

Table 1: Hyperparameter settings for each experiment. The set of examples that make up the prior ( $\pi$ ), including the target ( $z^*$ ), and the other examples in the training set ( $D_-$ ) are always drawn from the same data distribution, except for the experiment performed in Figure 9.

| Dataset  | Experiment                       | Clipping norm $C$ | Sampling probability $q$ | Update steps $T$ | Model architecture $\theta$                            | Training dataset size ( $ D_-  + 1$ ) | Prior size      |
|----------|----------------------------------|-------------------|--------------------------|------------------|--|---------------------------------------|-----------------|
| CIFAR-10 | Figure 1b, Figure 6b, Figure 12b | 1                 | 1                        | 100              | WRN-28-10  | 500                                   | -               |
|          | Figure 2f                        | 1                 | 1                        | 100              | WRN-28-10  | 500                                   | 10              |
|          | Figure 1a, Figure 6a, Figure 12a | 0.1               | 1                        | 100              | MLP (784 $\rightarrow$ 10 $\rightarrow$ 10)            | 1,000                                 | -               |
| MNIST    | Figure 3f                        | 0.1               | 1                        | 100              | MLP (784 $\rightarrow$ 10 $\rightarrow$ 10)            | 1,000                                 | $2^1 - 2^{11}$  |
|          | Figure 4f                        | 0.1               | 0.02                     | 1,000            | MLP (784 $\rightarrow$ 10 $\rightarrow$ 10)            | 500                                   | 10              |
|          | Figure 5f                        | 1                 | 0.01-0.99                | 100              | MLP (784 $\rightarrow$ 10 $\rightarrow$ 10)            | 1,000                                 | 10              |
|          | Figure 7f                        | 0.1               | 1                        | 100              | MLP (784 $\rightarrow$ 10, 100, 1000 $\rightarrow$ 10) | 1,000                                 | 10              |
|          | Figure 8f                        | 0.1               | 1                        | 100              | MLP (784 $\rightarrow$ 10 $\rightarrow$ 10)            | 1,000                                 | $2^1, 2^3, 2^7$ |
|          | Figure 9f                        | 0.1               | 1                        | 100              | MLP (784 $\rightarrow$ 10 $\rightarrow$ 10)            | 1,000                                 | 10              |
|          | Figure 10f                       | 0.1               | 1                        | 100              | MLP (784 $\rightarrow$ 10 $\rightarrow$ 10)            | 5, 129, 1,000                         | 10              |
|          | Figure 11f                       | 1.0               | 1                        | -                | -  | -                                     | 10              |
|          | Figure 13f                       | 0.1, 1            | 1                        | 100              | MLP (784 $\rightarrow$ 10 $\rightarrow$ 10)            | 1,000                                 | 10              |
|          | Figure 14f                       | 1                 | 0.01-0.99                | 100              | MLP (784 $\rightarrow$ 10 $\rightarrow$ 10)            | 1,000                                 | 2, 10, 100      |
|          | Figure 15f                       | 1                 | 0.01-0.99                | 100              | MLP (784 $\rightarrow$ 10 $\rightarrow$ 10)            | 1,000                                 | 2, 10, 100      |

## 518 A Experimental details

519 We detail the experimental settings used throughout the paper, and specific hyperparameters used for  
 520 the various attacks we investigate. The exact configurations for each experiment are given in Table 1.  
 521 We vary many experimental hyperparameters to investigate their effect on reconstruction, however,  
 522 the default setting is described next.

523 For MNIST experiments we use a two layer MLP with hidden width 10 and eLU activations. The  
 524 attacks we design in this work perform equally well on all common activation functions, however it  
 525 is well known that the model-based attack (Balle et al., 2022) performs poorly on piece-wise linear  
 526 activations like ReLU. We set  $|D_-| = 999$  (and so the training set size is  $|D_- \cup \{z^*\}| = 1,000$ ) and  
 527 train with full-batch DP-SGD for  $T = 100$  steps. For each  $\epsilon$ , we select the learning rate by sweeping  
 528 over a range of values between 0.001 and 100; we do not use any momentum in optimization. We set  
 529  $C = 0.1$ ,  $\delta = 10^{-5}$  and adjust the noise scale  $\sigma$  for a given target  $\epsilon$ . The accuracy of this model is  
 530 over 90% for  $\forall \epsilon \geq 10$ , however we emphasize that our experiments on MNIST are meant to primarily  
 531 investigate the tightness of our reconstruction upper bounds. We set the size of the prior  $\pi$  to ten,  
 532 meaning the baseline probability of successful reconstruction is 10%.

533 For the CIFAR-10 dataset, we use a Wide-ResNet (Zagoruyko & Komodakis, 2016) model with  
 534 28 layers and width factor 10 (denoted as WRN-28-10), group normalization, and eLU activations.  
 535 We align with the set-up of De et al. (2022), who fine-tune a WRN-28-10 model from ImageNet to  
 536 CIFAR-10. However, because the model-based attack is highly expensive, we only fine-tune the final  
 537 layer. We set  $|D_-| = 499$  (and so the training set size is  $|D_- \cup \{z^*\}| = 500$ ) and train with full-batch  
 538 DP-SGD for  $T = 100$  steps; again we sweep over the choice of learning rate for each value of  $\epsilon$ . We  
 539 set  $C = 1$ ,  $\delta = 10^{-5}$  and adjust the noise scale  $\sigma$  for a given target  $\epsilon$ . The accuracy of this model is  
 540 over 89% for  $\forall \epsilon \geq 10$ , which is close to the state-of-the-art results given by De et al. (2022), who  
 541 achieve 94.2% with the same fine-tuning setting at  $\epsilon = 8$  (with a substantially larger training set  
 542 size). Again, we set the size of the prior  $\pi$  to ten, meaning the baseline probability of successful  
 543 reconstruction is 10%.

544 For the gradient-based and model-based attack we generate 1,000 reconstructions and for prior-aware  
 545 attack experiments we generate 10,000 reconstructions from which we estimate a lower bound for  
 546 probability of successful reconstruction. That is, for experiments in Section 2 repeat the attack 1,000  
 547 times for targets randomly sampled from base dataset (MNIST or CIFAR-10), and for all other  
 548 experiments we repeat the attack 10,000 times for targets randomly sampled from the prior, which  
 549 is itself sampled from the base dataset (MNIST or CIFAR-10). We now give experimental details  
 550 specific to the various attacks used throughout the paper. Note that for attack results, we report 95%  
 551 confidence intervals around our lower bound estimate, however, in many cases these intervals are so  
 552 tight it renders them invisible to the eye.

553 **Model-based attack details.** For the model-based attack given by Balle et al. (2022), we train  
 554 40K shadow models, and as stated above, construct a test set by training a further 1,000 models  
 555 on 1,000 different targets (and  $D_-$ ) from which we evaluate our reconstructions. We use the same  
 556 architecture for the ReCoNN network and optimization hyperparameters as described in the MNIST  
 557 and CIFAR-10 experiments in Balle et al. (2022), and refer the interested reader there for details.

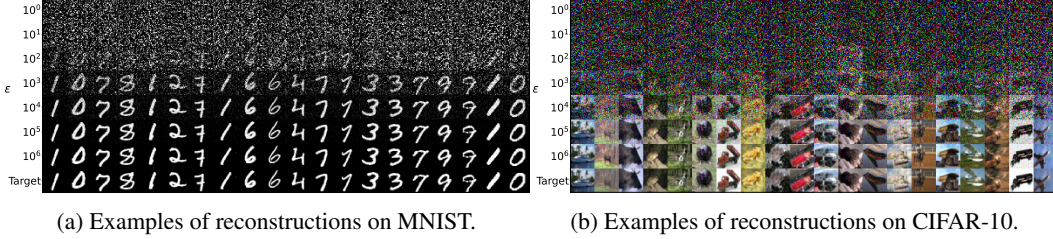


Figure 6: We give qualitative examples of reconstructions in Figure 6a and Figure 6b for the gradient-based reconstruction attack described in Section 2

558 **Gradient-based attack details.** Our optimization hyperparameters are the same for both MNIST  
 559 and CIFAR-10. We initialize a  $\hat{z}$  from uniform noise and optimize it with respect to the loss given in  
 560 Equation (1) for 1M steps of gradient descent with a learning rate of 0.01. We found that the loss  
 561 occasionally diverges and it is useful to have random restarts of the optimization process; we set the  
 562 number of random restarts to five. Note we assume that the label of  $z^*$  is known to the adversary.  
 563 This is a standard assumption in training data reconstruction attacks on federated learning, as Zhao  
 564 et al. (2020) demonstrated the label of the target can be inferred given access to gradients. If we did  
 565 not make this assumption, we can run the attack by exhaustively searching over all possible labels.  
 566 For the datasets we consider, this would increase the cost of the attack by a factor of ten. We evaluate  
 567 the attack using the same 1,000 targets used to evaluate the model-based attack.

