c = \check{G}_{\tilde{S}} - \check{G}_{\tilde{S},c} = \begin{bmatrix} \nabla_S \mathcal{L}(U^t S^t V^{t,\top}) - \nabla_S \mathcal{L}_c(U^t S^t V^{t,\top}) & 0 \\ 0 & 0 \end{bmatrix}. \quad (39)$$

651 Using the Riemmanian gradient, we can equivalently write

$$\check{V}_c = [U^\top | 0] (F(\tilde{W}_r) - F_c(\tilde{W}_r)) \begin{bmatrix} V \\ 0 \end{bmatrix} = \tilde{U}^\top \begin{bmatrix} I & 0 \\ 0 & 0 \end{bmatrix} (F_c(\tilde{W}_r) - F(\tilde{W}_r)) \begin{bmatrix} I & 0 \\ 0 & 0 \end{bmatrix} \tilde{V}.$$

652 Remember the simplified variance corrected local coefficient update, given by

$$\begin{aligned} \tilde{S}_c^{s+1} &= \tilde{S}_c^s + \lambda \tilde{U}^\top \left( F_c(\tilde{W}_{r,c}^s) + \begin{bmatrix} I & 0 \\ 0 & 0 \end{bmatrix} (F_c(\tilde{W}_r) - F(\tilde{W}_r)) \begin{bmatrix} I & 0 \\ 0 & 0 \end{bmatrix} \right) \tilde{V} \\ &= \tilde{S}_c^s + \lambda \tilde{U}^\top \left( F_c(\tilde{W}_{r,c}^s) \right) \tilde{V} + \check{V}_c. \end{aligned} \quad (40)$$

### 653 H.1 Global loss descent for FeDLRT with simplified variance correction

654 In the following we provide proof for a global loss descent for Algorithm 5, i.e. using the local  
 655 coefficient update with variance correction (40).

656 **Theorem 11.** (Restatement of Theorem 4) Under Assumption 1, if the local learning rate  $0 < \lambda \leq$   
 657  $\frac{1}{12Ls_*}$ , then Algorithm 5 leads to the global loss descent

$$\mathcal{L}(W_r^{t+1}) - \mathcal{L}(W_r^t) \leq -s_* \lambda (1 - \delta^2 - 12s_* \lambda L + \delta^2 s_* \lambda) \left\| \tilde{U}^\top F(\tilde{W}_r) \tilde{V} \right\|^2 + L\vartheta, \quad (41)$$

658 with  $W_r^t = U^t S^t V^{t,\top}$  and  $W_r^{t+1} = U^{t+1} S^{t+1} V^{t+1,\top}$ .

659 *Proof.* We split the adjusted coefficient update in (40) into the non-augmented  $r \times r$  matrix  $S$  and  
 660 the tree off-diagonal blocks given by the augmentation  $\hat{S}$ :

$$\hat{S} = \tilde{S} - \begin{bmatrix} S & 0 \\ 0 & 0 \end{bmatrix}. \quad (42)$$

661 Analogously to the proof of Theorem 2, we consider

$$\begin{aligned} \mathcal{L}(\tilde{W}_r^*) - \mathcal{L}(\tilde{W}_r) &\leq \left\langle \tilde{W}_r^* - \tilde{W}_r, F(\tilde{W}_r) \right\rangle + \frac{L}{2} \left\| \tilde{W}_r^* - \tilde{W}_r \right\|^2 \\ &= \left\langle \tilde{U} \tilde{S}^* \tilde{V}^\top - \tilde{U} \tilde{S} \tilde{V}^\top, F(\tilde{W}_r) \right\rangle + \frac{L}{2} \left\| \tilde{U} \tilde{S}^* \tilde{V}^\top - \tilde{U} \tilde{S} \tilde{V}^\top \right\|^2 \\ &= \left\langle \tilde{S}^* - \tilde{S}, \tilde{U}^\top F(\tilde{W}_r) \tilde{V} \right\rangle + \frac{L}{2} \left\| \tilde{S}^* - \tilde{S} \right\|^2 \\ &= \left\langle \tilde{S}^* - \tilde{S}, -\nabla_{\tilde{S}} \mathcal{L}(\tilde{W}_r) \right\rangle + \frac{L}{2} \left\| \tilde{S}^* - \tilde{S} \right\|^2, \end{aligned}$$

662 where the transformation uses orthonormality of  $\tilde{U}$  and  $\tilde{V}$  and definition of the projected gradient.

663 We split the right hand side in terms corresponding to augmented terms  $\hat{S}$  and non-augmented terms  
 664  $S$  according to (42), i.e.,

$$\left\langle S^* - S, -\nabla_S \mathcal{L}(\tilde{W}_r) \right\rangle + \frac{L}{2} \|S^* - S\|^2, \quad (43)$$

665 which is treated exactly as in the proof of Theorem 2, and the augmented terms

$$\left\langle \hat{S}^* - \hat{S}, -\nabla_{\hat{S}} \mathcal{L}(\tilde{W}_r) \right\rangle + \frac{L}{2} \left\| \hat{S}^* - \hat{S} \right\|^2. \quad (44)$$

666 First we bound the term (43). Remember that  $\widehat{S} = 0$  at the start of the local iterations due to  
 667 orthonormality of  $\widetilde{U}, \widetilde{V}$ . The coefficient update (40) for  $S$  reads

$$S_c^{s+1} = S_c^s + \lambda U^\top \left( F_c(\widetilde{W}_{r,c}^s) - F_c(\widetilde{W}_r) + F(\widetilde{W}_r) \right) V. \quad (45)$$

