
Foreseeing Privacy Threats from Gradient Inversion Through the Lens of Angular Lipschitz Smoothness

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Recent works proposed server-side *input recovery attacks* in federated learning
2 (FL), in which an honest-but-curious server can recover clients' data (e.g., images)
3 using shared model gradients, thus raising doubts regarding the safety of FL.
4 However, the attack methods are typically demonstrated on only a few models or
5 focus heavily on the reconstruction of a single image, which is easier than that of a
6 batch (multiple images). Thus, in this study, we systematically re-evaluated state-
7 of-the-art (SOTA) attack methods on a variety of models in the context of batch
8 reconstruction. For a broad spectrum of models, we considered two types of model
9 variations: *implicit* (i.e., *without any change in architecture*) and *explicit* (i.e., *with*
10 *architectural changes*). Motivated by the re-evaluation results that the quality of
11 reconstructed image batch differs per model, we propose *angular Lipschitz constant*
12 *of a model gradient function with respect to an input* as a measure that explains
13 the vulnerability of a model against input recovery attacks. The prototype of the
14 proposed measure is derived from our theorem on the convergence of attackers'
15 gradient matching optimization, and re-designed into the *scale-invariant* form to
16 prevent trivial server-side loss scaling trick. We demonstrated the predictability of
17 the proposed measure on the vulnerability under recovery attacks by empirically
18 showing its strong monotonic correlation with not only loss drop during gradient-
19 matching optimization but also the quality of the reconstructed image batch. We
20 expect our measure to be a key factor for developing client-side defensive strategies
21 against privacy threats in our proposed realistic FL setting called *black-box* setting,
22 where the server deliberately conceals global model information from clients
23 excluding model gradients.

24 1 Introduction

25 Federated learning (FL) is a cooperative machine learning between clients as local trainers and
26 a central server as a global aggregator [14, 21]. Participants in FL cannot access raw data from
27 others and only communicate with one another through gradients, which were believed to leak little
28 information of the original data in the past.

29 However, recent studies [31, 30, 6, 26, 12] challenge *inverting* gradients back to original data,
30 suggesting that there is potential for an honest-but-curious server to attack by sneakily recovering
31 clients' data from gradients in FL. Their algorithms, so-called gradient inversion attacks, aim at
32 optimizing input variables (e.g., images) to match the given gradients under the condition of fixed
33 model weights. For better reconstruction quality, state-of-the-art (SOTA) attacks assume that both
34 batch normalization (BN) [11] layers' statistics and private labels are known [6, 26, 12, 8]. However,

35 they are demonstrated on a limited range of global models. Thus, we systematically re-evaluated
36 SOTA gradient inversion attacks on a variety of models in the context of *batch (or multiple images)*
37 *reconstruction*, the recovery of input batch from the averaged gradients over itself, which is more
38 difficult to solve than single image reconstruction, the recovery of single image from its gradient. In
39 this paper, two kinds of model variations are considered, namely *implicit* and *explicit*.

40 Implicit model variations refer to a collection of different models with the same architecture. In this
41 paper, we consider two types of implicit model variations: *BN modes* and *training epochs*.

42 • As mentioned previously, SOTA gradient inversion attack methods are demonstrated on models
43 with BN layers to assume shared BN statistics. Note that there are two modes of a BN layer,
44 namely, *train mode* and *eval mode*. In the reality of FL, the server can choose any mode among
45 them. Therefore, we re-evaluated SOTA attacks by considering both modes of BN. This paper is
46 the first to consider BN modes for the evaluation of gradient inversion attacks. We empirically
47 found that the quality of reconstructed batch significantly changes by switching BN modes even
48 for the same model weights.

49 • By reflecting the reality that clients can encounter global model from the server at any time, we
50 consider models with different training epochs for the re-evaluation. This scheme extends the
51 scope of previous works’ training epoch choices of *black-and-white* manner: zero training epoch
52 (untrained) and maximum training epochs (fully trained). We empirically found that the best
53 reconstruction result was usually found at earlier training epochs, not untrained nor fully trained,
54 thus raising the need to expand the evaluation criterion for attack methods.

55 Meanwhile, explicit model variations are more straightforward than implicit model variations as they
56 only involve architectural changes. In this study, we consider two types of explicit model variations:
57 *skip connections* and *channel size*.

58 • Residual networks (ResNets) [9] are frequently employed in previous works [26, 6, 31, 12] even
59 for batch reconstruction, while networks *without skip connection* are introduced for only for the
60 recovery of single image from its gradient [6]. Therefore, we explored how a skip connection
61 affects the quality of SOTA gradient inversion attacks in the context of batch reconstruction. Our
62 empirical findings suggest that models without skip connection are more robust against the gradient
63 inversion attack than residual networks.

64 • The reconstruction quality is known to increase with the number of channels, but this property is
65 demonstrated on single image reconstruction [30, 6]. Thus, we recap how the number of channels
66 affects the attack quality in the context of batch reconstruction.

67 By re-evaluating SOTA attacks in a variety of models, we found that the vulnerability against gradient
68 inversion attack significantly differs per model, implying the need of more strict evaluation criteria
69 for attack methods. Then, clients are required to judge whether a shared model from the server is safe
70 or not *before sending* locally computed gradients back for their privacy. In this study, we consider
71 two settings on the transparency of global model information to clients: *white-box* and *black-box*.
72 In a white-box setting, clients have an absolute control over global model such as the server; thus,
73 clients can directly apply SOTA attacks to the model to assess its vulnerability.

74 On the other hand, a *black-box* setting only allows clients control over model gradients to restrict
75 access to the global model possibly due to companies’ secrets. For the client-side measurement of
76 privacy leakage in this practical and difficult setting, we propose *angular Lipschitz constant of model*
77 *gradients with respect to an input* as a predictive measure for the quality of reconstructed samples
78 inverted from model gradients.

79 This measure is derived from our theorem in Sec. 4 that an attacker’s gradient matching loss function
80 drops more abruptly with a smaller L in a particular range, where L is Lipschitz constant of model
81 gradients with respect to an input. However, using L as a measure for privacy leakage would be
82 inappropriate as L can be any nonnegative value by loss function scaling. Therefore, inspired by
83 scale-invariant cosine similarity loss function, we propose the angular Lipschitz constant, a *loss*
84 *scaling-invariant* alternative to L . We experimentally found that both measure monotonically correlates

85 with not only total loss drop during an attacker’s optimization but also the reconstruction quality
 86 than the norm of gradients. These findings are expected to support the construction of client-side
 87 defense algorithms particularly for *black-box* setting, where only model gradients are given to clients
 88 as minimal information of the model as described in Fig. 5.

