

# Geminio: Language-Guided Gradient Inversion Attacks in Federated Learning

Junjie Shan<sup>1</sup> Ziqi Zhao<sup>1</sup> Jialin Lu<sup>1</sup> Rui Zhang<sup>2</sup> Siu Ming Yiu<sup>1</sup> Ka-Ho Chow<sup>1\*</sup>

<sup>1</sup>School of Computing and Data Science, The University of Hong Kong

<sup>2</sup>Department of Computing, The Hong Kong Polytechnic University

## Abstract

Foundation models that bridge vision and language have made significant progress, inspiring numerous life-enriching applications. However, their potential for misuse to introduce new threats remains largely unexplored. This paper reveals that vision-language models (VLMs) can be exploited to overcome longstanding limitations in gradient inversion attacks (GIAs) within federated learning (FL), where an FL server reconstructs private data samples from gradients shared by victim clients. Current GIAs face challenges in reconstructing high-resolution images, especially when the victim has a large local data batch. While focusing reconstruction on valuable samples rather than the entire batch is promising, existing methods lack the flexibility to allow attackers to specify their target data. In this paper, we introduce Geminio<sup>1</sup>, the first approach to transform GIAs into semantically meaningful, targeted attacks. Geminio enables a brand new privacy attack experience: attackers can describe, in natural language, the types of data they consider valuable, and Geminio will prioritize reconstruction to focus on those high-value samples. This is achieved by leveraging a pretrained VLM to guide the optimization of a malicious global model that, when shared with and optimized by a victim, retains only gradients of samples that match the attacker-specified query. Extensive experiments demonstrate Geminio’s effectiveness in pinpointing and reconstructing targeted samples, with high success rates across complex datasets under FL and large batch sizes and showing resilience against existing defenses.

## 1. Introduction

Federated learning (FL) is a privacy-enhancing technology for training machine learning models on data distributed across multiple clients [23, 34]. By enabling clients to share gradients rather than raw data with a coordinating server, FL has demonstrated transformative potential in privacy-

\*Corresponding Author: kachow@cs.hku.hk

<sup>1</sup>“Geminio” is a spell from the Harry Potter series that allows the caster to describe an object and obtain a duplicate of it.

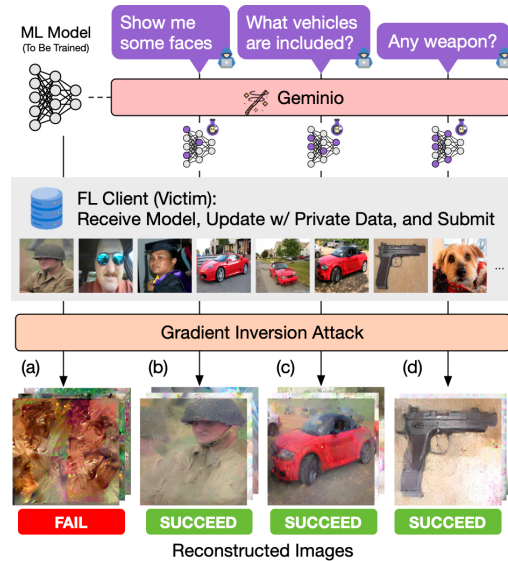


Figure 1. Geminio enables the attacker (a malicious server in FL) to describe what kind of data is valuable to them and prioritize gradient inversion to recover those images from a large data batch.

sensitive domains [13, 22, 45]. However, FL is vulnerable to various malicious attacks, with gradient inversion attacks (GIAs) posing a particularly critical threat [19, 43]. These attacks enable a malicious FL server to reconstruct private data samples from the gradients shared by a victim client, leading to privacy breaches and decelerating FL’s adoption.

GIAs face a longstanding challenge: they can only reconstruct images from gradients produced by a small batch of data [10, 31]. This limitation exists because GIAs rely on searching for data that reproduces the victim-submitted gradients, and the search space expands exponentially with batch size. Thus, much research has focused on whether this limitation is fundamental, as existing GIAs struggle with practical batch sizes. While some approaches incorporate image priors (e.g., spatial smoothness [14]) to facilitate the search, a performance gap still remains, as illustrated in Figure 1(a), with images reconstructed from a batch of 128 samples. Another direction has been to narrow the scope to reconstruct only a subset of samples. While promising, existing methods lack a semantically meaningful way for

the adversary to specify which samples are preferred and can only target, e.g., outliers [35] or images with particular brightness levels [11]. This raises an intriguing question: can reconstruction efforts be prioritized toward the data samples that truly matter most to the adversary? If so, how can we allow the adversary to specify their preferences in a meaningful, flexible, and generic way?