668 Then we can readily apply Theorem 2 to obtain the bound

$$\left\langle S^* - S, -\nabla_S \mathcal{L}(\widetilde{W}_r) \right\rangle + \frac{L}{2} \|S^* - S\|^2 \leq -(s_* \lambda)(1 - 12(s_* \lambda)L) \left\| U^\top F(\widetilde{W}_r) V \right\|^2. \quad (46)$$

669 Next, we bound (44), starting with the first term:

$$\begin{aligned} \left\langle \widehat{S}^* - \widehat{S}, -\nabla_{\widehat{S}} \mathcal{L}(\widetilde{W}_r) \right\rangle &\stackrel{(I)}{=} \left\langle \widehat{S}^* - 0, -\nabla_{\widehat{S}} \mathcal{L}(\widetilde{W}_r) \right\rangle \\ &= \left\langle -\frac{\lambda}{C} \sum_{c=1}^C \sum_{s=0}^{s_*-1} \nabla_{\widehat{S}} \mathcal{L}_c(\widetilde{W}_{r,c}^s), -\nabla_{\widehat{S}} \mathcal{L}(\widetilde{W}_r) \right\rangle \\ &= \frac{\lambda}{C} \sum_{c=1}^C \sum_{s=0}^{s_*-1} \left\langle \nabla_{\widehat{S}} \mathcal{L}_c(\widetilde{W}_{r,c}^s), \nabla_{\widehat{S}} \mathcal{L}(\widetilde{W}_r) \right\rangle \\ &\leq \frac{\lambda}{C} \sum_{c=1}^C \sum_{s=0}^{s_*-1} \left\| \nabla_{\widehat{S}} \mathcal{L}_c(\widetilde{W}_{r,c}^s) \right\| \left\| \nabla_{\widehat{S}} \mathcal{L}(\widetilde{W}_r) \right\| \\ &\stackrel{(II)}{\leq} \frac{\lambda}{C} \sum_{c=1}^C \sum_{s=0}^{s_*-1} \delta^2 \left\| \nabla_{\widehat{S}} \mathcal{L}(\widetilde{W}_r) \right\| \left\| \nabla_{\widehat{S}} \mathcal{L}(\widetilde{W}_r) \right\| \\ &= \delta^2 s_* \lambda \left\| \nabla_{\widehat{S}} \mathcal{L}(\widetilde{W}_r) \right\|^2 = \delta^2 s_* \lambda \left\| \widetilde{U}^\top F(\widetilde{W}_r) \widetilde{V} \right\|^2, \end{aligned}$$

670 where we use  $\widehat{S} = 0$  in (I), and Assumption 1 in (II). Next, we bound the second term

$$\begin{aligned} \frac{L}{2} \left\| \widehat{S}^* - \widehat{S} \right\|^2 &= \frac{L}{2} \left\| -\frac{\lambda}{C} \sum_{c=1}^C \sum_{s=0}^{s_*-1} \nabla_{\widehat{S}} \mathcal{L}(\widetilde{W}_{r,c}^s) \right\|^2 \\ &\stackrel{(I)}{\leq} \frac{L}{2} \lambda^2 \frac{1}{C} \sum_{c=1}^C \left\| \sum_{s=0}^{s_*-1} \nabla_{\widehat{S}} \mathcal{L}(\widetilde{W}_{r,c}^s) \right\|^2 \\ &\stackrel{(I)}{\leq} \frac{L}{2} s_* \lambda^2 \frac{1}{C} \sum_{c=1}^C \sum_{s=0}^{s_*-1} \left\| \nabla_{\widehat{S}} \mathcal{L}(\widetilde{W}_{r,c}^s) \right\|^2 \\ &\leq s_* \frac{L}{2} \delta^2 \lambda^2 \frac{1}{C} \sum_{c=1}^C \sum_{s=0}^{s_*-1} \left\| \nabla_{\widehat{S}} \mathcal{L}(\widetilde{W}_r) \right\|^2 \\ &\leq \frac{L}{2} \delta^2 (s_* \lambda)^2 \left\| \nabla_{\widehat{S}} \mathcal{L}(\widetilde{W}_r) \right\|^2 = \frac{L}{2} \delta^2 (s_* \lambda)^2 \left\| \widetilde{U}^\top F(\widetilde{W}_r) \widetilde{V} \right\|^2, \end{aligned}$$

671 where we used Jensen's inequality in (I) again Assumption 1. We combine the bound on the  
 672 non-augmented terms (46) and the two bounds above for the augmented terms to