89 2 Prior Art in the Gradient Inversion Attack

90 Given the neural network function $f_w : \mathbb{R}^{b \times d} \rightarrow \mathbb{R}^{b \times c}$ (w, b, d, c being the model weights, batch size,
 91 image size, and the number of classes, respectively), and the gradient $g^* = \frac{\partial \mathcal{L}(f_w(x^*), y^*)}{\partial w}$ computed
 92 with ground truth input batch $(x^*, y^*) \in \mathbb{R}^{b \times d} \times \mathbb{R}^b$ (x^*, y^* being the image batch, and corresponding
 93 label batch) and the loss function $\mathcal{L} : \mathbb{R}^{b \times c} \times \mathbb{R}^b \rightarrow \mathbb{R}$ (e.g., cross-entropy loss), the goal of gradient
 94 inversion attack is to reconstruct an image batch $x \in \mathbb{R}^{b \times d}$, a resemblance of ground truth image
 95 batch x^* . In the context of federated learning (FL), f_w is the global model, and g^* is the gradient
 96 computed from a client. Then, a honest-but-curious server aims to recover the client’s private data x^* .

97 A general method to tackle the problem of inverting gradients is to solve an optimization problem
 98 formulated as follows:

$$\arg \min_{x, y} \mathcal{L}_{grad} \left(\frac{\partial \mathcal{L}(f_w(x), y)}{\partial w}, \frac{\partial \mathcal{L}(f_w(x^*), y^*)}{\partial w} \right) + \alpha_{prior} \mathcal{R}_{prior}(x), \quad (1)$$

99 where $\mathcal{L}_{grad} : \mathbb{R}^N \times \mathbb{R}^N \rightarrow \mathbb{R}$ (N is the size of weights w) is the loss function for gradient matching
 100 (which closes the distance between current gradients and target gradients), $\mathcal{R}_{prior} : \mathbb{R}^{b \times d} \rightarrow \mathbb{R}$ is the
 101 regularization loss for image prior, with α_{prior} being its coefficient.

102 Prior to the advent of packages for automatic differentiation, the gradient term $g = \frac{\partial \mathcal{L}(f_w(x), y)}{\partial w}$ was
 103 computed as a function of (x, y) in a closed form. For the computation to be tractable, \mathcal{L}_{grad} was
 104 set to a squared loss ($\mathcal{L}(g, g^*) = \|g - g^*\|_2^2$), and f_w was also slightly modified from the original
 105 design of contemporary neural networks. For example, ReLU activation functions were replaced with
 106 Sigmoid, and all the strides in convolution modules were excluded from the original ResNet in [31].
 107 Consequently, the choice of f_w was limited.

108 Currently, with the advantages of automatic differentiation [22] and advanced deep learning opti-
 109 mization algorithms [13, 23, 5], solving for optimization problem in (1) becomes tractable for most
 110 contemporary deep neural networks without the need for modification. Further, the gradient matching
 111 loss is selected in a broad range from cosine similarity loss ($\mathcal{L}(g, g^*) = 1 - \frac{\langle g, g^* \rangle}{\|g\| \|g^*\|}$) [6, 12, 10, 26]
 112 to L2 loss ($\mathcal{L}(g, g^*) = \|g - g^*\|_2^2$) [31, 29, 26]. The liberation from the limited choice of loss
 113 functions and neural network architectures became the trigger of state-of-the-art attack methods.

114 State-of-the-art attack methods provide several assumptions which enable the baseline, which is only
 115 gradient-based, to be expanded.

116 First, the server is supposed to know the private labels of clients’ images. Currently, estimating x
 117 and y becomes a sequential process, in which y is estimated first, after which x is approximated with
 118 the estimated $y = y_{approx}^*$ given. Rather than jointly learning x and y in (1), prior works suggest
 119 estimating y directly by seeing the gradients from ground truth data g^* before optimization [26, 29].
 120 Therefore, the problem of estimating labels from gradients is separated from the original optimization
 121 problem in (1) [3, 25, 16] and some works, which focus on reconstruction of images rather than
 122 labels, assume that private labels are known [6, 12].

123 Second, the local batch statistics $\{\mu_l(x^*; w), \sigma_l^2(x^*; w)\}_{l=1}^M$ ($\mu_l(x^*; w)$, $\sigma_l(x^*; w)$, and M being
 124 the batch mean of the l^{th} batch normalization layer, batch standard deviation of the l^{th} batch
 125 normalization (BN) layer, and number of the BN layers, respectively), computed with client’s data
 126 batch, is given to the server. This assumption reflects a naive approach of a FL algorithm called
 127 FedAvg [21] on the global model with BN layers [19, 17]. When $\{\mu_l(x^*; w), \sigma_l(x^*; w)^2\}_{l=1}^M$ is
 128 shared from a client to the server for the update of population statistics in the global model’s BN
 129 layers, the server as an honest-but-curious adversary would work to add up the batch statistics
 130 matching loss term to (1) to ensure a stronger attack.

131 Then, optimization problem in (1) can be rewritten by considering both assumptions mentioned
 132 previously as follows:

$$\arg \min_x \mathcal{L}_{grad} \left(\frac{\partial \mathcal{L}(f_w(x), y^*)}{\partial w}, \frac{\partial \mathcal{L}(f_w(x^*), y^*)}{\partial w} \right) + \alpha_{prior} \mathcal{R}_{prior}(x) + \alpha_{BN} \sum_{l=1}^M \mathcal{R}_{BN}((\mu_l, \sigma_l^2), (\mu_l^*, \sigma_l^{*2})) \quad (2)$$

133 , where \mathcal{R}_{BN} is the BN statistics matching loss and α_{BN} being its coefficient with $(\mu_l, \sigma_l) =$
 134 $(\mu_l(x; w), \sigma_l^2(x; w))$ and $(\mu_l^*, \sigma_l^*) = (\mu_l(x^*; w), \sigma_l^2(x^*; w))$.

