
The Interpolated MVU Mechanism For Communication-efficient Private Federated Learning

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 We consider private federated learning (FL), where a server aggregates differentially
2 private gradient updates from a large number of clients in order to train a machine
3 learning model. The main challenge is balancing privacy with both classification
4 accuracy of the learned model as well as the amount of communication between
5 the clients and server. In this work, we build on a recently proposed method for
6 communication-efficient private FL—the MVU mechanism—by introducing a new
7 interpolation mechanism that can accommodate a more efficient privacy analysis.
8 The result is the new Interpolated MVU mechanism that provides SOTA results on
9 communication-efficient private FL on a variety of datasets.

10 1 Introduction

11 Machine-learned models leak information about their training data [26]. Private training methods
12 have been developed to train models that provide rigorous guarantees quantifying the amount of
13 information leaked [1, 8, 25]. *Federated learning* (FL) builds on private training to collaboratively
14 train a model among many devices while keeping the data at each device private [20]. To accomplish
15 this, (cross-device) FL requires that devices communicate updates to a server coordinating the training.
16 These updates can be privatized using a differentially private mechanism such as DP-SGD [1] by
17 injecting a controlled amount of noise into the gradient, or update direction, at each step.

18 To reduce communication overhead in FL, it is also of interest to compress updates before they are
19 transmitted to the server, and lossy compression can also be seen as a way of injecting noise into
20 updates. Most previous work has addressed the challenges of privacy and compression separately, first
21 applying a DP mechanism to privatize the response, and then compressing before transmitting [2, 12].

22 Recent work [7] introduces the *minimum-variance unbiased* (MVU) mechanism for jointly com-
23 pressing and ensuring privacy, and experimentally demonstrates that this can lead to better utility-
24 compression trade-offs than other methods which first privatize and then compress. The core of
25 MVU consists of a private mechanism that works for a finite number of scalar inputs. If the input is a
26 bounded continuous scalar, then the solution is to dither to this finite set before applying the core
27 mechanism, and this is further extended to vectors by privacy composition over all coordinates via
28 Rényi DP [21]. Empirically, the MVU mechanism achieves state-of-the-art performance in the local
29 DP setting for both distributed mean estimation and federated learning [7]. However, the analysis
30 in [7] does not benefit from randomization introduced by dithering, and furthermore the extension to
31 vectors leads to suboptimal privacy composition for the L_2 geometry, which is often of interest (*e.g.*,
32 working with L_2 -bounded update vectors such as in DP-SGD).

33 **Contributions.** Building on a simplified version of the MVU mechanism with only a single scalar
34 input, we propose the *interpolated MVU* (I-MVU) mechanism—a more natural interpolation mecha-
35 nism to extend MVU to continuous inputs. By its discrete nature, the MVU mechanism can be viewed

36 as sampling from a particular categorical distribution, and hence can be expressed in exponential
 37 family form. The proposed I-MVU mechanism handles continuous inputs by *interpolating* the natural
 38 exponential family parameters, rather than directly interpolating the probabilities as in dithering. We
 39 introduce a new analysis technique and, by further exploiting special properties of the exponential
 40 family, obtain a tight privacy analysis for the vector extension under L_2 geometry. Experimentally,
 41 we find that under both client-level and user-level DP settings, the I-MVU mechanism provides better
 42 privacy-utility trade-off than SignSGD [17] and MVU [7] at an extremely low communication budget
 43 of *one bit per gradient dimension*. Moreover, I-MVU achieves close to the same performance as the
 44 standard non-compressed Gaussian mechanism [1] for similar levels of (ϵ, δ) -DP.

45 2 Background and Related Work

46 **Differential privacy.** The framework of differential privacy [10] allows rigorous reasoning of
 47 privacy leakage through a mechanism \mathcal{M} applied to a dataset \mathcal{D} . We say that \mathcal{M} is (ϵ, δ) -differentially
 48 private, denoted (ϵ, δ) -DP, if for any \mathcal{D} , any $\mathbf{x} \in \mathcal{D}$ and any output set O , we have:

$$e^{-\epsilon} \mathbb{P}(\mathcal{M}(\mathcal{D} \setminus \mathbf{x}) \in O) - \delta \leq \mathbb{P}(\mathcal{M}(\mathcal{D}) \in O) \leq e^{\epsilon} \mathbb{P}(\mathcal{M}(\mathcal{D} \setminus \mathbf{x}) \in O) + \delta.$$

49 More generally, the framework of DP seeks to bound the difference in distribution between $\mathcal{M}(\mathcal{D})$
 50 and $\mathcal{M}(\mathcal{D} \setminus \mathbf{x})$ so that a single record \mathbf{x} will not affect the output of the mechanism \mathcal{M} significantly.

51 A useful variant of DP is Rényi differential privacy (RDP) [21], which instead bounds the Rényi
 52 divergence [23] between $\mathcal{M}(\mathcal{D})$ and $\mathcal{M}(\mathcal{D} \setminus \mathbf{x})$ by some ϵ . Formally, we say that \mathcal{M} is (α, ϵ) -RDP if

$$D_{\alpha}(\mathcal{M}(\mathcal{D}) \parallel \mathcal{M}(\mathcal{D} \setminus \mathbf{x})) \leq \epsilon \quad \text{and} \quad D_{\alpha}(\mathcal{M}(\mathcal{D} \setminus \mathbf{x}) \parallel \mathcal{M}(\mathcal{D})) \leq \epsilon,$$

53 where D_{α} denotes the order- α Rényi divergence [21]. Importantly, Rényi DP supports composition
 54 of mechanisms in a simple manner: If $\mathcal{M}_1, \dots, \mathcal{M}_T$ are mechanisms with \mathcal{M}_t being (α, ϵ_t) -RDP
 55 for $t = 1, \dots, T$, then the composition of the T mechanisms is $(\alpha, \sum_{t=1}^T \epsilon_t)$ -RDP. Another useful
 56 property of RDP is its conversion to (ϵ, δ) -DP [3]: If \mathcal{M} is $(\alpha, \epsilon_{\alpha})$ -RDP for $\alpha > 1$ then it is also
 57 (ϵ, δ) -DP for any $0 < \delta < 1$ with

$$\epsilon = \epsilon_{\alpha} + \log \left(\frac{\alpha - 1}{\alpha} \right) - \frac{\log \delta + \log \alpha}{\alpha - 1}. \quad (1)$$

58 **Federated learning with differential privacy.** Federated learning (FL) [18, 20] allows distributed
 59 training of ML models across multiple clients without centralized data storage. A server coordinates
 60 training by acquiring model updates from clients, aggregating them, and then transmitting an updated
 61 model back to the clients, with the process repeating until convergence. One promise of FL is data
 62 privacy since the updates are computed locally on each client using their own data, and hence no
 63 client data is ever explicitly transmitted to the server (or anyone else) throughout the training process.
 64 In spite of this, a recent line of work showed that despite the clients never explicitly sharing their
 65 data, it is possible to reconstruct training samples from the model updates in a process called *gradient*
 66 *inversion* [11, 31, 32]. This vulnerability remains even if a large number of clients participate in a
 67 round using secure aggregation [11, 16, 30].