568 **Prior-aware attack details.** The prior-aware attacks given in Algorithm 2 (and in Algorithm 3)  
 569 have no specific hyper-parameters that need to be set. As stated, the attack proceeds by summing the  
 570 inner-product defined in Section 3.3 over all training steps for each sample in the prior and selecting  
 571 the sample that maximizes this sum as the reconstruction. One practical note is that we found it  
 572 useful to normalize privatized gradients such that the privatized gradient containing the target will be  
 573 sampled from a Gaussian with unit mean instead of  $C^2$ , which will be sensitive to choice of  $C$  and  
 574 can lead to numerical precision issues.

575 **Estimating  $\gamma$  details.** As described in Section 3,  $\nu$  is instantiated as  $\mathcal{N}(0, \sigma^2 I)$ , a  $T$ -dimensional  
 576 isotropic Gaussian distribution with zero mean, and  $\mu$  is given by  $\sum_{w \in \{0,1\}^T} p(w) \mathcal{N}(w, \sigma^2 I)$ , a  
 577 mixture of  $T$ -dimensional isotropic Gaussian distributions with means in  $\{0, 1\}^T$  sampled according  
 578 to  $B(q, T)$ . Throughout all experiments, we use 1M independent Gaussian samples to compute the  
 579 estimation of  $\gamma$  given by the procedure in Algorithm 1, and because we use a discrete prior of size  
 580  $|\pi|$ , the base probability of reconstruction success,  $\kappa$ , is given as  $1/|\pi|$ .

## 581 B Visualization of reconstruction attacks on MNIST and CIFAR-10

582 In Figure 6, we give a selection of examples for the gradient-based reconstruction attack presented in  
 583 Section 2 and plotted in Figure 1

## 584 C Does the model size make a difference to the prior-aware attack?

585 Our results on MNIST and CIFAR-10 suggest that the model size does not impact the tightness of  
 586 our reconstruction attack (lower bound on probability of success); the MLP model used for MNIST  
 587 has 7,960 trainable parameters, while the WRN-28-10 model used for CIFAR-10 has 36.5M. We  
 588 systematically evaluate the impact of the model size on our prior-aware attack by increasing the  
 589 size of the MLP hidden layer by factors of ten, creating models with 7,960, 79,600, and 796,000  
 590 parameters. Results are given in Figure 7, where we observe almost no difference in terms of attack  
 591 success between the different model sizes.

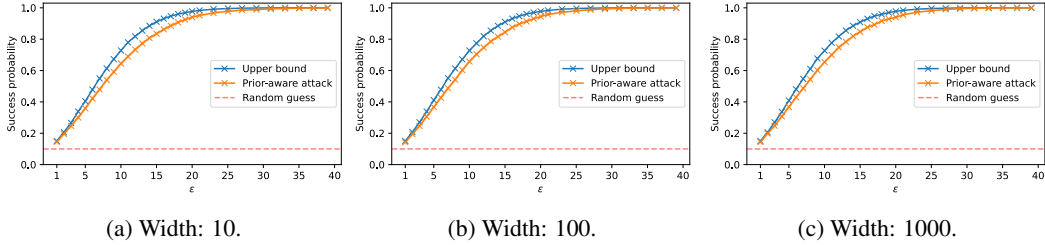


Figure 7: Comparison of model sizes on reconstruction by varying the hidden layer width in a two layer MLP.

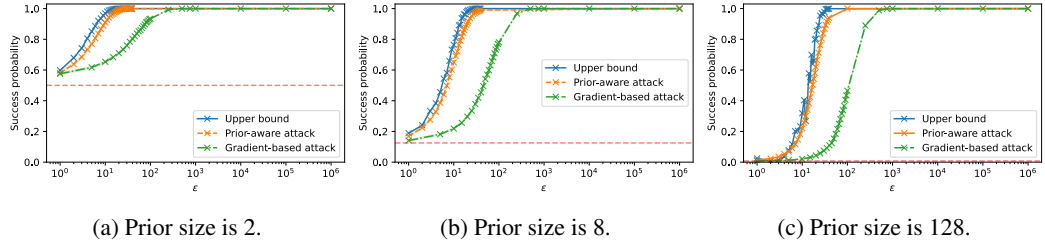


Figure 8: Comparison of prior-aware and gradient-based attack for different prior sizes.

## 592 D Comparing the gradient-based attack with the prior-aware attack

593 Our experiments have mainly been concerned with measuring how DP affects an adversary’s ability to  
 594 infer which point was included in training, given that they have access to all possible points that could  
 595 have been included, in the form of a discrete prior. This experimental set-up departs from Figure 1,  
 596 where we assumed the adversary does not have access to a prior set, and so cannot run the prior-aware  
 597 attack as described in Algorithm 2. Following on from results in Section 3.4, we transform these  
 598 gradient-based attack experimental findings into a probability of successful reconstruction by running  
 599 a post-processing conversion, allowing us to measure how the assumption of adversarial access to the  
 600 discrete prior affects reconstruction success. We run the post-processing conversion in the following  
 601 way: Given a target sample  $z^*$  and a reconstruction  $\hat{z}$  found through optimizing the gradient based  
 602 loss in Equation (1), we construct a prior consisting of  $z^*$  and  $n - 1$  randomly selected points from  
 603 the MNIST dataset, where  $n = 10$ . We then measure the  $L_2$  distance between  $\hat{z}$  and every point in  
 604 this constructed prior, and assign reconstruction a success if the smallest distance is with respect to  
 605  $z^*$ . For each target  $z^*$ , we repeat this procedure 1,000 times, with different random selections of size  
 606  $n - 1$ , and overall report the average reconstruction success over 1,000 different targets.

607 This allows us to compare the gradient-based attack (which is prior “unaware”) directly to our  
 608 prior-aware attack. Results are shown in Figure 8, where we vary the size of the prior between 2,  
 609 8, and 128. In all cases, we see an order of magnitude difference between the gradient-based and  
 610 prior-aware attack in terms of reconstruction success. This suggests that if we assume the adversary  
 611 does not have prior knowledge of the possible set of target points, the minimum value of  $\epsilon$  necessary  
 612 to protect against reconstruction attacks increases.

## 613 E Effects of the threat model and prior distribution on reconstruction

614 The ability to reconstruct a training data point will naturally depend on the threat model in which  
 615 the security game is instantiated. So far, we have limited our investigation to align with the standard  
 616 adversary assumptions in the DP threat model. We have also limited ourselves to a setting where the  
 617 prior is sampled from the same base distribution as  $D$ . These choices will change the performance  
 618 of our attack, which is what we measure next.