135 By solving the optimization problem in (2), high resolution images (e.g. ImageNet [4]) with a
 136 batch size of up to 40 can be constructed in [26]. However, f_w is only considered for three models:
 137 ImageNet pre-trained ResNet18 model, ImageNet pre-trained ResNet50 model, and MOCO v2 [2]
 138 pre-trained ResNet50 model fine-tuned with ImageNet. However, there are various choices of f_w .
 139 Although a broad spectrum of f_w choices is introduced in [6] (e.g., increasing channel size, models
 140 with or without skip connection), the authors of the work verified the effect from model variations on
 141 single image reconstruction as well as considered the optimization problem of the form (1) rather
 142 than (2). Thus, in this paper, we recap how model variations considered in [6] affect reconstruction
 143 of multiple images in a batch by solving optimization problem of the form (2) to achieve a better
 144 quality of reconstructed samples.

145 3 Re-evaluation of SOTA Gradient Attacks on a Broad Spectrum of Models

146 Prior works in gradient inversion attacks properly select limited range of models with vulnerability
 147 under the proposed attack methods to demonstrate their effectiveness [26, 12, 30, 6]. Therefore, this
 148 study aims to re-evaluate state-of-the-art attack methods on a broad spectrum of models. The target
 149 of our evaluation is attack methods that can solve the optimization problem of the form (2) assuming
 150 that the server as an honest-but-curious attacker desires to reconstruct multiple private images from
 151 batch gradients given, which is rarely studied previously. The model variations we considered are
 152 twofold: *implicit* and *explicit*.

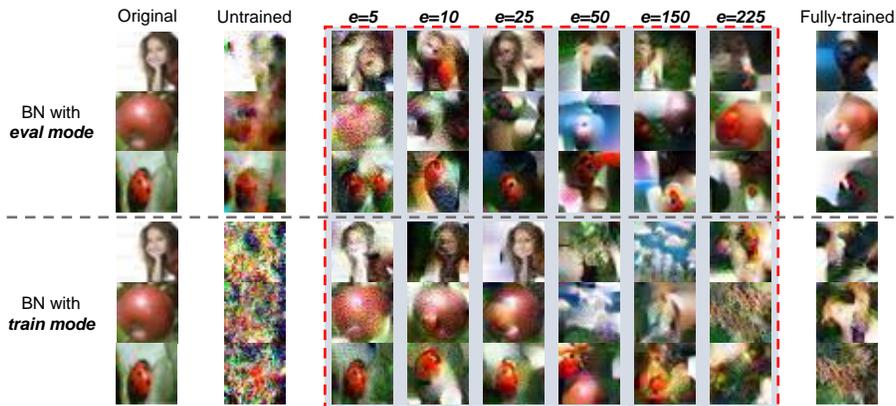


Figure 1: **Visualization of reconstructed images from implicit model variations of ResNet18.** Here e denotes training epochs. Then, “Untrained” means $e = 0$, and “Fully-trained” means $e = 300$ as the ResNet18 model is trained on CIFAR100 training set up to 300 epochs. Reconstructed images in the red dotted line box come from our choices of e . Original images (a woman image, an apple image, a beetle image) were randomly sampled from the CIFAR100 validation set.

153 3.1 Implicit model variation: BN modes and training epochs

154 While *explicit model variation* refers to an architectural change such as increasing channel sizes of
 155 the model, as suggested in [6], *implicit model variation* is invisible in the architectural level. However,
 156 changes arise internally within the same architecture such as applying different weights with different
 157 training epochs or switching the mode of normalization layers (e.g., switching between train and
 158 eval modes for BN). This is the first work to introduce the concept of implicit model variation. More

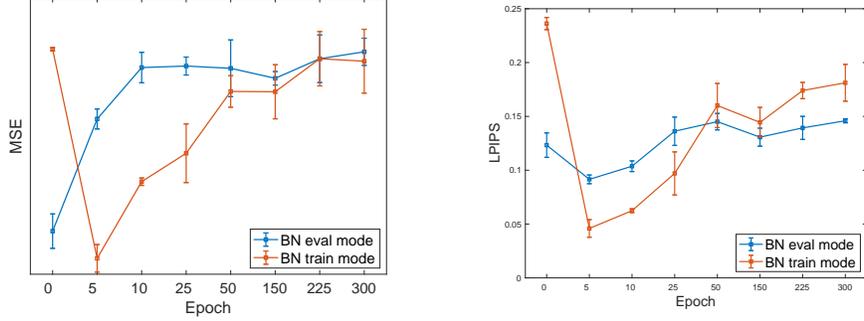


Figure 2: **Plotting the quality of reconstructed samples from implicit model variations of ResNet18 in terms of MSE (\downarrow , left) and LPIPS (\downarrow , right).**

159 specifically, this is the first time implicit model variation is considered for the evaluation of gradient
 160 inversion attacks. Interestingly, we experimentally found that the reconstruction quality ranges in a
 161 broad spectrum over *implicit model variations*.

162 3.1.1 BN modes: motivation

163 State-of-the-art gradient inversion attack methods elevate the quality of reconstructed samples by
 164 introducing batch statistics matching loss to the original problem of gradient matching as in (2).
 165 Therefore, we adopt a global model with BN layers to realize shared batch statistics in FL. BN
 166 layer has two modes of operation: train mode and eval mode [11]. However, recent works have
 167 not specified which mode is set for their demonstration while the malicious server, at least as an
 168 honest-but-curious attacker, can send a global model with BN layers set to any mode. Therefore,
 169 this study considers both BN train mode and BN eval mode for the re-evaluation of SOTA gradient
 170 attacks. Our re-evaluation results show that reconstruction results from different BN modes can be
 171 significantly different from each other even in terms of the same model weights as in Tab. 1.