68 Differential privacy is a principled method to ensure data privacy in FL as it provides provable
 69 guarantees against data reconstruction from the output of a private mechanism [4, 13, 27]. To apply
 70 DP to FL training, given a client update \mathbf{x} , the client instead sends $\mathcal{M}(\mathbf{x})$ to the server. For a
 71 given round, the client's privacy leakage can be computed in terms of *local DP* if the privatized
 72 update $\mathcal{M}(\mathbf{x})$ is revealed to the server, or *global DP* if secure aggregation is applied to aggregate the
 73 privatized updates before revealing it to the server. The total privacy leakage throughout training can
 74 then be computed via RDP composition and conversion to (ϵ, δ) -DP via Equation 1.

75 **Communication-efficient private mechanisms.** Since model updates in FL are high-dimensional
 76 vectors of size equal to the number of model parameters, it is also important in practice to compress
 77 these updates for communication efficiency. This requirement combined with privacy has led to a
 78 series of prior work that designed communication-efficient private mechanisms with application to
 79 FL [2, 6, 7, 9, 12, 17, 24, 29]. However, compressing the model update often leads to higher variance
 80 and/or biasedness [7, 9], and as a result the model's performance is subpar compared to ones trained
 81 using non-compressed DP mechanisms such as the Gaussian mechanism [7, 17]. In this work, we

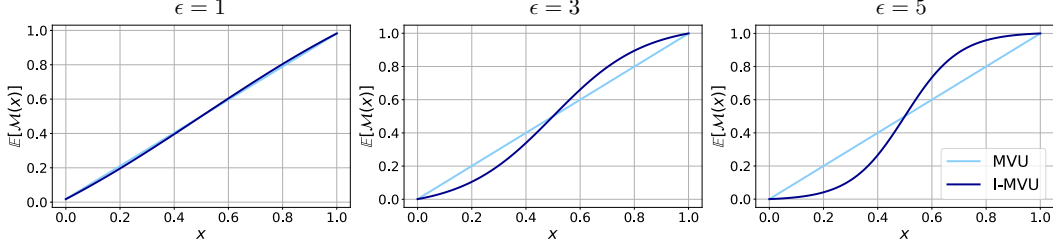


Figure 1: Plot showing the expected value of $\mathcal{M}(x)$ for different values of $x \in [0, 1]$ for the MVU and I-MVU mechanism. While MVU is unbiased in the entire interval $[0, 1]$, I-MVU incurs some bias for $x \neq 0, 1$, especially at higher values of the DP parameter ϵ .

82 drastically reduce this performance gap and show that replacing the Gaussian mechanism with the
 83 proposed interpolated MVU mechanism leads to the same test performance at equal privacy cost
 84 when using *one-bit output per coordinate*.

85 3 Interpolated MVU Mechanism

86 We introduce the *interpolated MVU mechanism*—a communication-efficient differentially-private
 87 mechanism with application to private FL training. We also present a novel privacy analysis technique
 88 for privatizing vectors with L_2 geometry, leading to a drastic improvement in the privacy-utility
 89 trade-off over the previously proposed MVU mechanism [7]. To this end, we begin by defining the
 90 problem of private-compression and recalling the MVU mechanism.

91 **Problem description.** Consider the private-compression problem of transmitting a vector $\mathbf{x} \in \mathbb{R}^d$
 92 with bounded L_2 -norm privately using at most bd bits, with b small enough so that the entire vector
 93 \mathbf{x} can be transmitted efficiently. One can reduce this problem to a scalar one by considering how
 94 to privately compress $x \in [0, 1]$ using at most b bits, and then scaling the vector \mathbf{x} appropriately to
 95 $[0, 1]^d$ and applying the scalar mechanism coordinate-wise.

96 **The minimum variance unbiased (MVU) mechanism [7]** solves the private-compression problem
 97 by first discretizing the interval $[0, 1]$ into B_{in} points $\mathcal{X} = \{x_1 = 0, x_2, \dots, x_{B_{\text{in}}} = 1\}$ with
 98 $x_i := (i - 1)/(B_{\text{in}} - 1)$. If $x = x_i$, the mechanism samples $j \sim \text{Categorical}(\mathbf{p}_i)$ using a probability
 99 vector $\mathbf{p}_i \in \Delta^{B_{\text{out}}-1}$ and outputs $\mathcal{M}(x) = a_j \in \mathbb{R}$ where $\{a_1, \dots, a_{B_{\text{out}}}\}$ is a pre-determined output
 100 alphabet. The probability vectors $\mathbf{p}_1, \dots, \mathbf{p}_{B_{\text{in}}}$ and output alphabet $\{a_1, \dots, a_{B_{\text{out}}}\}$ are designed so
 101 that the mechanism satisfies the following three properties:

- 102 1. ϵ -Differential Privacy: $e^{-\epsilon} \mathbf{p}_{i',j} \leq \mathbf{p}_{i,j} \leq e^{\epsilon} \mathbf{p}_{i',j}$ for all $i \neq i'$ and all j .
- 103 2. Unbiasedness: $\sum_{j=1}^{B_{\text{out}}} a_j \mathbf{p}_{i,j} = x_i$ for all i .
- 104 3. Minimum variance: $\sum_{i=1}^{B_{\text{in}}} \text{Var}(\mathcal{M}(x_i))$ is minimal among all mechanisms satisfying 1 and 2.

105 The MVU mechanism can then be applied to all $x \in [0, 1]$ by *randomly dithering* x to the nearest
 106 x_i and x_{i+1} such that the dithering is unbiased in expectation. One can also view this dithering
 107 procedure as linearly interpolating between \mathbf{p}_i and \mathbf{p}_{i+1} . It is straightforward to generalize the
 108 mechanism to any bounded x by scaling it to $[0, 1]$ and then applying the MVU mechanism.

109 For a d -dimensional vector \mathbf{x} , the MVU mechanism can be applied independently to each coordinate
 110 and the privacy cost is $d\epsilon$ by composition if $\mathbf{x} \in [0, 1]^d$ (or in general, if $\|\mathbf{x}\|_{\infty}$ is bounded). However,
 111 the privacy analysis becomes much more complicated for L_2 -norm bounded vectors—as is often
 112 the case for DP-SGD training [1]. We address this problem by expressing the MVU mechanism in
 113 exponential family form and interpolating in the natural parameter space, allowing us to use special
 114 properties of exponential family distributions to derive tight privacy analysis for the L_2 geometry.