$$\begin{aligned} \mathcal{L}(\widetilde{W}_r^*) - \mathcal{L}(\widetilde{W}_r) &\leq \left\langle \widetilde{W}_r^* - \widetilde{W}_r, F(\widetilde{W}_r) \right\rangle + \frac{L}{2} \left\| \widetilde{W}_r^* - \widetilde{W}_r \right\|^2 \\ &\leq -(s_* \lambda)(1 - 12(s_* \lambda)L) \left\| U^\top F(\widetilde{W}_r) V \right\|^2 + \delta s_* \lambda \left\| \widetilde{U}^\top F(\widetilde{W}_r) \widetilde{V} \right\|^2 + \delta (s_* \lambda)^2 \left\| \widetilde{U}^\top F(\widetilde{W}_r) \widetilde{V} \right\|^2 \\ &\stackrel{(I)}{\leq} -(s_* \lambda)(1 - 12(s_* \lambda)L) \left\| \widetilde{U}^\top F(\widetilde{W}_r) \widetilde{V} \right\|^2 + \delta s_* \lambda \left\| \widetilde{U}^\top F(\widetilde{W}_r) \widetilde{V} \right\|^2 + \delta (s_* \lambda)^2 \left\| \widetilde{U}^\top F(\widetilde{W}_r) \widetilde{V} \right\|^2 \\ &= -(s_* \lambda)(1 - \delta^2 - 12(s_* \lambda)L + \delta^2 (s_* \lambda)) \left\| \widetilde{U}^\top F(\widetilde{W}_r) \widetilde{V} \right\|^2, \end{aligned}$$

673 where we use in (I)  $\left\|U^\top F(\widetilde{W}_r)V\right\| \leq \left\|\widetilde{U}^\top F(\widetilde{W}_r)\widetilde{V}\right\|$ . Using Equation (36), we can conclude the  
 674 proof:

$$\begin{aligned} & \mathcal{L}(U^{t+1}S^{t+1}V^{t+1,\top}) - \mathcal{L}(U^tS^tV^{t,\top}) \\ & \leq -(s_*\lambda)(1 - \delta^2 - 12(s_*\lambda)L + \delta^2(s_*\lambda)) \left\|\widetilde{U}^\top F(\widetilde{W}_r)\widetilde{V}\right\|^2 + L\vartheta. \end{aligned}$$

675

□

## 676 H.2 Global convergence of FeDLRT with simplified variance correction

677 **Corollary 2.** (Restatement of Corollary 1) Under Assumption 1, Algorithm 5 guarantees for the  
 678 learning rate  $\lambda \leq \frac{1}{s_*(12L+\delta^2)}$

$$\min_{t=1,\dots,T} \left\|\nabla_{\widetilde{S}}\mathcal{L}(W_r^t)\right\|^2 \leq \frac{96L}{T} (\mathcal{L}(W_r^1) - \mathcal{L}(W_r^{T+1})) + 96L^2\vartheta, \quad (47)$$

679 with  $W_r^t = U^tS^tV^{t,\top}$ ,  $W_r^1 = U^1S^1V^{1,\top}$ . and  $W_r^{T+1} = U^{T+1}S^{T+1}V^{T+1,\top}$ .

680 *Proof.* Consider Theorem 4,

$$\mathcal{L}(W_r^{t+1}) - \mathcal{L}(W_r^t) \leq -(s_*\lambda)(1 - \delta^2 - 12(s_*\lambda)L + \delta^2(s_*\lambda)) \left\|\widetilde{U}^\top F(\widetilde{W}_r)\widetilde{V}\right\|^2 + L\vartheta$$

681 and assume that  $\lambda s_* = \frac{1}{(12L+\delta^2)}$ , i.e.  $\lambda = \frac{1}{s_*(12L+\delta^2)} \leq \frac{1}{Ls_*}$ , which obeys the learning rate  
 682 requirement of Theorem 2. Plugging this learning rate into (38) gives

$$\left\|\nabla_{\widetilde{S}}\mathcal{L}(W_r^t)\right\|^2 \leq 96L (\mathcal{L}(W_r^t) - \mathcal{L}(W_r^{t+1}) + L\vartheta),$$

683 where we use  $(\frac{1}{4} - \delta^2) \leq \frac{1}{4}$  and  $\frac{1}{(12L+\delta^2)} \leq \frac{1}{12L}$ . Averaging from  $t = 1$  to  $t = T$  yields

$$\begin{aligned} \min_{t=1,\dots,T} \left\|\nabla_{\widetilde{S}}\mathcal{L}(W_r^t)\right\|^2 & \leq \frac{1}{T} \sum_{t=1}^T \left\|\widetilde{U}^\top F(\widetilde{W}_r)\widetilde{V}\right\|^2 \\ & \leq \frac{96L}{T} (\mathcal{L}(W_r^{t=1}) - \mathcal{L}(W_r^{t=T+1})) + 96L^2\vartheta, \end{aligned}$$

684 which concludes the proof. □

685 **I NeurIPS review**

686 **I.1 Paper Decision**

687 While the reviewers agreed that the work has interesting contributions and found merits in them, they  
688 raised several issues that are worth addressing in a careful and thorough manner. These include the  
689 exposition of the manuscript, scope of the numerical experiments and presentation of the numerical  
690 results, and possible extensions of the proposed approach to other settings. As the revision will be  
691 extensive and thus requires another round of review, and in view of the fact that NeurIPS can only  
692 accommodate one round of review, I regrettably have to reject the manuscript at this point.

693 **I.2 Official Comment by the Authors**

694 Wrap up statement of the discussion period

695 As the discussion period approaches its final day, we would like to thank the reviewers again for their  
696 feedback and hope we have clarified their remarks.

697 Overall, the reviewers pointed out that the submission proposes a sound algorithm solving the increas-  
698 ingly relevant and important problem of automatic compression for distributed and edge computing,  
699 while providing a robust theoretical foundation with global convergence proofs. The algorithm is  
700 evaluated on multiple datasets and network architectures (convolutional layers, transformers, and  
701 fully connected layers) and test problems. The reviewers found the paper to be well written.