Epoch (e)	MSE \downarrow	PSNR \uparrow	LPIPS \downarrow
0	0.8499 \pm 0.1996	12.8833 \pm 1.241	0.1233 \pm 0.0227
5	1.5033 \pm 0.1157	10.4366 \pm 0.43	0.0915 \pm 0.0081
10	1.7985 \pm 0.1766	9.87 \pm 0.2749	0.1037 \pm 0.0099
25	1.8072 \pm 0.1042	10.02 \pm 0.5716	0.1362 \pm 0.0263
50	1.7941 \pm 0.3291	9.8666 \pm 0.5507	0.1451 \pm 0.0153
150	1.7361 \pm 0.0783	10.34 \pm 0.3732	0.1307 \pm 0.0167
225	1.8495 \pm 0.2759	10.1866 \pm 0.4878	0.1393 \pm 0.0214
300	1.8899 \pm 0.1575	9.75 \pm 0.4313	0.1459 \pm 0.0038

(a) BN with eval mode

Epoch (e)	MSE \downarrow	PSNR \uparrow	LPIPS \downarrow
0	1.9045 \pm 0.0195	8.9265 \pm 0.0665	0.2362 \pm 0.0113
5	0.6921 \pm 0.1601	14.9733 \pm 0.9168	0.0459 \pm 0.0164
10	1.1367 \pm 0.0434	12.5 \pm 0.2861	0.0624 \pm 0.0035
25	1.3015 \pm 0.3402	11.6733 \pm 1.4027	0.097 \pm 0.04
50	1.66 \pm 0.1831	10.0433 \pm 0.6354	0.1601 \pm 0.041
150	1.6581 \pm 0.3138	10.2066 \pm 1.1074	0.1444 \pm 0.0279
225	1.8497 \pm 0.315	9.49 \pm 1.008	0.174 \pm 0.0151
300	1.8353 \pm 0.3703	9.4333 \pm 0.8832	0.1812 \pm 0.0342

(b) BN with train mode

Table 1: Quantitative comparison between reconstruction results for 50 CIFAR100 images from ResNet18 model with BN set to (a) eval mode and (b) train mode. MSE (\downarrow), PSNR (\uparrow), and LPIPS (\downarrow) are used as evaluation metrics. We highlight the best performance for each column in **bold**.

172 3.1.2 Training epochs: motivation

173 In a scenario of FL, a client can participate at any time during training. Then, a client can encounter
 174 the global model with arbitrary performance. This fact contradicts previous works' experimental
 175 setup, where the global model is chosen in a dichotomous manner: an untrained (or initialized) model
 176 or a model fully trained on the training set [6, 26]. Therefore, we re-evaluated SOTA inversion
 177 attacks on models with a broad spectrum of training epochs. We empirically found that the best
 178 reconstruction quality is usually obtained at earlier training epochs.

179 3.1.3 BN modes and training epochs: experimental results

180 **Setup** We trained a ResNet18 model on CIFAR100 [15] training set for 300 epochs using SGD
 181 optimizer with initial learning rate 0.1, momentum 0.9, and learning rate decay 0.1 applied when

182 $e = 150$ and $e = 225$ for the training epoch e . During training, we saved checkpoints of model
 183 weights when $e \in \{0, 5, 10, 25, 50, 150, 225, 300\}$ to consider the models from different training
 184 epochs. We oversampled model weights before the first learning decay ($0 < e < 150$) to cover the
 185 whole set of dynamically changing model weights in the beginning of training. On the other hand,
 186 hyperparameters and loss function choices for input reconstruction attacks are borrowed from [10].

187 **Results** As expected from their difference in batch statistics computation, BN with train mode
 188 and BN with eval mode show different reconstruction results both qualitatively (see Fig. 1.) and
 189 quantitatively (see Fig. 2 and Tab. 1.). When BN is set to eval mode, partial information (e.g. colors
 190 or shapes) is barely leaked in reconstructed images only for the cases $e = 0$ and $e = 5$ as described
 191 in Fig. 1 and Fig. 2. On the other hand, for BN with train mode, the quality of reconstructed images
 192 were sufficient enough to identify the object in each image only for the cases $e = 5, 10, 25$. Unlike
 193 the BN mode set to eval mode, it is remarkable that reconstructed images from BN with train mode
 194 in Fig. 1 are noisy images for $e = 0$. For the cases $e \geq 50$, input reconstruction failed for both BN
 195 modes and reconstructed images even from the same target gradients look significantly different
 196 for different BN modes. However, both BN with train mode and BN with eval mode have similar
 197 reconstruction quality in terms of both mean squared error (MSE) and Learned Perceptual Image
 198 Patch Similarity (LPIPS) [28] in Fig. 2 and Tab. 1. Therefore, *in the early stage of training, a global*
 199 *model would be privacy threatening with high probability.*

200 3.2 Explicit model variation: skip connection and channel size

201 Explicit model variations involve *change in architecture level* like removing skip connections in
 202 residual blocks or increasing the number of channels in convolution module, which are the kinds
 203 considered in previous works but on single image reconstruction. Therefore, we re-explore the effect
 204 of skip connection and channel size on the model’s vulnerability against gradient inversion attack
 205 but in the context of batch reconstruction. Skip connection helps information flow both forward
 206 and backward through the network, thus input reconstruction is expected to be easier for residual
 207 networks but harder for models without skip connection [9]. On the other hand, increasing channel
 208 size implies increasing dimension of gradients, which is the capacity of gradients to store information.
 209 Therefore, we expect that more information about input would be compressed in gradients when the
 210 number of channels increase.



Figure 3: **Visualization of reconstructed images from ConvNet on CIFAR100.** In each image block, the images at the positions of the red, pink, and green borders denote the original image, the reconstruction with BN (*eval mode*), and the reconstruction with BN (*train mode*), respectively. Original images were randomly sampled from CIFAR100 validation set.

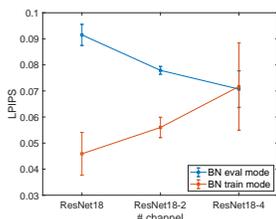


Figure 4: **Plotting the best reconstruction quality in terms of LPIPS (↓) among model variations through training epochs for ResNet18, ResNet18-2, ResNet18-4 models with BN *eval* (orange) and *train* (blue) modes.** $e = 5$ or $e = 10$ usually result in the best reconstruction quality. As channel size increases, the reconstruction quality increases for BN *eval* but decreases for BN *train*.

211 3.2.1 Skip connection and channel size: experimental results

212 **Setup** Instead of ResNet18, we trained a ConvNet model, ResNet18-2 model, and ResNet18-4
 213 model for explicit model variations. ConvNet, which is introduced in [6] for the first time, is a
 214 convolutional neural network without skip connection and ResNet18-2 and ResNet18-4 being ResNet
 215 with channel size doubled and quadrupled, respectively. Note that we apply implicit variations
 216 considered in the Sec. 3.1 to the models. Training conditions and hyperparameters for both model
 217 training and attack methods are kept the same with the setup in the previous section.