115 **Interpolated MVU mechanism.** As mentioned above, by dithering an input $x \in [x_i, x_{i+1}]$ to x_i or
 116 x_{i+1} , for general inputs $x \notin \mathcal{X}$, the MVU can be seen as linearly interpolating between the probability
 117 vectors \mathbf{p}_i and \mathbf{p}_{i+1} . Here we improve upon MVU by introducing a better form of interpolation. The

118 pmf for the categorical distribution with natural parameter $\boldsymbol{\eta}$ can be written as:

$$\mathbb{P}(j|\boldsymbol{\eta}) = \exp(\mathbf{e}_j^\top \boldsymbol{\eta} - A(\boldsymbol{\eta})), \quad A(\boldsymbol{\eta}) = \log \left(\sum_j \exp(\boldsymbol{\eta}_j) \right) \quad (2)$$

119 where \mathbf{e}_j is the j -th standard basis vector. Note that if $\mathbf{p} \in \Delta^{B_{\text{out}}-1}$ then its natural parameter is
 120 $\boldsymbol{\eta} = \log \mathbf{p}$. To define the *interpolated MVU* (I-MVU) mechanism, let $\mathbf{p}_1, \mathbf{p}_2 \in \Delta^{B_{\text{out}}-1}$ be sampling
 121 probability vectors obtained from the MVU mechanism with $B_{\text{in}} = 2$ and let $\boldsymbol{\eta}_i = \log \mathbf{p}_i$ for $i = 1, 2$
 122 be the natural parameters. Given $x \in [0, 1]$, the I-MVU mechanism samples $j \sim \mathbb{P}(\cdot|\boldsymbol{\eta}(x))$ according
 123 to Equation 2 and outputs a_j from the MVU output alphabet, where

$$\boldsymbol{\eta}(x) = (1-x)\boldsymbol{\eta}_1 + x\boldsymbol{\eta}_2. \quad (3)$$

124 In other words, instead of linearly interpolating between \mathbf{p}_1 and \mathbf{p}_2 to construct the sampling
 125 probability vector for x , we interpolate in the natural parameter space of the categorical distribution.
 126 Doing so incurs some bias¹ when $x \notin \{0, 1\}$; see Figure 1 for a plot illustrating this phenomenon.
 127 Nevertheless, this bias is small in comparison to the noise induced by differential privacy, and we
 128 show empirically that it does not affect the performance of the I-MVU mechanism for FL training.

129 **Input scaling.** One way to extend I-MVU to arbitrary bounded ranges is to first scale the input
 130 to $[0, 1]$ and then apply the mechanism as usual. However, note that the interpolation scheme in
 131 Equation 3 is in fact well-defined for any $x \in \mathbb{R}$, and hence the scaled input does not need to be strictly
 132 in the range $[0, 1]$. We leverage this property by introducing a scaling factor β : For $u \in [-C, C]$, the
 133 β -scaled I-MVU mechanism is defined as

$$\mathcal{M}_\beta(u) = \mathcal{M} \left(\frac{1}{2} + \frac{\beta u}{2C} \right),$$

134 where \mathcal{M} is the plain I-MVU mechanism. Note that this scaling effectively ensures that the input
 135 to \mathcal{M} is in the range $[(1-\beta)/2, (1+\beta)/2]$, with $\beta = 1$ corresponding to scaling the input to $[0, 1]$.
 136 For vectors \mathbf{u} with $\|\mathbf{u}\|_2 \leq C$, the β -scaled input $\mathbf{x} = \frac{1}{2} + \frac{\beta \mathbf{u}}{2C}$ satisfies $\|\mathbf{x}\|_2 \leq \beta/2$

137 One advantage for using β -scaling is that if the distribution of u is highly concentrated near zero, then
 138 scaling with $\beta > 1$ ensures that the input to \mathcal{M} is more spread out in the range $[0, 1]$. This ensures
 139 that the input distribution more closely reflects the minimum variance requirement (property 3) for
 140 the MVU mechanism. For L_2 -norm bounded vectors \mathbf{u} this is especially true, where the distribution
 141 of coordinates of \mathbf{u} is likely concentrated near zero. Consequently, for compressing client updates
 142 with bounded L_2 -norm, β -scaling with a large β is essential for achieving good performance.

143 3.1 Privacy Analysis

144 We analyze privacy leakage of the I-MVU mechanism for L_2 -norm bounded vectors in terms of
 145 Rényi DP [21]. Our strategy is to first analyze the scalar mechanism and express its Rényi divergence
 146 for two differing inputs x_1 and x_2 as a function of $(x_2 - x_1)^2$ (Lemma 1). Then, by independently
 147 applying the mechanism across coordinates, we can sum the Rényi divergence across coordinates and
 148 upper bound the total RDP ϵ as a function of $\|\mathbf{x}_2 - \mathbf{x}_1\|_2^2$ (Theorem 1). Our analysis depends crucially
 149 on a measure of information known as *Fisher information*, which we define below for completeness.

150 **Definition 1.** Let f be the density function of a distribution parameterized by $x \in \mathbb{R}$. The Fisher
 151 information of x contained in a sample $Z \sim f(\cdot|x)$ is:

$$\mathcal{I}_Z(x) := \mathbb{E}_Z \left[\left(\frac{d}{dx} \log f(Z|x) \right)^2 \right]. \quad (4)$$

152 In our setting, the distribution $\mathbb{P}(\cdot|\boldsymbol{\eta}(x))$ is defined by the private data x , and Fisher information
 153 measures how much information is revealed about x through a sample $j \sim \mathbb{P}(\cdot|\boldsymbol{\eta}(x))$. It is noteworthy
 154 that such a reasoning has also been used to define Fisher information as a privacy metric [14].

155 **Lemma 1.** Let $M = \sup_{x \in \mathbb{R}} \mathcal{I}_Z(x)$ be an upper bound on the Fisher information of the mechanism
 156 \mathcal{M} . Then for any $x_1, x_2 \in \mathbb{R}$:

$$D_\alpha(\mathbb{P}(\cdot|\boldsymbol{\eta}(x_1)) \parallel \mathbb{P}(\cdot|\boldsymbol{\eta}(x_2))) \leq \alpha M (x_2 - x_1)^2 / 2. \quad (5)$$

¹In spite of this, we still name our mechanism I-MVU for its connection to the MVU mechanism.

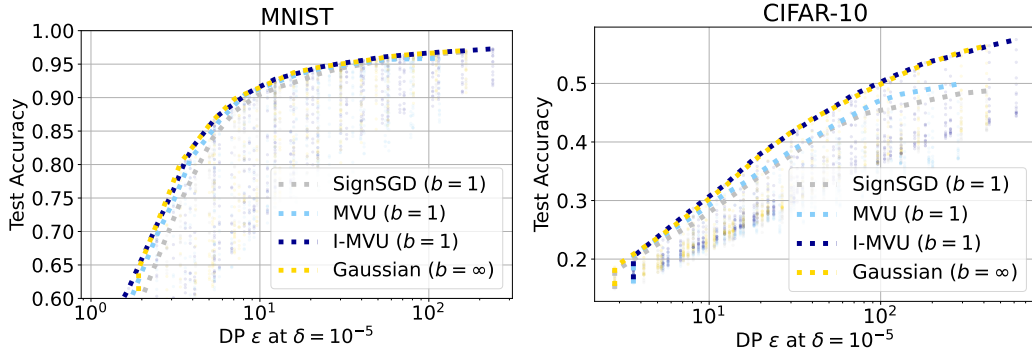


Figure 2: Privacy vs. accuracy plot for the client-level DP scenario on MNIST (left) and CIFAR-10 (right). Each point represents a single hyperparameter setting and the Pareto frontier is shown in dashed line. Across the entire range of ϵ , I-MVU consistently performs as well as the non-compressed Gaussian mechanism while requiring only one bit communication per update coordinate.