702 We received valuable, constructive feedback by the reviewers and summarize the rebuttal and  
703 discussions in the following:

- 704 • We are happy to have resolved some unclear statements in the test-case descriptions, and  
705 contribution statements.
- 706 • Prompted by reviewer VUf8, we have extended the algorithm description for Tensor valued  
707 layers, prominently featured in convolutional neural networks, where we apply the proposed  
708 method to Tucker-Factorized tensors, and demonstrate their viability in the general answer  
709 PDF.
- 710 • We have provided more convergence plots to illustrate the effect of variance correction in  
711 neural network training, in addition to the plots for least squares regression in the paper.
- 712 • We have explained the mechanics of the basis augmentation, prompted by reviewer rsLZ,  
713 and how the basis augmentation, which does not induce any approximation errors, provides  
714 a key ingredient for the analysis. We further clarified that indeed, the coefficient update step  
715 resembles an optimization step on the low-rank manifold.
- 716 • We have clarified the consequences of the convergence guarantee on the distance of the  
717 trained solution to the stationary point in the fruitful discussion with reviewer hYpR.
- 718 • In the fruitful discussion with reviewer VUf8, we have clarified some limitations of the  
719 proposed method in context of heterogeneous data, and partial participation. In summary,  
720 we have seen that the proposed method works well in the setting
  - 721 – homogeneous data, deterministic gradient (Main paper, Section 4.1 Homogeneous test;  
722 supplemented by Figure 4)
  - 723 – homogeneous data, stochastic (mini-batch) gradient (Main paper, Section 4.2 and  
724 appendix; supplemented by Figure 5-8)
  - 725 – heterogeneous data, deterministic gradient (Main paper, Section 4.1 Heterogeneous  
726 test; supplement 1)
- 727 Preliminary tests during the discussion period showed that more research is required to  
728 provide good results for heterogeneous data, stochastic (mini-batch) gradient.
- 729 • Finally, we point out how the method can be extended to a partial participation scenario,  
730 where not all clients are active at the same time.

731 We hope that our answers have satisfied the reviewers, and we thank them again for their feedback.

732 Kind regards,

733 Authors

734 **I.3 Review 1 - hYpR**

735 Summary: This paper proposes a low-rank scheme to reduce communication and computation cost in  
736 FL, while also reducing client drift.

737 Soundness: 3: good Presentation: 3: good Contribution: 3: good Strengths: The paper is well-written  
738 and seems to be solving a relevant and important problem.

739 Weaknesses: see below

740 Questions:

741 Section 2: In Fig 1, how is the initial trajectory of FedAvg and FedLin identical till FedAvg settles?

742 Section 3: lines 126-7: just my curiosity, but why are SVD and QR decomposition not GPU friendly?  
743 In Fig 3 caption, the sentence about cost drop after is unclear. The description following Theorem 3  
744 connects the non-zero bias in (12-13) with Fig 1. However, Fig 1 shows distance to solution, while  
745 in theorem 2, 3, these are gradient norms. Can we really say anything much about the convergence  
746 based on bias in gradient norm, since in the worst case, we can be arbitrarily far from any stationary  
747 point?

748 Section 4: heterogeneous test case - why do all clients have access to all the training points? Shouldn't  
749 the data be distributed across clients as in the homogeneous case?

750 Limitations: n/a

751 Flag For Ethics Review: No ethics review needed.

752 Rating: 6: Weak Accept: Technically solid, moderate-to-high impact paper, with no major concerns  
753 with respect to evaluation, resources, reproducibility, ethical considerations.

754 Confidence: 3: You are fairly confident in your assessment. It is possible that you did not understand  
755 some parts of the submission or that you are unfamiliar with some pieces of related work. Math/other  
756 details were not carefully checked.

757 Code Of Conduct: Yes

758 **I.3.1 Rebuttal by authors**

759 Rebuttal: We thank the reviewer for their review. Each of the questions is addressed below.

760 Regarding the trajectories in Fig.1: We would like to clarify that, in the early stage, the trajectories  
761 of FedAvg and FedLin reported in Fig. 1 are very close but not identical. The similarity of these  
762 two trajectories is due to the fact that the variance correction term, which is the distinguishing factor  
763 between FedAvg and FedLin, see Eq. (4), being insignificant in the early stage of training for this  
764 problem. The variance correction term corrects the local gradient directions of the clients to prevent  
765 stalling of the convergence process. In this problem, the distances between model parameters at early  
766 stage and the the local optima are much longer than the distance between local (client) optima and  
767 the global optimum of the federated problem. In this case, the local gradients are good estimates of  
768 the global federated gradients and thus the variance correction effect is insignificant, which leads to  
769 similar behavior of FedAvg and FedLin. As the model parameters approach the local/global optimum,  
770 the local gradients are no longer good estimates of the global gradient due to data heterogeneity, and  
771 the variance correction term eventually results in better convergence behavior of FedLin. The same  
772 argument holds for the proposed low-rank methods with and without variance correction.

773 Regarding GPU friendliness: We consider SVD and QR not as GPU friendly as other parts of the  
774 proposed algorithm. They are less GPU friendly because the underlying SVD and QR algorithms are  
775 inherently sequential. For example, in a QR decomposition, the orthogonal space is build by rotating  
776 each column vector of a matrix onto the orthogonal complement of the subspace spanned by existing  
777 vectors. This sequential iterative procedure makes massively parallel implementation nontrivial, as  
778 opposed to, e.g., batchwise network evaluations.