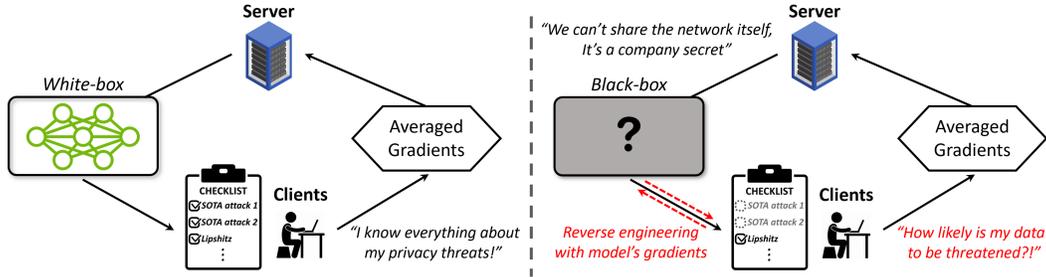
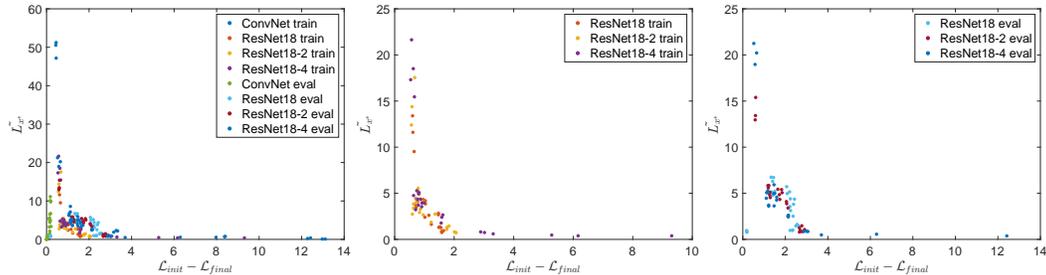


Figure 5: White-box (left) and black-box (right) FL settings.



(a) All models ($r_s = -0.78$) (b) ResNets, BN train ($r_s = -0.87$) (c) ResNets, BN eval ($r_s = -0.66$)

Figure 6: Proposed measure \tilde{L}_{x^*} is approximately a monotonic decreasing function with respect to $\mathcal{L}_{init} - \mathcal{L}_{final}$, the difference between initial (\mathcal{L}_{init}) and final losses (\mathcal{L}_{final}) among (a) all models, (b) ResNet models with BN train mode, and (c) ResNet models with BN eval mode considered in Sec. 3.

218 **Results** Reconstructed images from ConvNet models with the best quality, in terms of LPIPS,
 219 are listed in Fig. 3. For ConvNet models, reconstructed images, even with the best quality, are
 220 far from original images visually due to severe artifacts. Therefore, as expected from the role of
 221 skip connection in residual networks, a network without skip connection like ConvNet seems to be
 222 robust against input recovery attacks. Then, ConvNet models would be considered as global model
 223 candidates for privacy protection in FL despite of their worse performance than that of residual
 224 networks.

225 By contrast, the best averaged reconstruction results among the sampled training epochs $e \in$
 226 $\{0, 5, 10, 25, 50, 150, 225, 300\}$ are plotted in Fig. 4 for ResNet18, ResNet18-2, and ResNet18-
 227 4 models with varied BN modes. When BN is set to the eval mode, the reconstruction quality
 228 increases as the number of channels increases as expected. However, the reconstruction quality
 229 worsens as the number of channels increases for BN set to the train mode, which breaks the belief
 230 from previous works that increasing channel size makes input recovery attack easier [6, 30]. However,
 231 the reconstruction quality obtained with BN train mode is better than that with BN eval mode for all
 232 models considered except ResNet18-4, where their LPIPS range overlaps, implying that BN train
 233 mode is vulnerable against input recovery attacks than BN eval mode. The quantitative results for
 234 ConvNet, ResNet18-2, and ResNet18-4 are provided in Appendix A1.

235 4 Lipschitz Smoothness for Client-Side Privacy Leakage Detection

236 For privacy-preserving FL, choosing global model robust against any server-side input reconstruction
 237 attack method would be important. At the least, global model should be robust against well-known
 238 SOTA gradient inversion attack methods to alleviate clients' anxiety about any potential leakage from
 239 gradient sharing with the server. If clients can access the global model with the same level of a central
 240 server (*white-box*), applying SOTA attack methods directly to the global model with private images
 241 would be the best way for assessing whether or not the global model presents a risk to the client's
 242 privacy. However, in general, global model information would be opaque to clients due to company

243 secrets. As clients should communicate with the server via locally computed gradients, we suppose
 244 the *black-box* setting, where model gradients are given to clients as minimal information of the global
 245 model. Therefore, we provide a helpful measure for developing the system for clients to examine
 246 whether the given global model is safe in terms of privacy by using gradients computed with their
 247 self-controlled inputs. Note that *white-box* and *black-box* are described in Fig. 5.

248 4.1 Angular Lipschitz smoothness: motivation

249 If a function $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is Lipschitz smooth (or the derivative of f is Lipschitz continuous) with
 250 constant L , then the following holds: $\|\nabla f(x) - \nabla f(y)\| \leq L\|x - y\| \forall x, y \in \mathbb{R}^n$. The concept of
 251 Lipschitz smoothness or Lipschitz continuity is frequently employed to prove convergence theorem
 252 of gradient descent methods for optimization [24, 7, 27, 20, 18, 1]. This study employs the concept
 253 of Lipschitz smoothness to prove the following theorem in the context of gradient matching problem.