157 **Theorem 1.** Let M be the Fisher information constant from Lemma 1. Suppose that $\mathbf{x}_1, \mathbf{x}_2 \in \mathbb{R}^d$
 158 satisfy $\|\mathbf{x}_2 - \mathbf{x}_1\|_2 \leq C$. Then the I-MVU mechanism is $(\alpha, \alpha MC^2/2)$ -RDP for all $\alpha > 1$.

159 **Proof sketch.** We first derive the Taylor series expression for the Rényi divergence between
 160 $\mathbb{P}(\cdot|\boldsymbol{\eta}(x_1))$ and $\mathbb{P}(\cdot|\boldsymbol{\eta}(x_2))$. Since Rényi divergence is minimized and is equal to 0 when $x_1 = x_2$,
 161 the zeroth-order and first-order terms in the Taylor series are 0. The coefficient for the second-order
 162 term is given by the Fisher information $\mathcal{I}_Z(x_1)$ [15], and thus we give a numerical method to compute
 163 $M = \sup_{x \in \mathbb{R}} \mathcal{I}_Z(x)$ in Appendix B and use it in Equation 5 to bound the RDP ϵ . Full proofs of
 164 Lemma 1 and Theorem 1 are provided in Appendix A.

165 4 Experiments

166 We evaluate the I-MVU mechanism for federated learning under the local DP setting, *i.e.*, clients
 167 transmit the privately compressed model update $\mathcal{M}(\mathbf{x})$ to the server *before* aggregation. We consider
 168 private mechanisms that output *one bit per coordinate* of the update vector. This extreme level of
 169 compression reflects realistic constraints in FL and is very challenging for existing mechanisms.
 170 Previous work [7] found that the two most competitive baselines are the MVU mechanism with $b = 1$
 171 bit communication budget and SignSGD [17]. The latter applies the Gaussian mechanism for gradient
 172 perturbation [1] and then takes the sign of the perturbed gradient to obtain one-bit per coordinate.

173 4.1 Client-level DP

174 We first evaluate under the *client-level DP* setting on MNIST and CIFAR-10 [19]. Here, the privacy
 175 analysis guarantees that the learning algorithm is differentially private with respect to the removal of
 176 any client. We divide the training set among the clients with client sample size 1. Each client performs
 177 a single local gradient update in every FL round. This setting is equivalent to DP-SGD training [1]
 178 but with the Gaussian mechanism replaced by a communication-efficient private mechanism.

179 **Training details.** Following [7], we train a linear model on top of ScatterNet features [28], which
 180 remains to date the state-of-the-art DP model for MNIST and CIFAR-10 without leveraging any
 181 public data. We perform a grid search over hyperparameters such as number of update rounds, step
 182 size, gradient norm clip, and mechanism parameters σ (for Gaussian and SignSGD) and ϵ (for MVU
 183 and I-MVU). We use the same hyperparameter values reported in Tables 3 and 4 in [7].

184 **Result.** Figure 2 shows the privacy vs. test accuracy curve on MNIST (left) and CIFAR-10 (right).
 185 Privacy is measured in terms of (ϵ, δ) -DP at $\delta = 10^{-5}$. Each point in the scatter plot corresponds
 186 to a single hyperparameter setting and the dashed line shows the Pareto frontier of optimal privacy-
 187 accuracy trade-off. The yellow line corresponds to the standard Gaussian mechanism without
 188 compression, which attains the best test accuracy at any given privacy budget ϵ . Both SignSGD and
 189 MVU are competitive, achieving close to the same level of accuracy as the Gaussian mechanism,
 190 but a non-negligible gap remains, especially on CIFAR-10. In contrast, I-MVU attains nearly the
 191 same performance as the Gaussian mechanism at all values of ϵ on both MNIST and CIFAR-10.

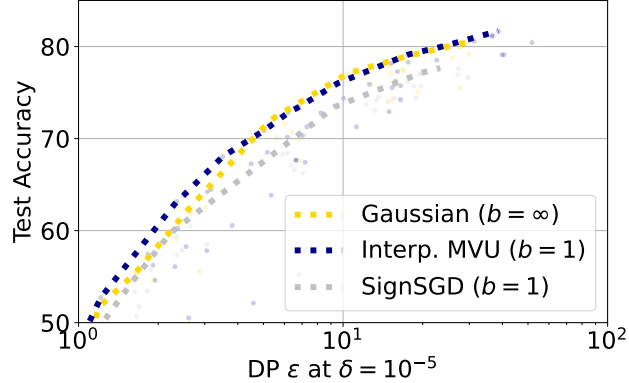


Figure 3: Privacy vs. accuracy plot for the sample-level DP scenario on FEMNIST. Each point represents a single hyperparameter setting and the Pareto frontier is shown in dashed line. I-MVU with one-bit communication budget per coordinate consistently performs better than SignSGD and is competitive with the non-compressed Gaussian baseline across the entire range of ϵ .

192 Since MVU and I-MVU are near-identical mechanisms, we argue that the performance gain comes
 193 primarily from the tight privacy analysis for L_2 geometry using Fisher information (Section 3.1).

194 4.2 Sample-level DP

195 Next, we evaluate under the *sample-level DP* setting on the FEMNIST dataset [5] for classifying
 196 written characters into 62 distinct classes. Privacy analysis guarantees that the learning algorithm is
 197 differentially private with respect to the removal of *any training sample from a client*. The dataset
 198 has a pre-defined train split with 3,500 clients, from which we randomly select 3,150 clients for
 199 training and the remaining 350 clients for testing. A set of 5 clients is selected in each training
 200 round, who then performs full batch gradient descent for a single local gradient update to compute
 201 the update vector. The update vector is privatized using a communication-efficient private mechanism
 202 and transmitted to the server.

203 **Training details.** We train a simple 4-layer convolutional network for classification. The model
 204 achieves 84% accuracy when trained non-privately. The client optimizer is SGD with a learning
 205 rate of 0.1 and no momentum. The server implements FedAvg [20] with a momentum of 0.9. We
 206 perform a grid search on the local and server learning rates, the clipping factor, the noise multiplier
 207 σ for both Gaussian and SignSGD baselines, and the ϵ and scale hyperparameters for I-MVU. The
 208 hyperparameter ranges are given in Tables 1, 2 and 3 in the appendix. In particular, SignSGD requires
 209 much lower server-side learning rates since the updates (in $\{\pm 1\}$) have higher magnitude.

210 **Result.** We show the privacy-accuracy trade-off for FEMNIST in Figure 3. Each point in the scatter
 211 plot represents a single hyperparameter setting and the Pareto frontier (dashed line) represents the
 212 optimal privacy-accuracy trade-off. The DP privacy budget ϵ is given at $\delta = 10^{-5}$. We observe
 213 that I-MVU (blue dashed line) performs better than SignSGD (silver dashed line) for the same
 214 communication budget of one bit per update coordinate across the entire range of considered privacy
 215 budgets ϵ . Moreover, I-MVU performs on par with the non-compressed Gaussian baseline (yellow
 216 dashed line), where clients perform local DP-SGD without compressing model updates.