In this paper, we empower gradient inversion attacks with a natural language interface and propose Geminio. It enables the FL server to provide a natural language query describing the data of interest, allowing Geminio to prioritize and reconstruct matching data samples. Taking the batch of images from a victim’s mobile phone in Figure 1 as an example, the adversary could submit queries like (b) “show me some faces” to retrieve images containing faces to see the victim or their friends, (c) “what vehicles are included?” to identify cars associated with the victim, or (d) “any weapon?” to detect if the victim owns a weapon. The query does not need to relate to the FL system’s ML task. By prioritizing reconstruction efforts, Geminio can pinpoint and retrieve targeted samples from large batches, offering high flexibility in defining valuable data. This capability is achieved by misusing pretrained vision-language models (VLMs) [25] to help craft a malicious global model. When shared and optimized by the victim client, gradients become dominated by samples that match the query. Existing reconstruction optimization algorithms [14, 39, 40, 49] can consume such gradients to recover high-quality, targeted data.

Our main contributions are summarized as follows. First, we explore the misuse of pretrained VLMs to bridge the gap in gradient inversion, enabling semantically meaningful, targeted attacks. We investigate the first natural language interface for the adversary to describe the data samples that truly matter and prioritize them for reconstruction. Second, we propose Geminio, which exploits a VLM to reshape the loss surface of a global model, so that once optimized locally by the victim, the gradients are dominated by the samples matching the query. This method complements existing reconstruction optimizations and can augment them as targeted attacks. Third, we reveal the limitations of current defenses, discuss potential design improvements, and highlight their shortcomings to motivate future work. Experiments were conducted across three datasets, five attack methods, and four defense mechanisms, and various configurations to assess the threat posed by Geminio. The source code of Geminio is available at <https://github.com/HKU-TASR/Geminio>.

## 2. Background

### 2.1. Gradient Inversion Attacks in FL

**Federated Learning.** Let  $F_\theta$  be the ML model trained via FL with a loss function  $\mathcal{L}$ . At each learning round  $t$ , the

FL server sends the current global model parameters  $\theta_t$  to FL clients. Under the FedSGD protocol [23], each client  $i$  samples a data batch  $\mathcal{B}_t^i$ , having pairs of input  $\mathbf{x}$  and label  $y$ , from its private dataset to optimize the received model and submit the gradients

$$\mathcal{G}(\mathcal{B}_t^i; \theta_t) = \frac{1}{|\mathcal{B}_t^i|} \sum_{(\mathbf{x}, y) \in \mathcal{B}_t^i} \nabla_{\theta_t} \mathcal{L}(F_{\theta_t}(\mathbf{x}); y) \quad (1)$$

to the server. Then, the server aggregates the gradients submitted by all clients to update the global model parameters for the next round. The FL protocol has different variations, such as randomly selecting a subset of clients to participate in each round or running several data batches locally before submitting the gradients to the server [23].

**Gradient Inversion Attacks.** The FL server that receives the gradients  $\mathcal{G}(\mathcal{B}_t^i; \theta_t)$  from a participating client  $i$  at round  $t$  can reconstruct the private data batch  $\mathcal{B}_t^i$  via gradient inversion attacks. As the batch size should be transparent to the FL server for proper aggregation, it can randomly initialize a batch of data samples  $\bar{\mathcal{B}}$  and use the global model parameters  $\theta_t$  shared with the victim at the beginning of the learning round for reconstruction optimization:

$$\bar{\mathcal{B}}^* = \operatorname{argmin}_{\bar{\mathcal{B}}} [\delta(\mathcal{G}(\mathcal{B}_t^i; \theta_t), \mathcal{G}(\bar{\mathcal{B}}; \theta_t)) + \mathcal{R}(\bar{\mathcal{B}})], \quad (2)$$

where  $\delta$  is a distance function that measures the dissimilarity between two sets of gradients, and  $\mathcal{R}$  is a regularization function. The overarching idea is to optimize those random data samples in such a way that they can reproduce the gradients shared by the victim client.