619 **Prior type.** We measure how the choice of prior affects reconstruction in Figure 9. We train models  
 620 when the prior is from the same distribution as the rest of the training set (MNIST), and when the  
 621 prior is sampled random noise. Note, because the target point  $z^*$  is included in the prior, this means

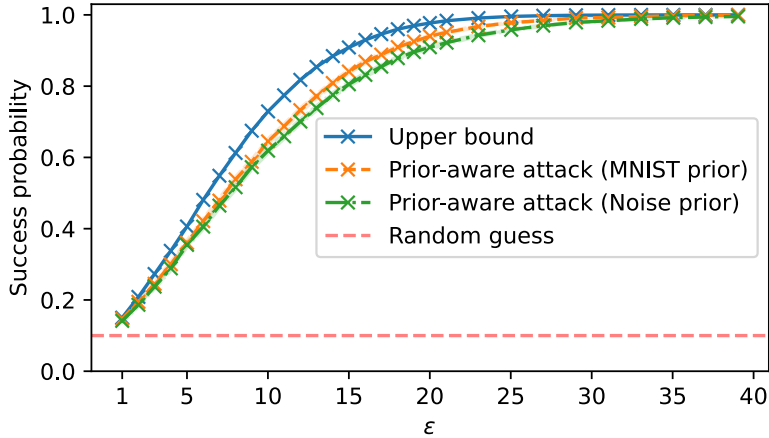


Figure 9: Comparison of how the choice of prior,  $\pi$ , affects reconstruction success. The prior is selected from a set of examples sampled from MNIST or uniform noise (that has the same intra-sample distance statistics as the MNIST prior).

622 we measure how reconstruction success changes when we change the distribution the target was  
 623 sampled from. One may expect that the choice of prior to make a difference to reconstruction success  
 624 if the attack relies on distinguishability between  $D$  and  $z^*$  with respect to some function operating  
 625 on points and model parameters (e.g. the difference in loss between points in  $D$  and  $z^*$ ). However,  
 626 we see that there is little difference between the two; both are close to the upper bound.

627 On reflection, this is expected as our objective is simply the sum of samples from a Gaussian, and  
 628 so the choice of prior may impact our probability of correct inference if this choice affects the  
 629 probability that a point will be clipped, or if points in the prior have correlated gradients. We explore  
 630 how different values of clipping,  $C$ , can change reconstruction success probability in Appendix [L](#).

631 **Knowledge of batch gradients.** The DP threat model assumes the adversary has knowledge of the  
 632 gradients of all samples other than the target  $z^*$ . Here, we measure how important this assumption  
 633 is to our attack. We compare the prior-aware attack (which maximizes  $\sum_{t=1}^T \langle \text{clip}_C(\nabla_{\theta_t} \ell(z_i)), \bar{g}_t \rangle$ )  
 634 against the attack that selects the  $z_i$  maximizing  $\sum_{t=1}^T \langle \text{clip}_C(\nabla_{\theta_t} \ell(z_i)), g_t \rangle$ , where the adversary  
 635 does not subtract the known gradients from the objective.

636 In Figure [10](#) we compare, in a full-batch setting, when  $|D_-|$  is small (set to 4), and see the attack does  
 637 perform worse when we do not deduct known gradients. However, the effect is more pronounced as  
 638  $|D_-|$  becomes larger, the attack completely fails when setting it to 128. This is somewhat expected,  
 639 as with a larger number of samples in a batch it is highly likely there are gradients correlated with  
 640 the  $z^*$  target gradient, masking out its individual contribution and introducing noise into the attack  
 641 objective.

## 642 F Improved prior-aware attack algorithm

643 As explained in Section [4](#), the prior-aware attack in Algorithm [2](#) does not account for the variance  
 644 introduced into the attack objective in mini-batch DP-SGD, and so we design a more efficient attack  
 645 specifically for the mini-batch setting. We give the pseudo-code for this improved prior-aware attack  
 646 in Algorithm [3](#).

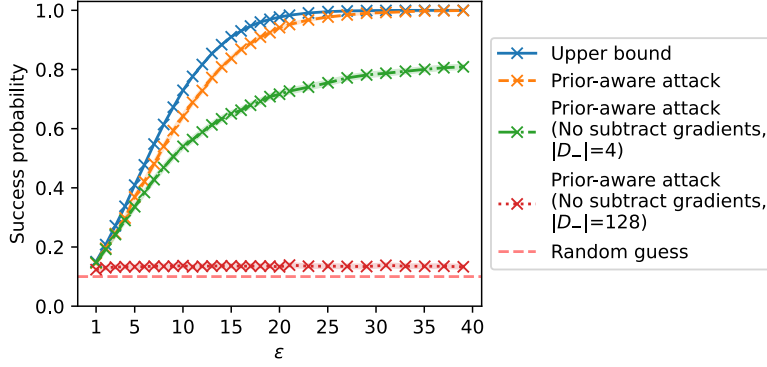


Figure 10: In line with the DP threat model, our attack in Algorithm 2 assumes the adversary can subtract known gradients from the privatized gradient. We measure what effect removing this assumption has on reconstruction success probability. When the size of the training set is small, removing this assumption has a minor effect, while reconstruction success drops to random with a larger training set size.

---

**Algorithm 3** Improved prior-aware attack

---

**Input:** Discrete prior  $\pi = \{z_1, \dots, z_n\}$ , Model parameters  $\{\theta_1, \theta_1, \dots, \theta_T\}$ , Privatized gradients (with known gradients subtracted)  $\{\bar{g}_1, \dots, \bar{g}_T\}$ , sampling probability  $q$ , function that takes the top  $qT$  values from a set of observed gradients  $top_{qT}$   
**Observations:**  $\mathcal{O} \leftarrow \{\}$   
**Output:** Reconstruction guess  $\hat{z} \in \pi$   
**for**  $i \in [1, 2, \dots, n]$  **do**  
     $\mathcal{R} \leftarrow \{\}$   
    **for**  $t \in [1, 2, \dots, T]$  **do**  
         $\mathcal{R}[t] \leftarrow \langle \text{clip}_C(\nabla_{\theta_t} \ell(\theta_t, z_i)), \bar{g}_t \rangle$   
    **end for**  
     $\mathcal{R} \leftarrow top_{qT}(\mathcal{R})$   
     $\mathcal{O}[i] \leftarrow sum(\mathcal{R})$   
**end for**  
 $\hat{i} \leftarrow \arg \max \mathcal{O}$   
**return**  $\hat{z} \leftarrow \pi[\hat{i}]$

---

647 **G Alternative variant of the prior-aware attack**

648 Here, we state an alternative attack that uses the log-likelihood to find out which point in the prior  
649 set is used for training. Assume we have  $T$  steps with clipping threshold  $C = 1$ , noise  $\sigma$ , and the  
650 sampling rate is  $q$ .

651 Let  $\bar{g}_1, \dots, \bar{g}_T$  be the observed gradients minus the gradient of the examples that are known to be in  
652 the batch and let  $l_1, \dots, l_T$  be the  $\ell_2$  norms of these gradients.

653 For each example  $z$  in the prior set let  $g_1^z, \dots, g_T^z$  be the clipped gradient of the example on the  
654 intermediate model. Also let  $l_1^z, \dots, l_T^z$  be the  $\ell_2$  norms of  $(\bar{g}_1 - g_1^z), \dots, (\bar{g}_T - g_T^z)$ .

655 Now we describe the optimal attack based on  $l_i^z$ . For each example  $z$ , calculate the following:

656  $s_z = \sum_{i \in [T]} \ln(1 - q + qe^{-\frac{(l_i^z)^2 + l_i^2}{2\sigma^2}})$ . It is easy to observe that this is the log probability of outputting  
657 the steps conditioned on  $z$  being used in the training set. Then since the prior is uniform over the  
658 prior set, we can choose the  $z$  with maximum  $s_z$  and report that as the example in the batch.

659 In fact, this attack could be extended to the non-uniform prior by choosing the example that maximizes  
660  $s_z \cdot p_z$ , where  $p_z$  is the original probability of  $z$ .