217 5 Conclusion

218 We proposed the Interpolated MVU (I-MVU) mechanism that drastically reduces the amount of
 219 uplink communication in (cross-device) FL while providing differential privacy guarantees. Our
 220 proposal builds on the recently introduced MVU mechanism to extend it to continuous-valued vectors
 221 with L_2 geometry using a more efficient privacy analysis. Under both client-level and sample-level
 222 local DP settings, I-MVU with an extreme compression level of one bit per update coordinate attains
 223 close to the performance of the non-compressed Gaussian mechanism. Given this strong empirical
 224 performance, we advocate for I-MVU as a practical tool for communication-efficient private FL.

225 References

- 226 [1] Martin Abadi, Andy Chu, Ian Goodfellow, H Brendan McMahan, Ilya Mironov, Kunal Talwar,
227 and Li Zhang. Deep learning with differential privacy. In *Proceedings of the 2016 ACM SIGSAC*
228 *conference on computer and communications security*, pages 308–318, 2016.
- 229 [2] Naman Agarwal, Peter Kairouz, and Ziyu Liu. The skellam mechanism for differentially private
230 federated learning. *Advances in Neural Information Processing Systems*, 34:5052–5064, 2021.
- 231 [3] Borja Balle, Gilles Barthe, Marco Gaboardi, Justin Hsu, and Tetsuya Sato. Hypothesis test-
232 ing interpretations and renyi differential privacy. In *International Conference on Artificial*
233 *Intelligence and Statistics*, pages 2496–2506. PMLR, 2020.
- 234 [4] Borja Balle, Giovanni Cherubin, and Jamie Hayes. Reconstructing training data with informed
235 adversaries. *arXiv preprint arXiv:2201.04845*, 2022.
- 236 [5] Sebastian Caldas, Sai Meher Karthik Duddu, Peter Wu, Tian Li, Jakub Konečný, H Brendan
237 McMahan, Virginia Smith, and Ameet Talwalkar. Leaf: A benchmark for federated settings.
238 *arXiv preprint arXiv:1812.01097*, 2018.
- 239 [6] Clément L Canonne, Gautam Kamath, and Thomas Steinke. The discrete gaussian for differen-
240 tial privacy. *Advances in Neural Information Processing Systems*, 33:15676–15688, 2020.
- 241 [7] Kamalika Chaudhuri, Chuan Guo, and Mike Rabbat. Privacy-aware compression for federated
242 data analysis. *arXiv preprint arXiv:2203.08134*, 2022.
- 243 [8] Kamalika Chaudhuri, Claire Monteleoni, and Anand D Sarwate. Differentially private empirical
244 risk minimization. *Journal of Machine Learning Research*, 12(3), 2011.
- 245 [9] Wei-Ning Chen, Ayfer Ozgur, and Peter Kairouz. The poisson binomial mechanism for unbiased
246 federated learning with secure aggregation. In *International Conference on Machine Learning*,
247 pages 3490–3506. PMLR, 2022.
- 248 [10] Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam Smith. Calibrating noise to
249 sensitivity in private data analysis. In *Theory of cryptography conference*, pages 265–284.
250 Springer, 2006.
- 251 [11] Jonas Geiping, Hartmut Bauermeister, Hannah Dröge, and Michael Moeller. Inverting gradients-
252 how easy is it to break privacy in federated learning? *Advances in Neural Information Processing*
253 *Systems*, 33:16937–16947, 2020.
- 254 [12] Antonious Girgis, Deepesh Data, Suhas Diggavi, Peter Kairouz, and Ananda Theertha Suresh.
255 Shuffled model of differential privacy in federated learning. In *International Conference on*
256 *Artificial Intelligence and Statistics*, pages 2521–2529. PMLR, 2021.
- 257 [13] Chuan Guo, Brian Karrer, Kamalika Chaudhuri, and Laurens van der Maaten. Bounding training
258 data reconstruction in private (deep) learning. *arXiv preprint arXiv:2201.12383*, 2022.
- 259 [14] Awni Hannun, Chuan Guo, and Laurens van der Maaten. Measuring data leakage in machine-
260 learning models with fisher information. In *Uncertainty in Artificial Intelligence*, pages 760–770.
261 PMLR, 2021.
- 262 [15] David Haussler and Manfred Opper. Mutual information, metric entropy and cumulative relative
263 entropy risk. *The Annals of Statistics*, 25(6):2451–2492, 1997.
- 264 [16] Jinwoo Jeon, Kangwook Lee, Sewoong Oh, Jungseul Ok, et al. Gradient inversion with
265 generative image prior. *Advances in Neural Information Processing Systems*, 34:29898–29908,
266 2021.
- 267 [17] Richeng Jin, Yufan Huang, Xiaofan He, Huaiyu Dai, and Tianfu Wu. Stochastic-sign sgd for
268 federated learning with theoretical guarantees. *arXiv preprint arXiv:2002.10940*, 2020.