254 **Theorem 1** (Monotonic decreasing loss function). Suppose $\nabla_w \mathcal{L}(f(x), y)$ is Lipschitz continuous
 255 with respect to x with constant L and $\mathcal{L}_{grad}^x = \|\nabla_w \mathcal{L}(f(x), y) - g^*\|_2^2$ is given as a gradient matching
 256 loss. Then, when gradient descent Δx is applied with step size $\mu > 0$ and $L > \epsilon$ for some $\epsilon > 0$, the
 257 following holds:

$$\mathcal{L}_{grad}^{x+\Delta x} \leq \mathcal{L}_{grad}^x - \frac{1}{L^2} \left\| \frac{\partial \mathcal{L}_{grad}^x}{\partial x} \right\|_2^2. \quad (3)$$

258
 259 Inequality (3) implies that gradient matching loss strictly decreases as the gradient descent steps
 260 unless the gradient term $\frac{\partial \mathcal{L}_{grad}^x}{\partial x}$ is zero (i.e. gradient matching loss already converges). Furthermore,
 261 a gradient descent with a small L (or large $\frac{1}{L^2}$) can accelerate the convergence of gradient matching
 262 optimization but with the premise that $L > \epsilon$ for $\epsilon > 0$. This premise is required to ensure the
 263 first-order Taylor approximation for $\nabla_w \mathcal{L}(f(x + \Delta x), y)$ in the proof in Appendix A2. Therefore, in
 264 a particular range of L (i.e., $L > \epsilon$), we hypothesize that a global model with smaller L experiences a
 265 sharper loss drop in gradient matching optimization. *We empirically found that L is not too small for*
 266 *most models, thus meeting the premise in reality.*

267 For the empirical verification of the hypothesis in the context of input reconstruction,
 268 we desire to compute Lipschitz smoothness constant locally around x^* , $L_{x^*, \epsilon} =$
 269 $\sup_{\|x-x^*\| < \epsilon, x \neq x^*} \frac{\|\nabla_w \mathcal{L}(f(x), y) - \nabla_w \mathcal{L}(f(x^*), y)\|}{\|x-x^*\|}$, with small ϵ , for the models considered in Sec-
 270 tions 3.1 and 3.2. Recent works on computing precise upper bound of L only focus on multi-layer
 271 perceptrons (MLP) due to the difficulty of computing L for normalization layers or residual layers.
 272 Therefore, $L_{x^*, \epsilon}$ is estimated as $\tilde{L}_{x^*} = \max_{n \neq 0} \frac{\|\nabla_w \mathcal{L}(f(x^*+n), y) - \nabla_w \mathcal{L}(f(x^*), y)\|}{\|n\|}$ by sampling 1,000
 273 noises (n) from the Gaussian distribution $\mathcal{N}(0, 0.001^2)$ in our experiments.

274 However, \tilde{L}_{x^*} can be any nonnegative value by scaling loss function \mathcal{L} . If \mathcal{L} is scaled by nonnegative
 275 scalar k , then \tilde{L}_{x^*} is scaled by k too, allowing \tilde{L}_{x^*} to be manipulated by the server using simple
 276 loss scaling. Therefore, inspired by the cosine similarity loss function, which is scale-invariant,
 277 we propose *the angular Lipschitz constant* $L_{x^*}^{cos} = \max_{n \neq 0} \frac{1 - cs(\nabla_w \mathcal{L}(f(x^*+n), y), \nabla_w \mathcal{L}(f(x^*), y))}{1 - cs(x^*, x^*+n)}$ (cs
 278 being the cosine similarity loss) as a loss scaling-invariant alternative to \tilde{L}_{x^*} . We find that $L_{x^*}^{cos}$ shows
 279 a strong monotonic correlation with the quality of reconstructed samples, demonstrating the potential
 280 of $L_{x^*}^{cos}$ to be imperative for client-side defense methods.

281 4.2 Angular Lipschitz smoothness: experimental results

282 We computed \tilde{L}_{x^*} and the attacker’s loss drop $\mathcal{L}_{init} - \mathcal{L}_{final}$ (\mathcal{L}_{init} and \mathcal{L}_{final} being the initial and
 283 final losses, respectively) for the models and input batches considered in Sec. 3 (Fig. 6a). We also
 284 quantified their correlation using the Spearman’s rank correlation coefficient r_s , which quantifies
 285 how two variables are in a monotonic relationship. $r_s = 1$ ($r_s = -1$) means that one variable is
 286 a completely monotonic increasing (decreasing) function with respect to the other one. Then, \tilde{L}_{x^*}
 287 is almost a monotonic decreasing function with respect to $\mathcal{L}_{init} - \mathcal{L}_{final}$ with $r_s = -0.78$, thus
 288 validating our hypothesis. For ResNet models with BN *train* (Fig. 6b), \tilde{L}_{x^*} and $\mathcal{L}_{init} - \mathcal{L}_{final}$
 289 show a stronger monotonic correlation than that for ResNet models with BN *eval* (Fig. 6c) with

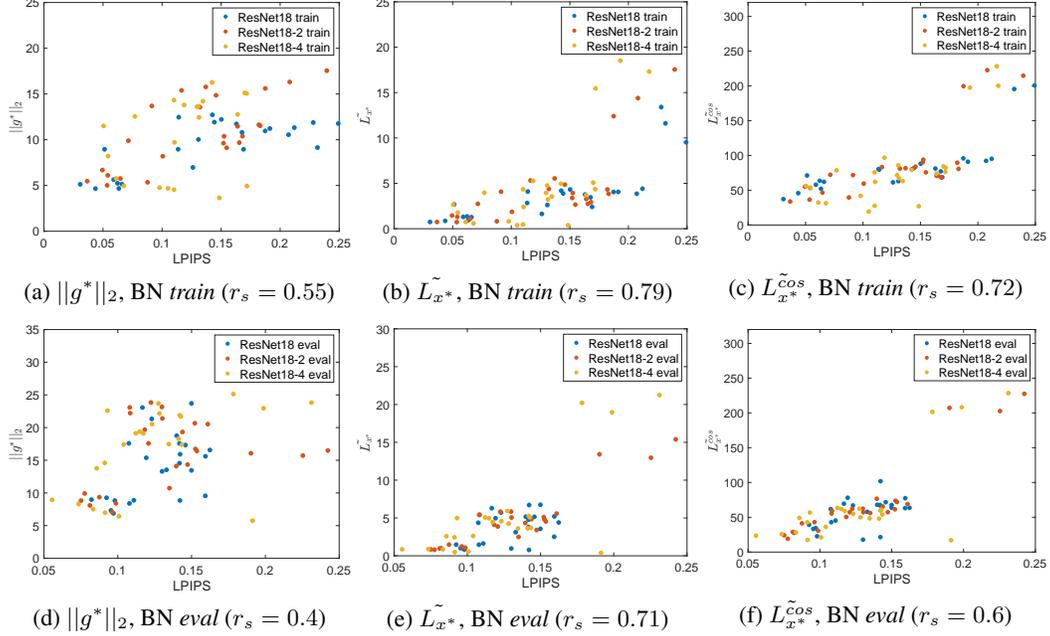


Figure 7: Comparison of $\|g^*\|_2$, L_{x^*} , and $L_{x^*}^{cos}$ in terms of the correlation between LPIPS of reconstructed samples for ResNet models with BN *train* (top) and BN *eval* (bottom)