### 2.2. Related Work

Early researches leverage the inherent shallow leakage in fully connected layers [47, 49] to reconstruct private data from gradient updates. Recent studies advance optimization and analytic methods to improve the accuracy and scalability of these attacks, but the underlying objective remains the same, which is to extract private information from gradients or model parameters [9, 10, 24, 27, 36–38].

**Reconstruction Optimization.** The reconstruction optimization in Equation 2 has been the primary focus of GIA advancements. These improvements target (i) enhanced distance functions, such as Euclidean distance [49] and cosine similarity [14, 39] to measure gradient alignment, and (ii) refined regularization methods to incorporate prior knowledge like spatial smoothness [14, 40]. Such approaches have successfully extended GIAs to support high-resolution image reconstructions, previously limited to toy datasets like MNIST [49]. However, these methods attempt to recover the entire private batch, struggling with practical batch sizes due to the vast search space involved [26].

**Narrowing the Reconstruction Scope.** Recognizing this limitation, recent work has explored reconstructing only a

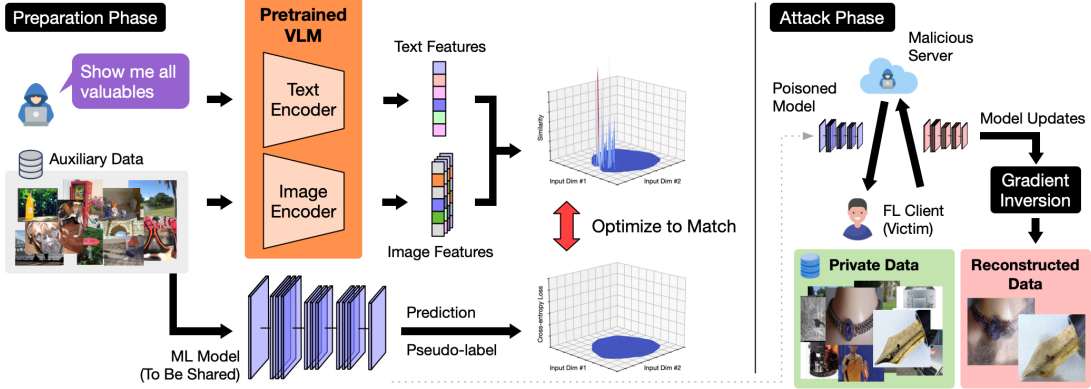


Figure 2. Geminio begins with a preparation phase. It receives a query from the attacker and uses a pretrained VLM and an auxiliary, unlabeled dataset to reshape the loss surface of the global model. Once the model is poisoned, we can send it to the client in the next FL round. The gradients received from the victim can be fed to existing reconstruction optimization methods to recover images of valuables.

subset of data samples in the private batch by manipulating the global model shared with the victim client. Abandon [1], Robbing [11], GradViT [15], LOKI [48], and SEER [12] require specific neural architectures in FL and adjust model parameters to retain only the gradients of selected private samples during the victim’s local training. These “trapped” samples are either random or satisfy simple conditions, such as brightness levels or average color intensity. Although this narrows the reconstruction scope, it provides the adversary with almost no meaningful control over which specific samples are recovered. The unusual neural architecture can be suspicious, even if it enables analytical reconstructions. Fishing [35] and GradFilt [44] improve this by setting certain model parameters at the output layer to very large values, causing the gradients to be dominated by one or all samples of a particular class, after which reconstruction optimization algorithms can be employed to recover them. However, this control remains restrictive, as the adversary can neither specify finer-grained sample characteristics within a class nor define conditions irrelevant to the FL system’s ML task. Also, the unnaturally large parameter values can be easily detected. Imperio [3] is the first method to leverage foundation models for implementing natural language-guided backdoor attacks. However, there has been no similar approach for GIA, leaving a gap in targeted attacks via natural language instructions. To address these limitations, this paper introduces Geminio, offering the first natural language interface for targeted GIAs.

**Defenses.** Encryption-based methods, such as homomorphic encryption [8], have been proposed to secure gradient confidentiality but are often computationally prohibitive [41] or can be circumvented if an active FL server modifies the FL protocol [2]. Gradient obfuscation techniques, such as differential privacy [30, 32, 33] and gradient pruning [46], allow FL clients to protect their data proactively. However, as our experiments reveal, these defenses

fail to mitigate the privacy threats posed by Geminio.