661 **H Comparison with Guo et al. (2022b)**

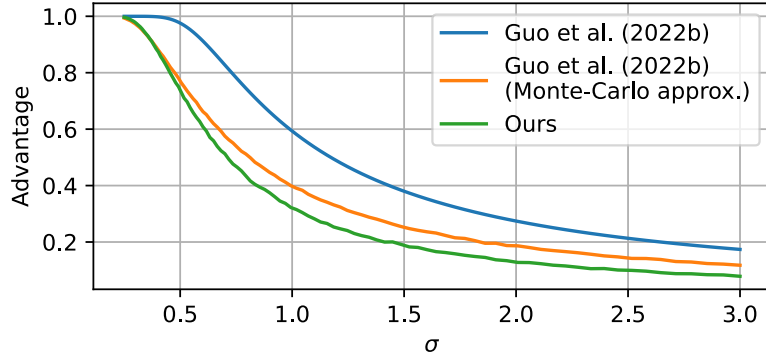


Figure 11: Comparison of our upper bound on advantage (Equation (4)) with Guo et al. (2022b) as function of  $\sigma$  for a uniform prior of size ten. We use a single step of DP-SGD with no mini-batch subsampling, and use 100,000 samples for Monte-Carlo approximation.

Table 2: Comparison of our upper bound on advantage (Equation (4)) with Guo et al. (2022b) and the Guo et al. (2022b) Monte-Carlo approximation (abbreviated to MC) as function of  $\sigma$  for a uniform prior size of ten and one hundred.

| Prior size | Method                  | Advantage upper bound |              |              |              |              |              |
|------------|-------------------------|-----------------------|--------------|--------------|--------------|--------------|--------------|
|            |                         | $\sigma$              |              |              |              |              |              |
|            |                         | 0.5                   | 1            | 1.5          | 2            | 2.5          | 3            |
| 10         | Guo et al. (2022b)      | 0.976                 | 0.593        | 0.380        | 0.274        | 0.213        | 0.174        |
|            | Guo et al. (2022b) (MC) | 0.771                 | 0.397        | 0.257        | 0.184        | 0.144        | 0.118        |
|            | Ours                    | <b>0.737</b>          | <b>0.322</b> | <b>0.189</b> | <b>0.128</b> | <b>0.099</b> | <b>0.080</b> |
| 100        | Guo et al. (2022b)      | 0.861                 | 0.346        | 0.195        | 0.131        | 0.097        | 0.076        |
|            | Guo et al. (2022b) (MC) | 0.549                 | 0.210        | 0.120        | 0.081        | 0.062        | 0.049        |
|            | Ours                    | <b>0.362</b>          | <b>0.077</b> | <b>0.035</b> | <b>0.024</b> | <b>0.018</b> | <b>0.012</b> |

662 Recently, Guo et al. (2022b) have analyzed reconstruction of discrete training data. They note that  
 663 DP bounds the mutual information shared between training data and learned parameters, and use  
 664 Fano’s inequality to convert this into a bound on reconstruction success. In particular, they define the  
 665 advantage of the adversary as

$$\text{Adv} := \frac{p_{\text{adversary success}} - p_{\pi}^{\max}}{1 - p_{\pi}^{\max}} \in [0, 1]. \quad (4)$$

666 where  $p_{\pi}^{\max}$  is the maximum sampling probability from the prior,  $\pi$ , and  $p_{\text{adversary success}}$  is the probabil-  
 667 ity that the adversary is successful at inferring which point in the prior was included in training. They  
 668 then bound the advantage by lower bounding the adversary’s error  $t := 1 - p_{\text{adversary success}}$  and by  
 669 appealing to Fano’s inequality they show this can be done by finding the smallest  $t \in [0, 1]$  satisfying

$$\begin{aligned} f(t) := & H(\pi) - I(\pi; w) + t \log t + (1 - t) \log(1 - t) \\ & - t \log(|\pi| - 1) \leq 0, \end{aligned} \quad (5)$$

670 where  $w$  is output of the private mechanism,  $H(\pi)$  is the entropy of the prior, and  $I(\pi; w)$  is the  
 671 mutual information between the prior and output of the private mechanism. For an  $(\alpha, \epsilon)$ -RDP  
 672 mechanism,  $I(\pi; w) \leq \epsilon$ , and so  $I(\pi; w)$  can be replaced by  $\epsilon$  in Equation (5). However, Guo et al.  
 673 (2022b) show that for the Gaussian mechanism, this can be improved upon either by using a Monte-Carlo  
 674 approximation of  $I(\pi; w)$  — this involves approximating the KL divergence between a Gaussian and  
 675 a Gaussian mixture — or by showing that  $I(\pi; w) \leq -\sum_{i=1}^{|\pi|} p_{\pi}^i \log \left( p_{\pi}^i + (1 - p_{\pi}^i) \exp \left( \frac{-\Delta^2}{2\sigma^2} \right) \right)$ ,



676 where  $\Delta$  is the sensitivity of the mechanism, and  $p_\pi^i$  is the probability of selecting the  $i$ th element  
 677 from the prior. We use a uniform prior in all our experiments and so  $H(\pi) = -\log(\frac{1}{|\pi|})$  and  
 678  $p_\pi^i = p_\pi^{\max} = \frac{1}{|\pi|}$ .

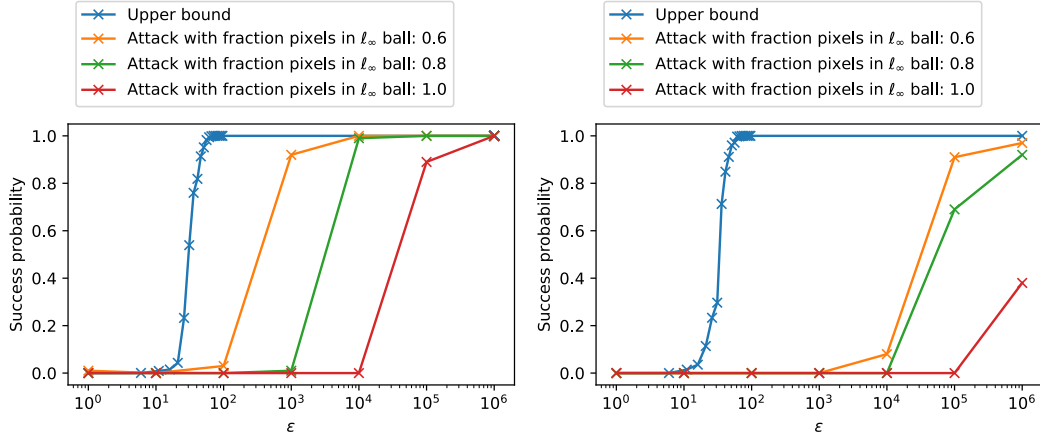
679 We convert our bound on success probability to advantage and compare with the [Guo et al. \(2022b\)](#)  
 680 upper bound (and its Monte-Carlo approximation) in Figure [11](#) and Table [2](#), and note our bound is  
 681 tighter.

## 682 I Experiments with *very* small priors (aka. experiments where the adversary 683 has no background knowledge about the target)

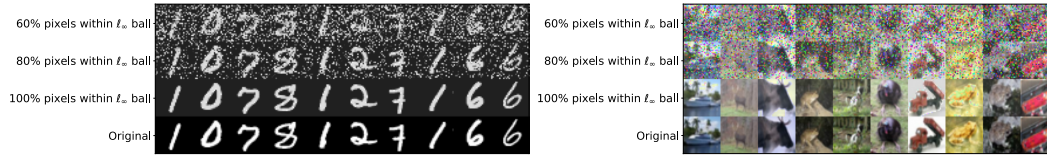
684 Our experiments in Section [3](#) and Section [4](#) were conducted with an adversary who has side informa-  
 685 tion about the target point. Here, we reduce the amount of background knowledge the adversary has  
 686 about the target, and measure how this affects the reconstruction upper bound and attack success.