290 $r_s = -0.85$. As in Tab. 1 and Fig. 1, reconstructed samples are closer to their original images in BN
 291 *train mode*, thus L_{x^*} , which is computed around the ground truth x^* , seems to fit more to BN *train*
 292 while L should be estimated around the solution from the attack method rather than x^* for the case of
 293 BN *eval*. However, clients cannot access to the solution from the the attack method in the *black-box*
 294 setting. The plot of L_{x^*} and $\mathcal{L}_{init} - \mathcal{L}_{final}$ for the ConvNet models is provided in Appendix A3.

295 5 Limitations and Future Work

296 Our hypothesis can be extended to the correlation between Lipschitz constant and the quality of
 297 reconstructed samples, rather than loss drop. Zero gradient matching loss does not mean complete
 298 recovery of original images due to the existence of *twin data* [30], two different data input with
 299 identical model gradients. However, we empirically found that both L_{x^*} and $L_{x^*}^{cos}$ show positive
 300 monotonic correlations with the quality of reconstructed samples, in terms of LPIPS (lower value is
 301 better) (Fig. 7). In particular, they beat the baseline measure, the norm of given gradients ($\|g^*\|_2$),
 302 which was implicitly believed to be the amount of information within the gradients in previous works,
 303 by a wide margin, in terms of r_s . Therefore, we expect L_{x^*} and $L_{x^*}^{cos}$ to be the key factors for
 304 developing future client-side defense strategies.

305 6 Conclusions

306 Here, we re-evaluated the SOTA attack method on a broad spectrum of models in the context of
 307 batch reconstruction, which is rarely studied in previous works. We considered model variations
 308 of two types: *implicit*, which changes in model weights or BN modes within the same architecture,
 309 and *explicit*, with changes in architecture. The re-evaluation results indicate that the quality of the
 310 reconstruction attack varies depending on the implicit or explicit model changes. Therefore, inspired
 311 by our theorem related to the convergence of gradient matching optimization and scale-invariance
 312 of the cosine similarity loss function, we propose an explainable and predictive measure for privacy
 313 leakage, an angular Lipschitz constant L^{cos} , which is invariant to trivial loss scaling attacks from
 314 malicious servers. We empirically find that L^{cos} shows a strong monotonic correlation with the
 315 quality of reconstructed samples, thus expecting the potential of L^{cos} to be a key factor for clients'
 316 defense strategies in a *black-box* setting, where only model gradients are given as minimal information
 317 about the global model to clients.

318 **References**

- 319 [1] H. H. Bauschke, J. Bolte, and M. Teboulle. A descent lemma beyond lipschitz gradient continuity:
320 first-order methods revisited and applications. *Mathematics of Operations Research*, 42(2):330–348, 2017.
- 321 [2] X. Chen, H. Fan, R. Girshick, and K. He. Improved baselines with momentum contrastive learning. *arXiv*
322 *preprint arXiv:2003.04297*, 2020.
- 323 [3] T. Dang, O. Thakkar, S. Ramaswamy, R. Mathews, P. Chin, and F. Beaufays. Revealing and protecting
324 labels in distributed training. *Advances in Neural Information Processing Systems*, 34, 2021.
- 325 [4] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. Imagenet: A large-scale hierarchical image
326 database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee,
327 2009.
- 328 [5] J. Duchi, E. Hazan, and Y. Singer. Adaptive subgradient methods for online learning and stochastic
329 optimization. *Journal of machine learning research*, 12(7), 2011.
- 330 [6] J. Geiping, H. Bauermeister, H. Dröge, and M. Moeller. Inverting gradients-how easy is it to break privacy
331 in federated learning? *Advances in Neural Information Processing Systems*, 33:16937–16947, 2020.
- 332 [7] H. Gouk, E. Frank, B. Pfahringer, and M. J. Cree. Regularisation of neural networks by enforcing lipschitz
333 continuity. *Machine Learning*, 110(2):393–416, 2021.
- 334 [8] A. Hatamizadeh, H. Yin, P. Molchanov, A. Myronenko, W. Li, P. Dogra, A. Feng, M. G. Flores,
335 J. Kautz, D. Xu, et al. Do gradient inversion attacks make federated learning unsafe? *arXiv preprint*
336 *arXiv:2202.06924*, 2022.
- 337 [9] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *Proceedings of the*
338 *IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- 339 [10] Y. Huang, S. Gupta, Z. Song, K. Li, and S. Arora. Evaluating gradient inversion attacks and defenses in
340 federated learning. *Advances in Neural Information Processing Systems*, 34, 2021.
- 341 [11] S. Ioffe and C. Szegedy. Batch normalization: Accelerating deep network training by reducing internal
342 covariate shift. In *International conference on machine learning*, pages 448–456. PMLR, 2015.
- 343 [12] J. Jeon, K. Lee, S. Oh, J. Ok, et al. Gradient inversion with generative image prior. *Advances in Neural*
344 *Information Processing Systems*, 34:29898–29908, 2021.
- 345 [13] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. In *ICLR (Poster)*, 2015.
- 346 [14] J. Konečný, H. B. McMahan, F. X. Yu, P. Richtárik, A. T. Suresh, and D. Bacon. Federated learning:
347 Strategies for improving communication efficiency. *arXiv preprint arXiv:1610.05492*, 2016.
- 348 [15] A. Krizhevsky. Learning multiple layers of features from tiny images. 2009.
- 349 [16] O. Li, J. Sun, X. Yang, W. Gao, H. Zhang, J. Xie, V. Smith, and C. Wang. Label leakage and protection in
350 two-party split learning. 2022.
- 351 [17] Q. Li, Y. Diao, Q. Chen, and B. He. Federated learning on non-iid data silos: An experimental study. 2021.
- 352 [18] X. Li and F. Orabona. On the convergence of stochastic gradient descent with adaptive stepsizes. In *The*
353 *22nd International Conference on Artificial Intelligence and Statistics*, pages 983–992. PMLR, 2019.
- 354 [19] X. Li, M. JIANG, X. Zhang, M. Kamp, and Q. Dou. Fedbn: Federated learning on non-iid features via
355 local batch normalization. In *International Conference on Learning Representations*, 2020.
- 356 [20] V. V. Mai and M. Johansson. Stability and convergence of stochastic gradient clipping: Beyond lipschitz
357 continuity and smoothness. In *International Conference on Machine Learning*, pages 7325–7335. PMLR,
358 2021.
- 359 [21] B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. y Arcas. Communication-efficient learning of
360 deep networks from decentralized data. In *Artificial intelligence and statistics*, pages 1273–1282. PMLR,
361 2017.
- 362 [22] A. Paszke, S. Gross, S. Chintala, G. Chanan, E. Yang, Z. DeVito, Z. Lin, A. Desmaison, L. Antiga, and
363 A. Lerer. Automatic differentiation in pytorch. 2017.