779 Regarding the caption of Fig. 3: As for the sentence in question in the caption of Fig. 3, we meant to  
780 state that, when the rank is below 200, the communication, computation, and memory costs of the  
781 FeDLRT are lower than the costs of the full rank FedLin method. Thank you for pointing our the  
782 potential confusion. We will clarify this in a revised version.

783 Regarding the description following Theorem 3: Thank you for this remark. The result in Theorem  
784 3 describes convergence to a stationary point by providing upper bounds on the norm of the loss  
785 function gradient. For a general loss function, it is indeed possible that a point with small gradient  
786 magnitude can be far from the stationary points. However, if we assume that the loss function is  
787 locally strongly convex in a neighborhood of a stationary point, then the gradient magnitude can be  
788 used to bound the distance to this stationary point in the neighborhood. Please see, for example, Eq.  
789 (4.12) in Bottou, Léon, Frank E. Curtis, and Jorge Nocedal, "Optimization methods for large-scale  
790 machine learning." SIAM review 60, no. 2 (2018): 223-311, for the estimate and Appendix B therein  
791 for the proof.

792 Regarding the heterogeneous data test case: Thank you for pointing out the ambiguity in the problem  
793 description. In the heterogeneous linear regression test case, each client performs regression to a  
794 different target function. Therefore, even though they share the same 10,000 locations sampled on ,  
795 the local objective functions are defined with different target functions . We will clarify the problem  
796 configuration in a revised version.

### 797 **I.3.2 Comment by reviewer**

798 Thanks to the authors for their response. I maintain my score. All the best!

### 799 **I.4 Review 2 - VUf8**

800 Summary: The paper introduces FeDLRT, a federated algorithm to train and truncate low-rank  
801 weights automatically. The algorithm is based on a distributed version of the dynamic low-rank  
802 training. This requires multiple communication rounds (3 at worst) between the server and all clients,  
803 where first the  $U$  and  $V$  basis are augmented on the server after the basis gradients are sent from  
804 the clients and aggregated by the server. Then the clients learn the coefficients  $S$  and eventually  
805 correct the variance. The server then aggregates  $S$ , compress, and update the basis. The algorithm  
806 has theoretical guarantees of global convergence.

807 Soundness: 3: good Presentation: 2: fair Contribution: 3: good Strengths: The algorithm is sound  
808 and has theoretical guarantees of global convergence. Automatic compression is an increasingly  
809 important research topic, especially for edge and distributed training.

810 Weaknesses: I personally found the paper hard to follow and to distinguish between the actual  
811 contributions and what is instead based on the existing literature. The appendix is helpful, but  
812 I suggest the authors restructure section 3 and divide it into the background for dynamical low-  
813 rank training and their actual contributions. The algorithm seems a federated porting of the DLRT  
814 algorithm, which in order to have guarantees requires at least double communication per round to  
815 have shared augmented bases. Also, a clear contributions section would be helpful.

816 The way the CIFAR10 dataset has been split across clients is quite naive (only a few clients) and  
817 homogeneous - this is not a standard practice in FL where the algorithms are generally tested  
818 in heterogeneous non-iid settings, for instance splitting data among clients, based on a Dirichlet  
819 distribution (see for instance <https://arxiv.org/abs/1909.06335> and <https://arxiv.org/abs/2003.00295>).

820 The algorithm seems to be working only in full participation mode (at each round it needs to commu-  
821 nicate with all the clients), so it is mainly made for cross-silo settings with a few always available  
822 clients rather than cross-device. Indeed it requires 2 (or even 3 in the worst case) communication  
823 rounds (broadcast and aggregate operations)

824 Questions: Experiments on computer vision datasets: in the main paper the authors present exper-  
825 iments by training only the classifier using their proposed method. It is unclear if the method can  
826 be extended to all layers to train them and automatically compress them to their optimal rank. It is  
827 unclear if the method can be extended to convolutional layers.

828 It would be interesting to see plots of the loss and accuracy on the CIFAR dataset (with heterogeneity)  
829 to check the actual speed of convergence of the method against baselines. Something similar to  
830 Figure 4, but at least for the CIFAR10 dataset and against baselines such as FedAVG, FedLin, and  
831 potentially also something more recent to tackle heterogeneity. Indeed, while the method proposed  
832 has a variance reduction correction, apparently for mitigating client-drift, it is unclear if it can handle  
833 and mitigate the effect of heterogeneity.

834 Could the algorithm be extended to avoid communicating twice, hence to work in, cross-device,  
835 realistic, and partial participation settings?

836 Limitations: The authors should dedicate more space to the limitations of their approach as they are  
837 not clearly expressed and, while sound, the work seems not ready to be a practical algorithm yet.

838 Flag For Ethics Review: No ethics review needed.

839 Rating: 6: Weak Accept: Technically solid, moderate-to-high impact paper, with no major concerns  
840 with respect to evaluation, resources, reproducibility, ethical considerations.

841 Confidence: 3: You are fairly confident in your assessment. It is possible that you did not understand  
842 some parts of the submission or that you are unfamiliar with some pieces of related work. Math/other  
843 details were not carefully checked.