687 We do this in the following set-up: Given a target  $z$ , we initialize our reconstruction from uniform  
 688 noise and optimize with the gradient-based reconstruction attack introduced in Section [2](#) to produce  
 689  $\hat{z}$ . We mark  $\hat{z}$  as a successful reconstruction of  $z$  if  $\frac{1}{d} \sum_{i=1}^d \mathbb{I}[|z[i] - \hat{z}[i]| < \delta] \geq \rho$ , where  $\rho \in [0, 1]$ ,  
 690  $d$  is the data dimensionality, and we set  $\delta = \frac{32}{255}$  in our experiments. If  $\rho = 1$  this means we mark  
 691 the reconstruction as successful if  $\|\hat{z} - z\|_\infty < \delta$ , and for  $\rho < 1$ , then at least a fraction  $\rho$  values  
 692 in  $\hat{z}$  must be within an  $\ell_\infty$  ball of radius  $\delta$  from  $z$ . Under the assumption the adversary has no  
 693 background knowledge of the target point, with  $\delta = \frac{32}{255}$  and a uniform prior, the prior probability of  
 694 reconstruction is given by  $(2 \times 32/256)^{d\rho}$  — if  $\rho = 1$ , for MNIST and CIFAR-10, this means the prior  
 695 probability of a successful reconstruction is  $9.66 \times 10^{-473}$  and  $2.96 \times 10^{-1850}$ , respectively.

696 We plot the reconstruction upper bound compared to the attack success for different values of  $\rho$  in  
 697 Figure [12](#). We also visualize the quality of reconstructions for different values of  $\rho$ . Even for  $\rho = 0.6$ ,  
 698 where 40% of the reconstruction pixels can take any value, and the remaining 60% are within an  
 699 absolute value of  $\frac{32}{255}$  from the target, one can easily identify that the reconstructions look visually  
 700 similar to the target.



(a) Comparison of reconstruction success under a *very* small prior for MNIST, where we judge a reconstruction as successful if at least  $\rho$  pixels are within an absolute distance of  $\frac{32}{255}$  of the target. (b) Comparison of reconstruction success under a *very* small prior for CIFAR-10, where we judge a reconstruction as successful if at least  $\rho$  pixels are within an absolute distance of  $\frac{32}{255}$  of the target.



(c) MNIST examples of reconstructions where at least  $\rho$  pixels are within an absolute distance of  $\frac{32}{255}$  of the target. (d) CIFAR-10 examples of reconstructions where at least  $\rho$  pixels are within an absolute distance of  $\frac{32}{255}$  of the target.

Figure 12: Comparison of reconstruction success under a *very* small prior. The prior probability of success for MNIST and CIFAR-10 are  $9.66 \times 10^{-473}$  and  $2.96 \times 10^{-1850}$ , respectively.

## 701 J Estimating $\kappa$ from samples

702 Here, we discuss how to estimate the base probability of reconstruction success,  $\kappa$ , if the adversary  
 703 can only sample from the prior distribution.

704 Let  $\hat{\pi}$  be the empirical distribution obtained by taking  $N$  independent samples from the prior and  
 705  $\hat{\kappa} = \kappa_{\hat{\pi}, \rho}(\eta)$  be the corresponding parameter for this discrete approximation to  $\pi$  – this can be  
 706 computed using the methods sketched in Section 3. Then we have the following concentration bound.

707 **Proposition 5.** *With probability  $1 - e^{-N\tau^2\kappa/2}$  we have*

$$\kappa \leq \frac{\hat{\kappa}}{1 - \tau}.$$

708 The proof is given in Appendix M

## 709 K Discussion on related work

710 Here, we give a more detailed discussion of relevant related work over what is surfaced in Section I  
 711 and Section 2.

712 **DP and reconstruction.** By construction, differential privacy bounds the success of a membership  
 713 inference attack, where the aim is to infer if a point  $z$  was in or out of the training set. While  
 714 the connection between membership inference and DP is well understood, less is known about the  
 715 relationship between training data reconstruction attacks and DP. A number of recent works have  
 716 begun to remedy this in the context of models trained with DP-SGD by studying the value of  $\epsilon$



717 required to thwart training data reconstruction attacks (Bhowmick et al., 2018; Balle et al., 2022;  
718 Guo et al., 2022a,b; Stock et al., 2022). Of course, because differential privacy bounds membership  
719 inference, it will also bound ones ability to reconstruct training data; if one cannot determine if  $z$  was  
720 used in training, they will not be able to reconstruct that point. These works are interested in both  
721 formalizing training data reconstruction attacks, and quantifying the necessary  $\epsilon$  required to bound its  
722 success. Most of these works share a common finding – the  $\epsilon$  value needed for this bound is much  
723 larger than the value required to protect against membership inference attacks ( $< 10$  in practice).  
724 If all other parameters in  $q\sqrt{T \log(\frac{1}{\delta})}/\epsilon$  remain fixed, one can see that a larger value of  $\epsilon$  reduces the  
725 scale of noise we add to gradients, which in turn results in models that achieve smaller generalization  
726 error than models trained with DP-SGD that protect against membership inference.

727 The claim that a protection against membership inference attacks also protects against training data  
728 reconstruction attacks glosses over many subtleties. For example, if  $z$  was not included in training it  
729 could still have a non-zero probability of reconstruction if samples that are close to  $z$  were included  
730 in training. Balle et al. (2022) take the approach of formalizing training reconstruction attacks in a  
731 Bayesian framework, where they compute a prior probability of reconstruction, and then find how  
732 much more information an adversary gains by observing the output of DP-SGD.

733 Balle et al. (2022) use an average-case definition of reconstruction over the output of a randomized  
734 mechanism. In contrast, Bhowmick et al. (2018) define a worst-case formalization, asking when  
735 should an adversary not be able to reconstruct a point of interest regardless of the output of the  
736 mechanism. Unfortunately, such worst-case guarantees are not attainable under DP-relaxations like  
737  $(\epsilon, \delta)$ -DP and RDP, because the privacy loss is not bounded; there is a small probability that the  
738 privacy loss will be high.

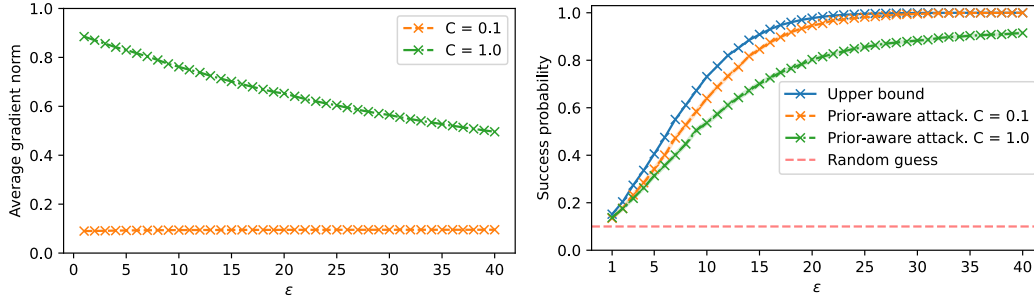
739 Stock et al. (2022) focus on bounding reconstruction for language tasks. They use the probability  
740 preservation guarantee from RDP to derive reconstruction bounds, showing that the length of a secret  
741 within a piece of text itself provides privacy. They translate this to show a smaller amount of DP  
742 noise is required to protect longer secrets.

743 While Balle et al. (2022) propose a Bayesian formalization for reconstruction error, Guo et al. (2022a)  
744 propose a frequentist definition. They show that if  $M$  is  $(2, \epsilon)$ -RDP, then the reconstruction MSE is  
745 lower bounded by  $\sum_{i=1}^d \text{diam}_i(\mathcal{Z})^2 / 4d(e^\epsilon - 1)$ , where  $\text{diam}_i(\mathcal{Z})$  is the diameter of the space  $\mathcal{Z}$  in the  $i$ th  
746 dimension.

747 **Gradient inversion attacks.** The early works of Wang et al. (2019) and Zhu et al. (2019) showed  
748 that one can invert single image representation from gradients of a deep neural network. Zhu et al.  
749 (2019) actually went beyond this and showed one can jointly reconstruct both the image and label  
750 representation. The idea is that given a target point  $z$ , a loss function  $\ell$ , and an observed gradient  
751 (wrt to model parameters  $\theta$ )  $g_z = \nabla_{\theta} \ell(\theta, z)$ , to construct a  $\hat{z}$  such that  $\hat{z} = \arg \min_{z'} \|g_{z'} - g_z\|$ .  
752 The expectation is that images that have similar gradients will be visually similar. By optimizing  
753 the above objective with gradient descent, Zhu et al. (2019) showed that one can construct visually  
754 accurate reconstruction on standard image benchmark datasets like CIFAR-10.