- 269 [18] Peter Kairouz, H Brendan McMahan, Brendan Avent, Aurélien Bellet, Mehdi Bennis, Ar-
270 jun Nitin Bhagoji, Keith Bonawitz, Zachary Charles, Graham Cormode, Rachel Cummings,
271 Rafael GL D'Oliveira, Hubert Eichner, Salim El Rouayheb, David Evans, Josh Gardner, Zachary
272 Garrett, Adrià Gascón, Badih Ghazi, Phillip B Gibbons, Marco Gruteser, Zaid Harchaoui,
273 Chaoyang He, Lie He, Zhouyuan Huo, Ben Hutchinson, Justin Hsu, Martin Jaggi, Tara Javidi,
274 Gauri Joshi, Mikhail Khodak, Jakub Konečný, Aleksandra Korolova, Farinaz Koushanfar, Sanmi
275 Koyejo, Tancrède Lepoint, Yang Liu, Prateek Mittal, Mehryar Mohri, Richard Nock, Ayfer
276 Özgür, Rasmus Pagh, Hang Qi, Daniel Ramage, Ramesh Raskar, Mariana Raykova, Dawn
277 Song, Weikang Song, Sebastian U Stich, Ziteng Sun, Ananda Theertha Suresh, Florian Tramèr,
278 Praneeth Vepakomma, Jianyu Wang, Li Xiong, Zheng Xu, Qiang Yang, X Yu Felix, Han Yu,
279 and Sen Zhao. Advances and open problems in federated learning. *Foundations and Trends in
280 Machine Learning*, 14(1–2):1–210, 2021.
- 281 [19] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images.
282 2009.
- 283 [20] Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas.
284 Communication-efficient learning of deep networks from decentralized data. In *Artificial
285 intelligence and statistics*, pages 1273–1282. PMLR, 2017.
- 286 [21] Ilya Mironov. Rényi differential privacy. In *2017 IEEE 30th computer security foundations
287 symposium (CSF)*, pages 263–275. IEEE, 2017.
- 288 [22] Frank Nielsen and Richard Nock. On Rényi and tsallis entropies and divergences for exponen-
289 tial families. *arXiv preprint arXiv:1105.3259*, 2011.
- 290 [23] Alfréd Rényi et al. On measures of entropy and information. In *Proceedings of the fourth
291 Berkeley symposium on mathematical statistics and probability*, volume 1. Berkeley, California,
292 USA, 1961.
- 293 [24] Abhin Shah, Wei-Ning Chen, Johannes Balle, Peter Kairouz, and Lucas Theis. Optimal
294 compression of locally differentially private mechanisms. In *International Conference on
295 Artificial Intelligence and Statistics*, pages 7680–7723. PMLR, 2022.
- 296 [25] Reza Shokri and Vitaly Shmatikov. Privacy-preserving deep learning. In *Proceedings of the
297 22nd ACM SIGSAC conference on computer and communications security*, pages 1310–1321,
298 2015.
- 299 [26] Reza Shokri, Marco Stronati, Congzheng Song, and Vitaly Shmatikov. Membership inference
300 attacks against machine learning models. In *2017 IEEE symposium on security and privacy
301 (SP)*, pages 3–18. IEEE, 2017.
- 302 [27] Pierre Stock, Igor Shilov, Ilya Mironov, and Alexandre Sablayrolles. Defending against
303 reconstruction attacks with Rényi differential privacy. *arXiv preprint arXiv:2202.07623*, 2022.
- 304 [28] Florian Tramèr and Dan Boneh. Differentially private learning needs better features (or much
305 more data). *arXiv preprint arXiv:2011.11660*, 2020.
- 306 [29] Aleksei Triastcyn, Matthias Reisser, and Christos Louizos. Dp-rec: Private & communication-
307 efficient federated learning. *arXiv preprint arXiv:2111.05454*, 2021.
- 308 [30] Hongxu Yin, Arun Mallya, Arash Vahdat, Jose M Alvarez, Jan Kautz, and Pavlo Molchanov.
309 See through gradients: Image batch recovery via gradinversion. In *Proceedings of the IEEE/CVF
310 Conference on Computer Vision and Pattern Recognition*, pages 16337–16346, 2021.
- 311 [31] Bo Zhao, Konda Reddy Mopuri, and Hakan Bilen. idlg: Improved deep leakage from gradients.
312 *arXiv preprint arXiv:2001.02610*, 2020.
- 313 [32] Ligeng Zhu, Zhijian Liu, and Song Han. Deep leakage from gradients. *Advances in neural
314 information processing systems*, 32, 2019.

315 **A Proofs**

316 **Lemma 1.** Let $M = \sup_{x \in \mathbb{R}} \mathcal{I}_Z(x)$ be an upper bound on the Fisher information of the mechanism
 317 \mathcal{M} . Then for any $x_1, x_2 \in \mathbb{R}$:

$$D_\alpha(\mathbb{P}(\cdot|\boldsymbol{\eta}(x_1)) \parallel \mathbb{P}(\cdot|\boldsymbol{\eta}(x_2))) \leq \alpha M(x_2 - x_1)^2/2.$$

318 *Proof.* We first derive an explicit form for the Fisher information. Let $f(\mathbf{z}; x)$ denote the pmf in
 319 Equation 2 for any $\mathbf{z} \in \{\mathbf{e}_1, \dots, \mathbf{e}_{B_{\text{out}}}\}$. The log pmf is:

$$\log f(\mathbf{z}; x) = \mathbf{z}^\top \boldsymbol{\eta}(x) - A(\boldsymbol{\eta}(x)) \quad (6)$$

320 Taking derivative with respect to x gives:

$$\begin{aligned} \frac{d}{dx} \log f(\mathbf{z}; x) &= (\mathbf{z} - \sigma(\boldsymbol{\eta}(x)))^\top (\boldsymbol{\eta}_2 - \boldsymbol{\eta}_1), \\ \left(\frac{d}{dx} \log f(\mathbf{z}; x) \right)^2 &= (\boldsymbol{\eta}_2 - \boldsymbol{\eta}_1)^\top (\mathbf{z} - \sigma(\boldsymbol{\eta}(x))) (\mathbf{z} - \sigma(\boldsymbol{\eta}(x)))^\top (\boldsymbol{\eta}_2 - \boldsymbol{\eta}_1), \end{aligned}$$

321 where σ denotes the sigmoid function. Taking expectation over \mathbf{z} gives the Fisher information:

$$\mathcal{I}_Z(x) = (\boldsymbol{\eta}_2 - \boldsymbol{\eta}_1)^\top U (\boldsymbol{\eta}_2 - \boldsymbol{\eta}_1), \quad (7)$$

322 with $U = \text{diag}(\sigma(\boldsymbol{\eta}(x))) - \sigma(\boldsymbol{\eta}(x))\sigma(\boldsymbol{\eta}(x))^\top$.

323 To derive the upper bound, we first define a function F for the Rényi divergence of the mechanism
 324 for a fixed x_1 and varying x_2 :

$$F_\alpha(x_2; x_1) = D_\alpha(\mathbb{P}(\cdot|\boldsymbol{\eta}(x_1)) \parallel \mathbb{P}(\cdot|\boldsymbol{\eta}(x_2))). \quad (8)$$

325 By Taylor's theorem, we can express F as:

$$F_\alpha(x_2; x_1) = F_\alpha(x_1; x_1) + (x_2 - x_1)F'_\alpha(x_1; x_1) + (x_2 - x_1)^2 F''_\alpha(\xi; x_1)/2$$

326 for some $\xi \in [x_1, x_2]$. Note that $F_\alpha(x_1; x_1) = 0$ and $F'_\alpha(x_1; x_1) = 0$ (since x_1 is the global
 327 minimizer of $F_\alpha(\cdot; x_1)$), so F is locally a quadratic function:

$$F_\alpha(x_2; x_1) = (x_2 - x_1)^2 F''_\alpha(\xi; x_1)/2. \quad (9)$$

328 Since f is the pmf of an exponential family distribution, we can use the closed form expression [22]
 329 for Rényi divergence of exponential family distributions to express F and its derivatives:

$$\begin{aligned} F_\alpha(\xi; x_1) &= \frac{1}{\alpha - 1} [A(\alpha\boldsymbol{\eta}(x_1) + (1 - \alpha)\boldsymbol{\eta}(\xi)) - \alpha A(\boldsymbol{\eta}(x_1)) - (1 - \alpha)A(\boldsymbol{\eta}(\xi))] \\ F'_\alpha(\xi; x_1) &= (\sigma(\boldsymbol{\eta}(\xi)) - \sigma(\alpha\boldsymbol{\eta}(x_1) + (1 - \alpha)\boldsymbol{\eta}(\xi)))^\top (\boldsymbol{\eta}_2 - \boldsymbol{\eta}_1) \\ F''_\alpha(\xi; x_1) &= (\boldsymbol{\eta}_2 - \boldsymbol{\eta}_1)^\top (V + (\alpha - 1)W) (\boldsymbol{\eta}_2 - \boldsymbol{\eta}_1) \end{aligned}$$