844 Code Of Conduct: Yes

#### 845 **I.4.1 Rebuttal by authors**

846 Rebuttal: We thank the reviewer for their review. To improve the presentation, we propose to  
847 restructure Section 2 and 3 by moving the description of the (non-federated) dynamical low-rank  
848 training from Section 3 to Section 2 as part of the background. The new Section 2 will include  
849 subsections on background for federated learning and variance correction, background for low-rank  
850 and dynamical low-rank training, as well as a standalone subsection on the contribution, which will  
851 be derived from the last paragraph of Section 2 in the current version. After this, the entire Section 3  
852 is dedicated to the proposed method and analysis.

853 We address each of the questions below.

854 Regarding compressing convolutions: We focused on the classifier since these layers are matrix-  
855 valued, and thus the proposed algorithm is directly applicable. We have extended the implementation  
856 of FeDLRT to train convolutional layers in a low-rank fashion as well. The results of FeDLRT applied  
857 to all layers (convolutions and classifiers) of VGG16 on CIFAR10 are reported in Fig. 2 in the general  
858 response PDF file. These results resemble the the ones in Fig. 7 with slightly different compression  
859 ratios, since more layers are now low rank. In the following paragraph, we give technical details in  
860 the extension of FeDLRT to compress convolutional layers.

861 To extend FeDLRT to convolutional layers, we follow the approach considered in, e.g.,  
862 (<https://arxiv.org/abs/2305.19059>) for (non-federated) DLRT, where a 2D convolution is interpreted  
863 as an order-4 tensor and factorized by using the Tucker decomposition. To this end, the Tucker bases  
864  $U_i \in \mathbb{R}^{n_i \times r_i}$  for  $i = 1, \dots, 4$ . replace the  $U$  and  $V$  bases in the matrix case, and the Tucker core  
865 tensor  $C \in \mathbb{R} \in \mathbb{R}^{r_1 \times \dots \times r_4}$  replaces the coefficient matrix  $S$ , to which the variance correction  
866 is applied. The analysis holds for the Tucker Tensor case, since Tucker Tensors have a manifold  
867 structure. In the proofs, we need to project onto all bases  $U_i$ . The compression step is performed  
868 with an truncated Tucker decomposition of the core tensor, instead of an SVD of the coefficient  
869 matrix. For intuition, one can also refer to the matrix case as the order-2 Tucker Tensor case. Remark  
870 that the bases are all updated simultaneously, thus the adaption to the tensor case does not require  
871 more communication rounds.

872 Regarding accuracy plots: We thank the reviewer for the constructive question. First, we provide  
873 in the general response PDF file, Fig. 1, a convergence plot for Resnet18 on CIFAR10 for the  
874 (homogeneous) test case reported in Fig. 5 of the main manuscript. One can see that the benefit of the  
875 variance correction term (in FeDLRT w/ var/cor and FedLin) mitigates the stalling of the convergence  
876 seen in the non-variance-corrected methods (FeDLRT w/o var/cor and FedAvg).

877 Regarding heterogeneous test cases: Prompted by this question, we conducted a preliminary study for  
878 a federated scenario with heterogeneous data on the client devices, drawn from a Dirichlet distribution.  
879 We found that the variance correction does not provide significant performance increase in scenarios  
880 with stochastic multi-batch gradient descent on clients and strong heterogeneity. Given the positive  
881 results for homogeneous data, we consider this challenge a relevant future research direction and  
882 will investigate the potential incorporation of more recent techniques into FeDLRT to tackle strongly  
883 heterogeneous data.

884 Regarding a modification to reduce the communication rounds: The FeDLRT algorithm and the  
885 convergence analysis require communication of the basis and optionally the variance correction



886 term prior to the client coefficient updates, therefore, both communication rounds are necessary.  
887 The variance corrected baseline, FedLin, considered in this work also requires two communication  
888 rounds. Moreover, we would like to emphasize that the total communication cost (including all  
889 communication rounds) per aggregation round of FeDLRT in practice is nearly an order of magnitude  
890 smaller than the full-rank baselines, e.g. FedLin or FedAvg, because FeDLRT only communicates  
891 part of the factors each round. See Fig. 3 of the manuscript and results in e.g. Fig. 5. We also remark  
892 that the variance correction benefits the convergence behavior, see, e.g. Fig. 1, and the two right  
893 panels of Fig. 4 in the main manuscript, as well as Fig. 1 in the general response PDF file. The  
894 superior convergence behavior implies that the proposed method reaches the target accuracy in fewer  
895 aggregation rounds, thus requiring fewer overall communication rounds.

896 A potential limitation of having two communication rounds, instead of one, is that latency differences  
897 of clients are more pronounced during hand-shakes. However, even for basic methods with single  
898 communication round, e.g. FedAvg, latency differences still pose a problem. To fully address this  
899 issue, one may need to extend the method non-trivially to accommodate asynchronous communication  
900 scenarios, which we find relevant as a future research direction.

901 On the other hand, allowing for partial participation is certainly possible in FeDLRT, as long as  
902 the active clients are consistent in all communication rounds within the same aggregation round.  
903 However, we have not been able to experiment in the partial participation configuration with many  
904 ( $> 100$ ) clients training relevant network architectures such as Resnet18 or VGG16, due to the  
905 constraint on the computation resources and the current implementation of FeDLRT. We agree that  
906 this is an important research direction and will make attempts to scale up the FeDLRT method.

#### 907 **I.4.2 Comment by reviewer**

908 Thank you for answering my questions and concerns.

909 Could you clarify (please be specific) how the algorithm could work in the partial participation case  
910 and if errors could arise (if the algorithm’s guarantees are broken), especially in the heterogeneous  
911 case?