755 Jeon et al. (2021); Yin et al. (2021); Jin et al. (2021); Huang et al. (2021); Geiping et al. (2020)  
756 proposed a number of improvements over the reconstruction algorithm used in Zhu et al. (2019): they  
757 showed how to reconstruct multiple training points in batched gradient descent, how to optimize  
758 against batch normalization statistics, and incorporate priors into the optimization procedure, amongst  
759 other improvements.

760 The aforementioned attacks assumed an adversary has access to gradients through intermediate  
761 model updates. Balle et al. (2022) instead investigate reconstruction attacks when adversary can  
762 only observe a model after it has finished training, and propose attacks against (parametric) ML  
763 models under this threat model. However, the attack they construct is computationally demanding as  
764 it involves retraining thousands of models. This computational bottleneck is also a factor in Haim  
765 et al. (2022), who also investigate training data reconstruction attacks where the adversary has access  
766 only to final model parameters.



(a) Average gradient norm (over all samples and steps) for different values of  $\epsilon$  at  $C = 0.1$  and  $C = 1$ . (b) Reconstruction success probability for different values of  $\epsilon$  at  $C = 0.1$  and  $C = 1$ .

Figure 13: Comparison of how reconstruction success is changes with the clipping norm,  $C$ . We see that if examples have a gradient norm smaller than  $C$ , and so are not clipped, reconstruction success probability becomes smaller.

## 767 L More experiments on the effect of DP-SGD hyperparameters

768 We extend on our investigation into the effect that DP-SGD hyperparameters have on reconstruction.  
 769 We begin by varying the clipping norm parameter,  $C$ , and measure the effect on reconstruction. Fol-  
 770 lowing this, we replicate our results from Section 4 (the effect hyperparameters have on reconstruction  
 771 at a fixed  $\epsilon$ ) across different values of  $\epsilon$  and prior sizes,  $|\pi|$ .

### 772 L.1 Effect of clipping norm

773 If we look again at our attack set-up in Algorithm 2, we see that in essence we are either summing a  
 774 set of samples only from a Gaussian centred at zero or a Gaussian centred at  $C^2$ . If the gradient of  
 775 the target point is not clipped, then this will reduce the sum of gradients when the target is included in  
 776 a batch, as the Gaussian will be centred at a value smaller than  $C^2$ . This will increase the probability  
 777 that the objective is not maximized by the target point.

778 We demonstrate how this changes the reconstruction success probability by training a model for 100  
 779 steps with a clipping norm of 0.1 or 1, and measuring the average gradient norm of all samples over  
 780 all steps. Results are shown in Figure 13. We see at  $C = 0.1$ , our attack is tight to the upper bound,  
 781 and the average gradient norm is 0.1 for all values of  $\epsilon$ ; all individual gradients are clipped. When  
 782  $C = 1$ , the average gradient norm decreases from 0.9 at  $\epsilon = 1$  to 0.5 at  $\epsilon = 40$ , and we see a larger  
 783 gap between upper and lower bounds. The fact that some gradients may not be clipped is not taken  
 784 into account by our theory used to compute upper bounds, and so we conjecture that the reduction in  
 785 reconstruction success is a real effect rather than a weakness of our attack.

786 We note that these findings chime with work on individual privacy accounting (Feldman & Zrnic,  
 787 2021; Yu et al., 2022; Ligett et al., 2017; Redberg & Wang, 2021). An individual sample’s privacy  
 788 loss is often much smaller than what is accounted for by DP bounds. These works use the gradient  
 789 norm of an individual sample to measure the true privacy loss, the claim is that if the gradient norm is  
 790 smaller than the clipping norm, the amount of noise added is too large, as the DP accountant assumes  
 791 all samples are clipped. Our experiments support the claim that there is a disparity in privacy loss  
 792 between samples whose gradients are and are not clipped.

### 793 L.2 More results on the effect of DP-SGD hyperparameters at a fixed $\epsilon$

794 In Section 4, we demonstrated that the success of a reconstruction attack cannot be captured only  
 795 by the  $(\epsilon, \delta)$  guarantee, when  $\epsilon = 4$  and the size of the prior,  $\pi$ , is set to ten. We now observe how  
 796 these results change across different  $\epsilon$  and  $|\pi|$ , where we again fix the number of updates to  $T = 100$ ,  
 797  $C = 1$ , vary  $q \in [0.01, 0.99]$ , and adjust  $\sigma$  accordingly.

798 Firstly, in Figure 14, we measure the upper and lower bound ((improved) prior-aware attack) on  
 799 the probability of successful reconstruction across different  $q$ . In all settings, we observe smaller  
 800 reconstruction success at smaller  $q$ , where the largest fluctuations in reconstruction success are for

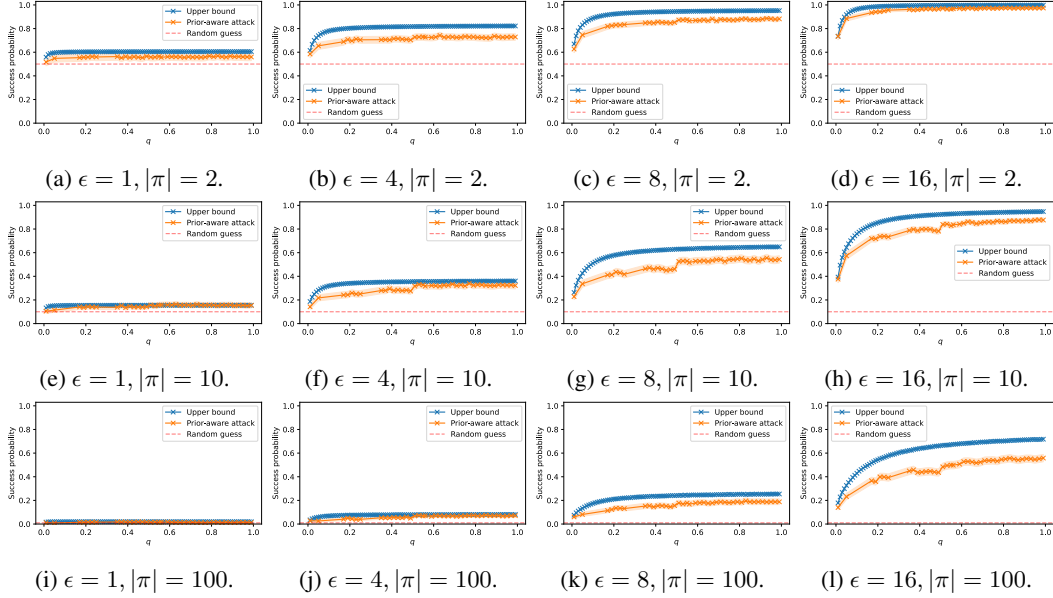


Figure 14: How the upper bound and (improved) prior-aware attack change as a function of  $q$  at a fixed value of  $\epsilon$  and prior size,  $|\pi|$ . The amount of privacy leaked through a reconstruction at a fixed value of  $\epsilon$  can change with different  $q$ .

larger values of  $\epsilon$ . We visualise this in another way by plotting  $\sigma$  against  $q$  and report the upper bound in Figure 15. Note that the color ranges in Figure 15 are independent across subfigures.

## M Proofs

Throughout the proofs we make some of our notation more succinct for convenience. For a probability distribution  $\omega$  we write  $\omega(E) = \mathbb{P}_\omega[E]$ , and rewrite  $\mathcal{B}_\kappa(\mu, \nu) = \sup\{\mathbb{P}_\mu[E] : E \text{ s.t. } \mathbb{P}_\nu[E] \leq \kappa\}$  as  $\sup_{\nu(E) \leq \kappa} \mu(E)$ . Given a distribution  $\omega$  and function  $\phi$  taking values in  $[0, 1]$  we also write  $\omega(\phi) = \mathbb{E}_{X \sim \omega}[\phi(X)]$ .