330 where $V = \text{diag}(\sigma(\boldsymbol{\eta}(\xi))) - \sigma(\boldsymbol{\eta}(\xi))\sigma(\boldsymbol{\eta}(\xi))^\top$, $W = \text{diag}(\sigma(\boldsymbol{\eta}(x_1))) - \sigma(\boldsymbol{\eta}(x_1))\sigma(\boldsymbol{\eta}(x_1))^\top$, and
 331 $x' = \alpha x_1 + (1 - \alpha)\xi$. Hence $F''_\alpha(\xi; x_1) = \mathcal{I}_Z(\xi) + (\alpha - 1)\mathcal{I}_Z(x_1)$ by Equation 7. Upper bounding
 332 $\mathcal{I}_Z(\xi)$ and $\mathcal{I}_Z(x_1)$ by $M := \sup_{x \in \mathbb{R}} \mathcal{I}_Z(x)$ and combining with Equation 9 gives the desired result.
 333 \square

334 **Theorem 1.** Let M be the Fisher information constant from Lemma 1. Suppose that $\mathbf{x}_1, \mathbf{x}_2 \in \mathbb{R}^d$
 335 satisfy $\|\mathbf{x}_2 - \mathbf{x}_1\|_2 \leq C$, then the I-MVU mechanism is $(\alpha, \alpha M C^2 / 2)$ -RDP for all $\alpha > 1$.

336 *Proof.* Let $\mathbf{a} = \mathcal{M}(\mathbf{x}) \in \{a_1, \dots, a_{B_{\text{out}}}\}^d$ be the output of the vector I-MVU mechanism that
 337 independently applies the scalar mechanism to each coordinate. Then:

$$\begin{aligned}
 D_\alpha(\mathcal{M}(\mathbf{x}_1) \parallel \mathcal{M}(\mathbf{x}_2)) &= \frac{1}{\alpha - 1} \log \sum_{\mathbf{a} \in \{a_1, \dots, a_{B_{\text{out}}}\}^d} \prod_{k=1}^d \frac{\mathbb{P}(\mathbf{a}_k | \boldsymbol{\eta}((\mathbf{x}_2)_k))^\alpha}{\mathbb{P}(\mathbf{a}_k | \boldsymbol{\eta}((\mathbf{x}_1)_k))^\alpha} \\
 &= \sum_{k=1}^d \frac{1}{\alpha - 1} \log \sum_{\mathbf{a}_k \in \{a_1, \dots, a_{B_{\text{out}}}\}} \frac{\mathbb{P}(\mathbf{a}_k | \boldsymbol{\eta}((\mathbf{x}_2)_k))^\alpha}{\mathbb{P}(\mathbf{a}_k | \boldsymbol{\eta}((\mathbf{x}_1)_k))^\alpha} \quad \text{by independence} \\
 &= \sum_{k=1}^d D_\alpha(\mathbb{P}(\cdot | \boldsymbol{\eta}((\mathbf{x}_1)_k)) \parallel \mathbb{P}(\cdot | \boldsymbol{\eta}((\mathbf{x}_2)_k))) \\
 &\leq \sum_{k=1}^d \alpha M((\mathbf{x}_2)_k - (\mathbf{x}_1)_k)^2 / 2 \quad \text{by Lemma 1} \\
 &= \alpha M \|\mathbf{x}_2 - \mathbf{x}_1\|_2^2 / 2.
 \end{aligned}$$

338

□

339 B Computing Fisher Information

340 In this section we describe a method for computing $M = \sup_{x \in \mathbb{R}} \mathcal{I}_Z(x)$. We first define a condi-
 341 tion for $\boldsymbol{\eta}_1, \boldsymbol{\eta}_2$ that allows us to reduce this problem to maximizing $\mathcal{I}_Z(x)$ over a bounded range
 342 $[x_{\min}, x_{\max}]$.

343 **Definition 2.** Two vectors $\boldsymbol{\eta}_1, \boldsymbol{\eta}_2 \in \mathbb{R}^B$ are said to be anadromic if for all $j = 1, \dots, B$, we have
 344 $(\boldsymbol{\eta}_1)_j = (\boldsymbol{\eta}_2)_{B-j+1}$.

345 The following technical lemma proves several useful properties that hold when $\boldsymbol{\eta}_1$ and $\boldsymbol{\eta}_2$ are
 346 anadromic.

347 **Lemma 2.** Suppose that $\boldsymbol{\eta}_1, \boldsymbol{\eta}_2 \in \mathbb{R}^B$ are anadromic. Let $\boldsymbol{\theta} = \boldsymbol{\eta}_2 - \boldsymbol{\eta}_1$ and suppose that $j^+ =$
 348 $\arg \max_j \boldsymbol{\theta}_j, j^- = \arg \min_j \boldsymbol{\theta}_j$ are unique. Then the following hold:

- 349 (i) $\boldsymbol{\theta}_j = -\boldsymbol{\theta}_{B-j+1}$ for all j , and hence $j^- = B - j^+ + 1$.
- 350 (ii) $\boldsymbol{\eta}(x)_j = \boldsymbol{\eta}(1-x)_{B-j+1}$ for all j .
- 351 (iii) $\sigma(\boldsymbol{\eta}(x)) \rightarrow \mathbf{e}_{j^+}$ as $x \rightarrow \infty$ and $\sigma(\boldsymbol{\eta}(x)) \rightarrow \mathbf{e}_{j^-}$ as $x \rightarrow -\infty$.
- 352 (iv) $\mathcal{I}_Z(x) = \mathcal{I}_Z(1-x)$ for all $x \in \mathbb{R}$.
- 353 (v) $x = 1/2$ is a stationary point for $\mathcal{I}_Z(x)$.
- 354 (vi) If $\sigma(\boldsymbol{\eta}(x))_{j^+} \geq 1/2$ then $\mathcal{I}_Z(x) \leq 4\boldsymbol{\theta}_{j^+}^2 \sigma(\boldsymbol{\eta}(x))_{j^+} (1 - \sigma(\boldsymbol{\eta}(x))_{j^+})$.

355 *Proof.* (i) Since $\boldsymbol{\eta}_1$ and $\boldsymbol{\eta}_2$ are anadromic,

$$\boldsymbol{\theta}_j = (\boldsymbol{\eta}_2)_j - (\boldsymbol{\eta}_1)_j = (\boldsymbol{\eta}_2)_j - (\boldsymbol{\eta}_2)_{B-j+1} = -((\boldsymbol{\eta}_2)_{B-j+1} - (\boldsymbol{\eta}_1)_{B-j+1}) = -\boldsymbol{\theta}_{B-j+1}.$$

356 In particular, $\arg \max_j \boldsymbol{\theta}_j = B - (\arg \min_j \boldsymbol{\theta}_j) + 1$.

357 (ii) $\boldsymbol{\eta}(x)_j = (1-x)(\boldsymbol{\eta}_1)_j + x(\boldsymbol{\eta}_2)_j = (1-x)(\boldsymbol{\eta}_1)_{B-j+1} + x(\boldsymbol{\eta}_2)_{B-j+1} = \boldsymbol{\eta}(1-x)_{B-j+1}$.