#### 912 **I.4.3 Comment by authors**

913 We start our answer with a description of FeDLRT without variance correction for the partial  
914 participation case and then discuss a potential direction for extending FeDLRT to handle data  
915 heterogeneity in the partial participation scenario.

916 When variance correction is turned off, FeDLRT can be applied to partial active clients by considering  
917 only a (potentially random) subset  $C^t \subset 1, \dots, C$  of clients within a global aggregation round in  
918 Algorithm 1 in the original manuscript. Specifically, at aggregation round  $t$ , only the clients  $C^t$  are  
919 taken into account in the broadcasting in lines 2 and 6, client operations in lines 3, 7, 8, 15, and the  
920 aggregations in lines 4 and 16.

921 Due to the partial participation, the gradients computed in line 4 are no longer the global loss gradients  
922 with respect to  $U$  and  $V$ , respectively. However, this does not break the mechanism and analysis of  
923 FeDLRT since the augmented basis is not required to come from augmenting the global gradients.

924 The set of active clients can vary for different aggregation rounds, but, for FeDLRT,  $C^t$  needs to  
925 remain constant within an aggregation round. This restriction is consistent to the scenario considered  
926 in most existing work on federated learning with partial participation.

927 As for the performance, we expect that the FeDLRT w/o variance correction described above to  
928 perform similarly as FedAvg in terms of final accuracy, but at a much lower communication and  
929 memory cost due to the low rank technique.

930 Based on the preliminary results on heterogeneous data discussed in the rebuttal, we do not expect  
931 the current variance correction scheme to provide significant advantages in the partial participation  
932 case with heterogeneous data. A potential approach to address this issue is to incorporate in the  
933 FeDLRT algorithm an advanced variance correction scheme, such as the FedVARP scheme proposed  
934 in <https://openreview.net/forum?id=HIWLLdUocx5>, which is tailored to the partial participation  
935 case with heterogeneous data.

936 We are happy to provide more details if there are further questions.

937 **I.4.4 Comment by reviewer**

938 Thank you for your responses. In its current state, I still believe the paper is borderline, as it does not  
939 seem intended for general heterogeneous federated learning. That said, this could be a bias on my  
940 part, as this is one of my main areas of expertise. The rest of the paper and the authors' responses  
941 are convincing, but I believe the authors should incorporate their explanations and additional details  
942 about my concerns into the main paper. I have the impression that the algorithm is more suited  
943 for distributed learning, where heterogeneity and partial participation are less of an issue, though  
944 compression could still be beneficial due to communication constraints.

945 I am raising my score, and I hope the authors will consider my concerns and suggestions in the final  
946 version of the paper as well as in future work.

947 **I.5 Review 3 - rsLZ**

948 Summary: This paper introduces FeDLRT (Federated Dynamical Low-Rank Training), an innovative  
949 method designed to enhance federated learning by incorporating a low-rank client optimization  
950 step and an optional variance correction mechanism. FeDLRT builds upon the dynamic low-rank  
951 approximation (DLRA) method, extending it to neural network training in a federated learning  
952 context. The key contributions of this work include the development of a basis update and Galerkin  
953 (BUG) splitting scheme that allows for the efficient and dynamic adjustment of the rank, ensuring  
954 client-wide manifold consistency, and minimizing communication costs.

955 Soundness: 2: fair Presentation: 3: good Contribution: 3: good Strengths: The paper presents a  
956 robust theoretical framework by building upon the dynamic low-rank approximation (DLRA) method  
957 and extending it to the federated learning context.

958 The dynamic adjustment of the rank through the BUG splitting scheme is an innovation. This approach  
959 not only ensures client-wide manifold consistency but also enables efficient basis augmentation and  
960 coefficient updates, leading to better utilization of communication resources. The optional variance  
961 correction mechanism adds another layer of robustness, addressing potential discrepancies in local  
962 updates and ensuring convergence.

963 The extensive evaluation on real datasets demonstrates the effectiveness of the proposed approach for  
964 federated dynamical low-rank training.

965 Weaknesses: The method requires two communication rounds—one for aggregating global basis  
966 gradients and another for locally updated coefficients, which might still be considered high in some  
967 federated learning scenarios. When the number of clients is large, each gradient and basis update can  
968 result in additional communication overhead. Can everything be done in one communication round?

969 The experiments are limited to ResNet18 on CIFAR-10. This method involves gradient calculation  
970 and local optimization on the client side, as well as incremental basis update and QR decomposition  
971 on the server side. Whether the model is valid when applied to higher-dimensional data or larger  
972 models such as RoBERTa or LLaMA, and large datasets like SST-2.

973 When updating the basis  $U$  and  $V$ , the effect of the upper triangular matrix  $R$  is ignored in the new  
974 incremental basis obtained by using QR decomposition. Will this affect the performance of the  
975 model? What is the error range caused by updating the coefficient matrix  $S$  with the new incremental  
976 basis?

977 When updating the incremental coefficient matrix  $S$  in this paper, using an update method similar to  
978 SGD will lead to the original parameter not being on the manifold after updating.

979 It is better to conduct the experiments with baselines. Otherwise it is difficult to justify the effective-  
980 ness of the proposed method.

981 Questions: Weaknesses

982 Limitations: No.

983 Flag For Ethics Review: No ethics review needed.

984 Rating: 3: Reject: For instance, a paper with technical flaws, weak evaluation, inadequate repro-  
985 ducibility and/or incompletely addressed ethical considerations.