### M.1 Proof of Theorem 2

We say that a pair of distributions  $(\mu, \nu)$  is *testable* if for all  $\kappa \in [0, 1]$  we have

$$\inf_{\nu(\phi) \leq \kappa} (1 - \mu(\phi)) = \inf_{\nu(E) \leq \kappa} (1 - \mu(E)) ,$$

where the infimum on the left is over all  $[0, 1]$ -valued measurable functions and the one on the right is over measurable events (i.e.  $\{0, 1\}$ -valued functions). The Neyman-Pearson lemma (see e.g. Lehmann & Romano (2005)) implies that this condition is satisfied whenever the statistical hypothesis problem of distinguishing between  $\mu$  and  $\nu$  admits a uniformly most powerful test. For example, this is the case for distributions on  $\mathbb{R}^d$  where the density ratio  $\mu/\nu$  is a continuous function.

**Theorem 6** (Formal version of Theorem 2). *Fix  $\pi$  and  $\rho$ . Suppose that for every fixed dataset  $D$ . there exists a distribution  $\mu_{D_\cdot}$  such that  $\sup_{z \in \text{supp}(\pi)} \mathcal{B}_\kappa(\mu_{D_\cdot}, \nu_{D_\cdot}) \leq \mathcal{B}_\kappa(\mu_{D_\cdot}, \nu_{D_\cdot})$  for all  $\kappa \in [0, 1]$ . If the pair  $(\mu, \nu)$  is testable, then  $M$  is  $(\eta, \gamma)$ -ReRo with*

$$\gamma = \sup_{D_\cdot} \sup_{\nu_{D_\cdot}(E) \leq \kappa_{\pi, \rho}(\eta)} \mu_{D_\cdot}(E) .$$

The following lemma from Dong et al. (2019) will be useful.

**Lemma 7.** *For any  $\mu$  and  $\nu$ , the function  $\kappa \mapsto \inf_{\nu(\phi) \leq \kappa} (1 - \mu(\phi))$  is convex in  $[0, 1]$ .*

**Lemma 8.** *For any testable pair  $(\mu, \nu)$ , the function  $\kappa \mapsto \sup_{\nu(E) \leq \kappa} \mu(E)$  is concave.*

821 *Proof.* By the testability assumption we have

$$\begin{aligned}
\sup_{\nu(E) \leq \kappa} \mu(E) &= \sup_{\nu(E) \leq \kappa} \mu(E) \\
&= \sup_{\nu(E) \leq \kappa} (1 - \mu(\bar{E})) \\
&= 1 - \inf_{\nu(E) \leq \kappa} \mu(\bar{E}) \\
&= 1 - \inf_{\nu(E) \leq \kappa} (1 - \mu(E)) \\
&= 1 - \inf_{\nu(\phi) \leq \kappa} (1 - \mu(\phi)) .
\end{aligned}$$

822 Concavity now follows from Lemma 7.  $\square$

823 *Proof of Theorem 6* Fix  $D$ . and let  $\kappa = \kappa_{\pi, \rho}(\eta)$  throughout. Let also  $\nu = \nu_{D_\cdot}$ ,  $\mu_z = \mu_{D_z}$ ,  $\nu^* = \nu_{D_\cdot}^*$   
824 and  $\mu^* = \mu_{D_\cdot}^*$ .

825 Expanding the probability of successful reconstruction, we get:

$$\begin{aligned}
\mathbb{P}_{Z \sim \pi, W \sim M(D_\cdot \cup \{Z\})}[\rho(Z, R(W)) \leq \eta] &= \mathbb{E}_{Z \sim \pi} \mathbb{P}_{W \sim M(D_\cdot \cup \{Z\})}[\rho(Z, R(W)) \leq \eta] \\
&= \mathbb{E}_{Z \sim \pi} \mathbb{E}_{W \sim M(D_\cdot \cup \{Z\})} \mathbb{I}[\rho(Z, R(W)) \leq \eta] \\
&= \mathbb{E}_{Z \sim \pi} \mathbb{E}_{W \sim \mu_Z} \mathbb{I}[\rho(Z, R(W)) \leq \eta] \\
&= \mathbb{E}_{Z \sim \pi} \mathbb{E}_{W \sim \nu} \left[ \frac{\mu_Z(W)}{\nu(W)} \mathbb{I}[\rho(Z, R(w)) \leq \eta] \right] .
\end{aligned}$$

826 Now fix  $z \in \text{supp}(\pi)$  and let  $\kappa_z = \mathbb{P}_{W \sim \nu}[\rho(z, R(W)) \leq \eta]$ . Using the assumption on  $\mu^*$  we get:

$$\begin{aligned}
\mathbb{E}_{W \sim \nu} \left[ \frac{\mu_z(W)}{\nu(W)} \mathbb{I}[\rho(z, R(w)) \leq \eta] \right] &\leq \sup_{\nu(E) \leq \kappa_z} \mathbb{E}_{W \sim \nu} \left[ \frac{\mu_z(W)}{\nu(W)} \mathbb{I}[W \in E] \right] && \text{(By definition of } \kappa) \\
&= \sup_{\nu(E) \leq \kappa_z} \mathbb{E}_{W \sim \mu_z} [\mathbb{I}[W \in E]] \\
&= \sup_{\nu(E) \leq \kappa_z} \mu_z(E) \\
&\leq \sup_{\nu^*(E) \leq \kappa_z} \mu^*(E) . && \text{(By definition of } \mu^* \text{ and } \nu^*)
\end{aligned}$$

827 Finally, using Lemma 8 and Jensen's inequality on the following gives the result:

$$\begin{aligned}
\mathbb{E}_{Z \sim \pi}[\kappa_Z] &= \mathbb{E}_{Z \sim \pi} \mathbb{P}_{W \sim \nu}[\rho(Z, R(W)) \leq \eta] \\
&= \mathbb{E}_{W \sim \nu} \mathbb{P}_{Z \sim \pi}[\rho(Z, R(W)) \leq \eta] \\
&\leq \mathbb{E}_{W \sim \nu} \kappa \\
&= \kappa . \quad \square
\end{aligned}$$

## 828 M.2 Proof of Corollary 3

829 Here we prove Corollary 3. We will use the following shorthand notation for convenience:  $\mu =$   
830  $\mathcal{N}(B(T, q), \sigma^2 I)$  and  $\nu = \mathcal{N}(0, \sigma^2 I)$ . To prove our result, we use the notion of  $TV_a$ .

**Definition 9** (Mahloujifar et al. (2022)). For two probability distributions  $\omega_1(\cdot)$  and  $\omega_2(\cdot)$ ,  $TV_a$  is defined as

$$TV_a(\omega_1, \omega_2) = \int |\omega_1(x) - a \cdot \omega_2(x)| dx.$$

831 Now we state the following lemma borrowed from Mahloujifar et al. (2022).

832 **Lemma 10** (Theorem 6 in Mahloujifar et al. (2022)). Let  $\nu_{D_\cdot}$ ,  $\mu_{D_z}$  be the output distribution of  
833 DP-SGD applied to  $D$ . and  $D_z$  respectively, with noise multiplier  $\sigma$ , sampling rate  $q$ . Then we have

$$TV_a(\nu_{D_\cdot}, \mu_{D_z}) \leq TV_a(\nu, \mu) .$$

834 Now, we state the following lemma that connects  $TV_a$  to blow-up function.

**Lemma 11** (Lemma 21 in [Zhu et al. \(2022\)](#)). *For any pair of distributions  $\omega_1, \omega_2$  we have*

$$\sup_{\omega_1(E) \leq \kappa} \omega_2(E) = \inf_{a > 1} \min \left\{ 0, a \cdot \kappa + \frac{TV_a(\omega_1, \omega_2) + 1 - a}{2}, \frac{2\kappa + TV_a(\omega_1, \omega_2) + a - 1}{2a} \right\}$$

835 Since  $TV_a(\nu_{D_-}, \mu_{D_z})$  is bounded by  $TV_a(\nu, \mu)$  for all  $a$ , therefore we have

$$\sup_{\nu_{D_-}(E) \leq \kappa} \mu_{D_z}(E) \leq \sup_{\nu(E) \leq \kappa} \mu(E) .$$