358 (iii) Let $\bar{\boldsymbol{\eta}} = (\boldsymbol{\eta}_1 + \boldsymbol{\eta}_2)/2$ so that $\boldsymbol{\eta}(x) = \bar{\boldsymbol{\eta}} + (x-1/2)\boldsymbol{\theta}$ for all $x \in \mathbb{R}$. It is clear that $\sigma(\boldsymbol{\eta}(x)) \rightarrow \mathbf{e}_{j^+}$
 359 as $x \rightarrow \infty$ since j^+ is unique. A similar argument shows that $\sigma(\boldsymbol{\eta}(x)) \rightarrow \mathbf{e}_{j^-}$ as $x \rightarrow -\infty$.

360 (iv) Using the expression of $\mathcal{I}_Z(x)$ in the proof of Lemma 1, we get

$$\begin{aligned}
 \mathcal{I}_Z(x) &= \sum_j \boldsymbol{\theta}_j^2 \sigma(\boldsymbol{\eta}(x))_j - \left(\sum_j \boldsymbol{\theta}_j \sigma(\boldsymbol{\eta}(x))_j \right)^2 \\
 &= \sum_j \boldsymbol{\theta}_{B-j+1}^2 \sigma(\boldsymbol{\eta}(1-x))_{B-j+1} - \left(\sum_j \boldsymbol{\theta}_{B-j+1} \sigma(\boldsymbol{\eta}(1-x))_{B-j+1} \right)^2 \quad \text{by (i) and (ii)} \\
 &= \mathcal{I}_Z(1-x).
 \end{aligned}$$

361 (v) Differentiating $\mathcal{I}_Z(x)$ and using the above argument gives:

$$\begin{aligned}\mathcal{I}'_Z(x) &= (\boldsymbol{\theta}^3)^\top \sigma(\boldsymbol{\eta}(x)) - 3\boldsymbol{\theta}^\top \sigma(\boldsymbol{\eta}(x))(\boldsymbol{\theta}^2)^\top \sigma(\boldsymbol{\eta}(x)) + 2(\boldsymbol{\theta}^\top \sigma(\boldsymbol{\eta}(x)))^3 \\ &= -(\boldsymbol{\theta}^3)^\top \sigma(\boldsymbol{\eta}(1-x)) + 3\boldsymbol{\theta}^\top \sigma(\boldsymbol{\eta}(1-x))(\boldsymbol{\theta}^2)^\top \sigma(\boldsymbol{\eta}(1-x)) - 2(\boldsymbol{\theta}^\top \sigma(\boldsymbol{\eta}(1-x)))^3 \\ &= -\mathcal{I}'_Z(1-x).\end{aligned}$$

362 Then $\mathcal{I}'_Z(1/2) = -\mathcal{I}'_Z(1/2)$, so $\mathcal{I}'_Z(1/2) = 0$ and $x = 1/2$ is a stationary point.

(vi) Using the fact that $0 \leq \boldsymbol{\theta}_j^2 \leq \boldsymbol{\theta}_{j+}^2$ for all j and

$$\sum_j \boldsymbol{\theta}_j \sigma(\boldsymbol{\eta}(x))_j \geq \boldsymbol{\theta}_{j+} \sigma(\boldsymbol{\eta}(x))_{j+} + \boldsymbol{\theta}_{j-} (1 - \sigma(\boldsymbol{\eta}(x))_{j+}) = \boldsymbol{\theta}_{j+} \sigma(\boldsymbol{\eta}(x))_{j+} - \boldsymbol{\theta}_{j+} (1 - \sigma(\boldsymbol{\eta}(x))_{j+}) \geq 0,$$

363 we have:

$$\begin{aligned}\mathcal{I}_Z(x) &= \sum_j \boldsymbol{\theta}_j^2 \sigma(\boldsymbol{\eta}(x))_j - \left(\sum_j \boldsymbol{\theta}_j \sigma(\boldsymbol{\eta}(x))_j \right)^2 \\ &\leq \boldsymbol{\theta}_{j+}^2 - (\boldsymbol{\theta}_{j+} \sigma(\boldsymbol{\eta}(x))_{j+} - \boldsymbol{\theta}_{j+} (1 - \sigma(\boldsymbol{\eta}(x))_{j+}))^2 \\ &= \boldsymbol{\theta}_{j+}^2 (1 - (2\sigma(\boldsymbol{\eta}(x))_{j+} - 1)^2) \\ &= 4\boldsymbol{\theta}_{j+}^2 \sigma(\boldsymbol{\eta}(x))_{j+} (1 - \sigma(\boldsymbol{\eta}(x))_{j+}).\end{aligned}$$

364

□

365 **Algorithm.** To use Lemma 2 to compute M , we first compute $I^* = \mathcal{I}_Z(1/2)$ since $x = 1/2$ is a
366 stationary point by Lemma 2(v). By setting

$$4\boldsymbol{\theta}_{j+}^2 \sigma(\boldsymbol{\eta}(x))_{j+} (1 - \sigma(\boldsymbol{\eta}(x))_{j+}) \leq I^*$$

367 and solving this quadratic equation for $\sigma(\boldsymbol{\eta}(x))_{j+}$, we can use the bound in Lemma 2(vi) to obtain
368 that $\mathcal{I}_Z(x) \leq I^*$ when $\sigma(\boldsymbol{\eta}(x))_{j+} \geq \left(1 + \sqrt{1 - I^*/\boldsymbol{\theta}_{j+}^2}\right)/2 \geq 1/2$. Since $\sigma(\boldsymbol{\eta}(x))_{j+} \rightarrow 1$ as
369 $x \rightarrow \infty$ by (iii), we can determine the value x_{\max} for which $\mathcal{I}_Z(x) \leq I^*$ when $x \geq x_{\max}$. By (iv),
370 $x_{\min} = 1 - x_{\max}$ satisfies $\mathcal{I}_Z(x) \leq I^*$ when $x \leq x_{\min}$. We can then do line search in $[x_{\min}, x_{\max}]$
371 (or equivalently, in $[1/2, x_{\max}]$ by Lemma 2(iv)) to obtain M .

372 C Hyperparameters for FEMNIST

Hyperparameter	Values
ε	0.25, 0.5, 0.75, 1, 2, 3, 5, 6, 7, 8, 9, 10
Server-side learning rate	0.5, 1, 2
Scaling factor β	32, 64, 128

Table 1: Hyperparameter range for I-MVU on FEMNIST.

Hyperparameter	Values
σ	0.6, 0.8, 1, 2, 4, 6, 8, 10, 16, 32, 64, 128
Server-side learning rate	0.5, 1, 2
Clipping factor	0.5, 1, 2

Table 2: Hyperparameter range for Gaussian on FEMNIST.

Hyperparameter	Values
σ	0.6, 0.8, 1, 2, 4, 6, 8, 10, 16, 32, 64, 128
Server-side learning rate	0.0001, 0.001, 0.01
Clipping factor	0.5, 1, 2

Table 3: Hyperparameter range for SignSGD on FEMNIST.