986 Confidence: 4: You are confident in your assessment, but not absolutely certain. It is unlikely, but not  
987 impossible, that you did not understand some parts of the submission or that you are unfamiliar with  
988 some pieces of related work.

989 Code Of Conduct: Yes

### 990 I.5.1 Rebuttal by Authors

991 We thank the reviewer for this review. Please find the answer to the questions below.

992 1. Regarding communication cost: The method requires two communication rounds, since the  
993 basis update and variance correction term need to be available to each (active) client before  
994 the client update starts. This being said, the proposed method communicates only parts of  
995 the weight matrix factors during each communication round, i.e. its total communication  
996 cost is significantly reduced compared to baseline methods, such as FedAvg and FedLin.  
997 Further, FedLin, the baseline with variance correction, also requires two communication  
998 rounds. We remark that the variance correction also benefits the convergence behavior, see,  
999 e.g. Fig. 1, and the two right panels of Fig. 4 in the main manuscript, as well as Fig. 1 in the  
1000 general response PDF file. The superior convergence behavior implies that the proposed  
1001 method reaches the target accuracy in fewer aggregation rounds, thus requiring fewer overall  
1002 communication rounds. In conclusion, we argue that the total number of communicated  
1003 floating point numbers is significantly reduced in FedDRLT, compared to the mentioned  
1004 baselines.

1005 2. Regarding experiments: In addition to ResNet18 on CIFAR10, we provide numerical results  
1006 for two convex test problems in the main manuscript, as well as AlexNet on CIFAR10,  
1007 VGG16 on CIFAR10, and a Vision Transformer on CIFAR100 in the appendix, thus dis-  
1008 cussing performance on convolutional networks and transformers, two of the most widely  
1009 used network architectures.

1010 The QR decomposition required in the basis augmentation step acts on a tall, but skinny  
1011  $n \times 2r$  matrix, thus requiring  $(2r)^2$  computational cost (typically  $n \gg r$ ), which is still smaller  
1012 than the  $n^2$  cost of multiplying a full-rank weight matrix with an input vector required in  
1013 the full-rank baseline methods, such as FedLin and FedAvg. Further, the QR decomposition  
1014 is performed once per aggregation round on the server, which has typically more compute  
1015 resources than the clients. We stress that the method aims to minimize total communication  
1016 and client compute costs, combined with preferred convergence behavior. Considering  
1017 Table 1, we remark that FedDRLT is (to the best of our knowledge) the only one with linear  
1018 dependence of the client compute cost on the layer dimensions.

1019 3. Regarding the basis update: The basis update extends the old basis by the span of the gradient  
1020 dynamics. Thus, the spans of the augmented bases obtained in the basis augmentation  
1021 step also contain the spans of original bases. Consequently, **no error is introduced by**  
1022 **augmenting the basis, and the training loss does not increase.** Intuitively the basis  
1023 augmentation can be seen as a conservative extension of the search space of the neural  
1024 network training: we allow to search for new coefficients in a manifold of twice the rank.

1025 In further detail, we refer to line 5 of Algorithm 1 (using Eq. (6)), where the basis update of  
1026  $U$  and  $V$  is performed. Due to the QR decomposition, we have  $\text{span}(\tilde{U}) = \text{span}([U^t, G_U])$ .  
1027 **The  $R$  matrix is not relevant to the construction of the new basis and thus can be**  
1028 **discarded in the algorithm.** However, since  $U^t$  is already orthonormal by construction,  
1029 we further have  $\tilde{U} = [U^t, \bar{U}]$  with  $U^t \perp \bar{U}$ , which implies that the upper half of  $R$  is a  
1030 unit matrix. This is indeed important since it yields the explicit expression of  $\tilde{S}$  in Lemma  
1031 1. As a consequence, the augmented low-rank representation  $\tilde{U}\tilde{S}\tilde{V}^\top$  is consistent with  
1032 the non-augmented representation  $USV^\top$ , i.e.  $\|\tilde{U}\tilde{S}\tilde{V}^\top - USV^\top\|_F = 0$ , which is a  
1033 requirement in the proof of Theorems 2 and 4.

1034 4. Regarding coefficient updates: The method is carefully constructed so that the coef-  
1035 ficient matrix update is an update within the manifold of rank  $2r$  matrices, because the  
1036 bases  $\tilde{U}$  and  $\tilde{V}$  remain constant in the client update steps. This not only implies that  
1037 the updates stay on the manifold, but that the proposed method is robust with respect to  
1038 the curvature of the low-rank manifold. We refer to Appendix D and specifically The-  
1039 orem 5 in the manuscript for a technical discussion of the robust optimization method

1040 that forms the foundation of this federated scheme. For further reading on why the  
1041 BUG scheme is a robust optimization method on manifolds, we would like to refer to  
1042 [<https://arxiv.org/pdf/2205.13571>, Section 4]. For a well-written geometric interpretation of  
1043 the method, we refer to (<https://arxiv.org/abs/1705.08521>).

1044 5. Baselines of experiments: We compare FeDLRT to the full-rank baselines, FedAvg and  
1045 FedLin, in all numerical experiments. We show that across all test cases, the FeDLRT  
1046 method confidently mirrors the convergence behavior of its full-rank counterpart, just as  
1047 estimated in Theorem 5. Meanwhile, FeDLRT dynamically compresses the model to reduce  
1048 communication bandwidth and the computational cost.