- 364 [23] S. Ruder. An overview of gradient descent optimization algorithms. *arXiv preprint arXiv:1609.04747*,
365 2016.
- 366 [24] S. Santurkar, D. Tsipras, A. Ilyas, and A. Madry. How does batch normalization help optimization?
367 *Advances in neural information processing systems*, 31, 2018.
- 368 [25] D. Ye, T. Zhu, S. Zhou, B. Liu, and W. Zhou. Label-only model inversion attack: The attack that requires
369 the least information. *arXiv preprint arXiv:2203.06555*, 2022.
- 370 [26] H. Yin, A. Mallya, A. Vahdat, J. M. Alvarez, J. Kautz, and P. Molchanov. See through gradients: Image
371 batch recovery via gradinversion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and
372 Pattern Recognition*, pages 16337–16346, 2021.
- 373 [27] J. Zeng, T. T.-K. Lau, S. Lin, and Y. Yao. Global convergence of block coordinate descent in deep learning.
374 In *International conference on machine learning*, pages 7313–7323. PMLR, 2019.
- 375 [28] R. Zhang, P. Isola, A. A. Efros, E. Shechtman, and O. Wang. The unreasonable effectiveness of deep
376 features as a perceptual metric. In *Proceedings of the IEEE conference on computer vision and pattern
377 recognition*, pages 586–595, 2018.
- 378 [29] B. Zhao, K. R. Mopuri, and H. Bilen. idlg: Improved deep leakage from gradients. *arXiv preprint
379 arXiv:2001.02610*, 2020.
- 380 [30] J. Zhu and M. B. Blaschko. R-gap: Recursive gradient attack on privacy. In *International Conference on
381 Learning Representations*, 2020.
- 382 [31] L. Zhu, Z. Liu, and S. Han. Deep leakage from gradients. *Advances in Neural Information Processing
383 Systems*, 32, 2019.

384 Checklist

385 The checklist follows the references. Please read the checklist guidelines carefully for information on
386 how to answer these questions. For each question, change the default **[TODO]** to **[Yes]**, **[No]**, or
387 **[N/A]**. You are strongly encouraged to include a **justification to your answer**, either by referencing
388 the appropriate section of your paper or providing a brief inline description. For example:

- 389 • Did you include the license to the code and datasets? **[Yes]** See Section ??.
- 390 • Did you include the license to the code and datasets? **[No]** The code and the data are
391 proprietary.
- 392 • Did you include the license to the code and datasets? **[N/A]**

393 Please do not modify the questions and only use the provided macros for your answers. Note that the
394 Checklist section does not count towards the page limit. In your paper, please delete this instructions
395 block and only keep the Checklist section heading above along with the questions/answers below.

396 1. For all authors...

- 397 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
398 contributions and scope? **[Yes]** We thoroughly listed the contributions of our work in
399 both abstract and introduction.
- 400 (b) Did you describe the limitations of your work? **[Yes]** See Sec. 5 for limitations of our
401 work and future research direction.
- 402 (c) Did you discuss any potential negative societal impacts of your work? **[N/A]**
- 403 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
404 them? **[Yes]**

405 2. If you are including theoretical results...

- 406 (a) Did you state the full set of assumptions of all theoretical results? **[Yes]** In Theorem 1,
407 all the assumptions are included in the statement. The details are described in Appendix
408 A2.

- 409 (b) Did you include complete proofs of all theoretical results? [Yes] We include complete
410 proofs in Appendix A2.
- 411 3. If you ran experiments...
- 412 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
413 mental results (either in the supplemental material or as a URL)? [Yes] The details for
414 experiments are partially described in Results section in Sec. 3 including Appendix.
- 415 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
416 were chosen)? [Yes] All the training details for experiments are described in Results
417 section in Sec. 3.
- 418 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
419 ments multiple times)? [Yes] We report standard deviation of our results in Tab. 1 and
420 Fig. 2.
- 421 (d) Did you include the total amount of compute and the type of resources used (e.g., type
422 of GPUs, internal cluster, or cloud provider)? [Yes] We report them in Appendix A4.
- 423 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 424 (a) If your work uses existing assets, did you cite the creators? [Yes] We cited the creators
425 for CIFAR100, ImageNet, ConvNet, and ResNet18.
- 426 (b) Did you mention the license of the assets? [N/A]
- 427 (c) Did you include any new assets either in the supplemental material or as a URL? [No]
- 428 (d) Did you discuss whether and how consent was obtained from people whose data you're
429 using/curating? [Yes] We use data widely used in the field of research related to our
430 work
- 431 (e) Did you discuss whether the data you are using/curating contains personally identifiable
432 information or offensive content? [N/A]
- 433 5. If you used crowdsourcing or conducted research with human subjects...
- 434 (a) Did you include the full text of instructions given to participants and screenshots, if
435 applicable? [N/A]
- 436 (b) Did you describe any potential participant risks, with links to Institutional Review
437 Board (IRB) approvals, if applicable? [N/A]
- 438 (c) Did you include the estimated hourly wage paid to participants and the total amount
439 spent on participant compensation? [N/A]