### 836 M.3 Proof of Proposition 5

837 Recall  $\kappa = \sup_{z_0 \in \mathcal{Z}} \mathbb{P}_{Z \sim \pi}[\rho(Z, z_0) \leq \eta]$  and  $\hat{\kappa} = \sup_{z_0 \in \mathcal{Z}} \mathbb{P}_{Z \sim \hat{\pi}}[\rho(Z, z_0) \leq \eta]$ . Let  $\kappa_z =$   
 838  $\mathbb{P}_{Z \sim \pi}[\rho(Z, z) \leq \eta]$  and  $\hat{\kappa}_z = \mathbb{P}_{Z \sim \hat{\pi}}[\rho(Z, z) \leq \eta]$ . Note  $\hat{\kappa}_z$  is the sum of  $N$  i.i.d. Bernoulli random  
 839 variables and  $\mathbb{E}_{\hat{\pi}}[\hat{\kappa}_z] = \kappa_z$ . Then, using a multiplicative Chernoff bound, we see that for a fixed  $z$   
 840 the following holds with probability at least  $1 - e^{-N\tau_z^2 \kappa/2}$ :

$$\kappa_z \leq \frac{\hat{\kappa}_z}{1 - \tau} .$$

841 Applying this to  $z^* = \arg \sup_{z_0 \in \mathcal{Z}} \mathbb{P}_{Z \sim \pi}[\rho(Z, z_0) \leq \eta]$  we get that the following holds with  
 842 probability at least  $1 - e^{-N\tau^2 \kappa/2}$ :

$$\kappa = \kappa_{z^*} \leq \frac{\hat{\kappa}_{z^*}}{1 - \tau} \leq \frac{\hat{\kappa}}{1 - \tau} .$$

### 843 M.4 Proof of Proposition 4

Let  $z = \frac{(r'_{N'} + r'_{N'-1})}{2}$ . Let  $E_1$  be the event that  $|\mathbb{P}[r > z] - \kappa| \geq \tau$ . By applying Chernoff-Hoeffding bound we have  $\mathbb{P}[E_1] \leq 2e^{-2N\tau^2}$ . Now note that since  $\mu$  is a Gaussian mixture, we can write  $\mu = \sum_{i \in [2T]} a_i \mu_i$  where each  $\mu_i$  is a Gaussian  $\mathcal{N}(c_i, \sigma)$  where  $|c_i|_2 \leq \sqrt{T}$ . Now let  $r_i = \mu_i(W)/\nu(W)$ . By holder, we have  $\mathbb{E}[r^2] \leq \sum a_i \mathbb{E}[r_i^2]$ . We also now that  $\mathbb{E}[r_i^2] \leq e^T$ , therefore,  $\mathbb{E}[r^2] \leq e^T$ . Now let  $E_2$  be the event that  $|\mathbb{E}[r \cdot I(r > z)] - \gamma'| \geq \tau$ . Since the second moment of  $r$  is bounded, the probability of  $E_2$  goes to zero as  $N$  increases. Therefore, almost surely we have

$$\sup_{\nu(E) \leq \kappa - \tau} \mu(E) - \tau \leq \lim_{N \rightarrow \infty} \gamma' \leq \sup_{\nu(E) \leq \kappa + \tau} \mu(E) + \tau.$$

Now by pushing  $\tau$  to 0 and using the fact that  $\mu$  and  $\nu$  are smooth we have

$$\lim_{N \rightarrow \infty} \gamma' = \lim_{\tau \rightarrow 0} \sup_{\nu(E) \leq \kappa + \tau} \mu(E) + \tau = \sup_{\nu(E) \leq \kappa} \mu(E).$$

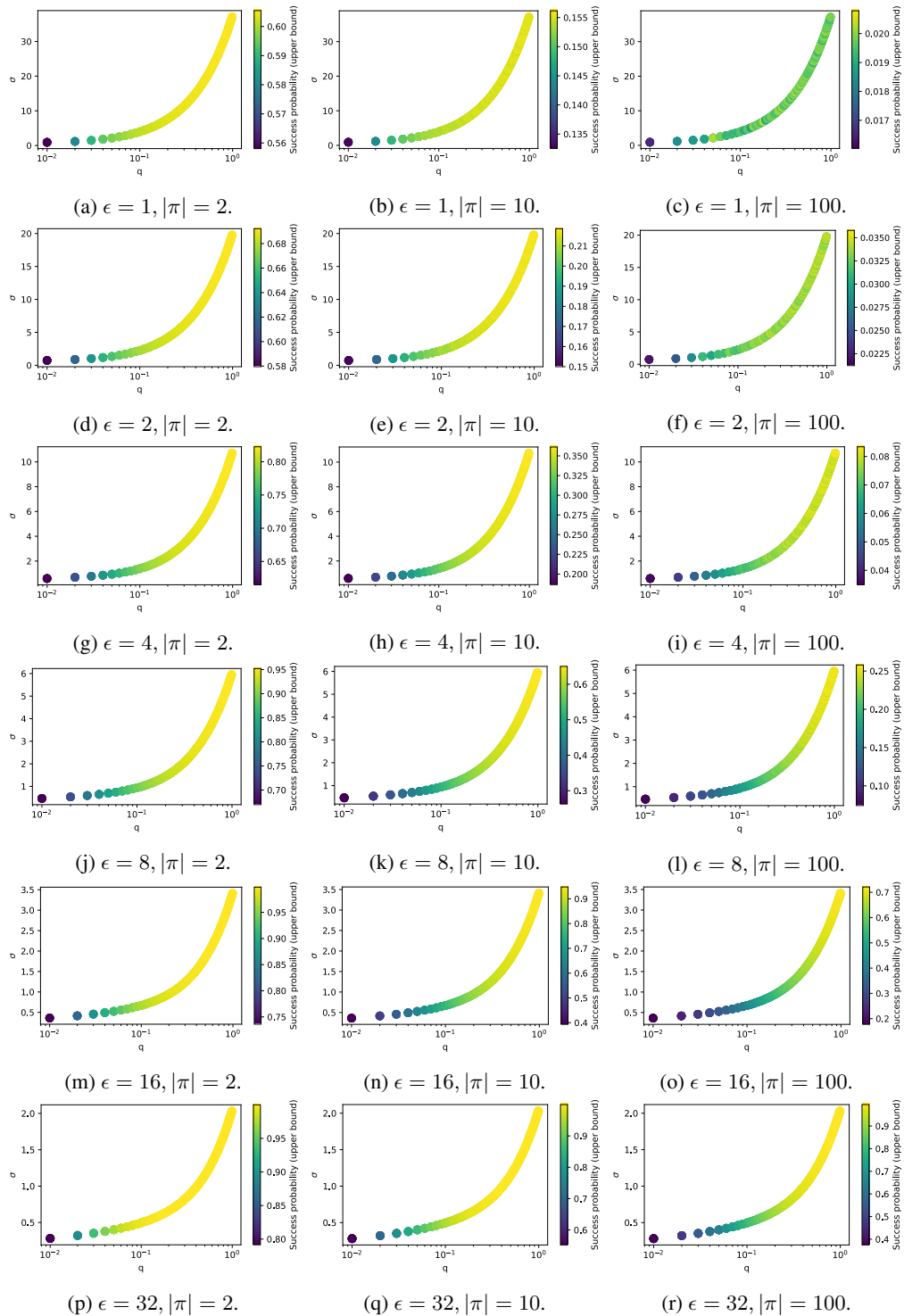


Figure 15: How the upper bound changes as a function of  $q$  and  $\sigma$  at a fixed value of  $\epsilon$  and prior size,  $|\pi|$ , and setting  $T = 100$ . The probability of a successful reconstruction can vary widely with different values of  $q$ . For example, at  $\epsilon = 32$  and  $|\pi| = 100$ , at  $q = 0.01$  the upper bound is 0.4 and at  $q = 0.99$  it is 1. Note that the color ranges are independent across subfigures.