### 2.3. Threat Model

Consistent with prior studies [35, 44], we consider an FL server acting as an active adversary who (i) can modify the model parameters before sharing them with FL clients but not altering the neural architecture, (ii) can read the gradients submitted by a victim client and attempt to reconstruct private data samples from them, (iii) can provide a natural language description of the characteristics of data it deems valuable, and (iv) possesses an auxiliary, unlabeled image dataset that may originate from a completely different domain (e.g., public datasets like ImageNet [4] or images scraped from the Internet). The FL clients adhere to the FedSGD protocol, optimizing the received model with a batch of private data. We will consider other FL scenarios in Section 4, such as Geminio under FedAvg [23], as well as client-side defenses like gradient obfuscation and model parameter inspection.

## 3. Methodology

Figure 2 gives an overview of Geminio. It consists of two phases. During the preparation phase, Geminio takes a query  $Q$  (e.g., “show me all valuables”) from the adversary to craft malicious global model parameters  $\Theta_Q$ . During the attack phase, those parameters that are pretended to be legitimate will be shared with the victim client, who optimizes them with its private data batch  $\mathcal{B}$  and uploads the gradients  $\mathcal{G}(\mathcal{B}; \Theta_Q)$  (from Equation 1) to the FL server. Then, any existing reconstruction optimization method can be directly applied to recover those private samples relevant to the query (e.g., the necklace and the fountain pen retrieved by InvertingGrad [14]).

The overarching idea of Geminio is to craft a malicious global model such that those private samples in the victim’s data batch matching the query will dominate the

submitted gradients. Consider the scenario that only one private sample  $(\mathbf{x}_{\text{target}}, y_{\text{target}}) \in \mathcal{B}$  matches the query as an example; the malicious global model should behave as follows when being optimized by the victim client:  $\|\nabla_{\Theta_Q} \mathcal{L}(F_{\Theta_Q}(\mathbf{x}); y)\| \ll \|\nabla_{\Theta_Q} \mathcal{L}(F_{\Theta_Q}(\mathbf{x}_{\text{target}}); y_{\text{target}})\|$  for all  $\mathbf{x} \neq \mathbf{x}_{\text{target}}$ . Then, the victim-submitted gradients become  $\mathcal{G}(\mathcal{B}; \Theta_Q) \approx \frac{1}{|\mathcal{B}|} \nabla_{\Theta_Q} \mathcal{L}(F_{\Theta_Q}(\mathbf{x}_{\text{target}}); y_{\text{target}})$ , and any existing reconstruction optimization method (see Equation 2) will recover  $\mathbf{x}_{\text{target}}$  from them. To achieve such a behavior, instead of directly optimizing how the global model should produce gradients, which involves second-order derivatives and is highly unstable, we could exploit the property that the per-sample gradient magnitude  $\|\nabla_{\Theta_Q} \mathcal{L}(F_{\Theta_Q}(\mathbf{x}); y)\|$  is proportional to the per-sample loss value  $\mathcal{L}(F_{\Theta_Q}(\mathbf{x}); y)$ . Geminio’s objective is to craft a malicious model that amplifies the loss value of matched samples while suppressing the others.

Given a query, crafting such a malicious model requires two ingredients: (i) a supervisor that guides how it should react when given an image and (ii) a training dataset. Both are challenging because the supervisor should be able to associate images with text data, and the FL server should not possess many (or any) data samples. In this regard, the key enabler of Geminio is a pretrained VLM.

### 3.1. VLM-Guided Loss Surface Reshaping

Given an auxiliary dataset  $\mathcal{A}$ , Geminio exploits a pretrained VLM to measure the similarity between each auxiliary image and the query. The top 3D surface plot in Figure 2 shows the similarity surface as a function of auxiliary images (projected onto a 2D space by PCA). Some images align with the query well, while others have close-to-zero relatedness. An untrained global model has a roughly flat loss surface (see the bottom 3D surface plot in Figure 2). We need to train the malicious global model to have a loss surface matching the aforementioned similarity surface such that those irrelevant samples will lead to a zero loss and matched ones will dominate.

A VLM comprises two components [25]: an image encoder  $\mathcal{V}_{\text{image}}$  and a text encoder  $\mathcal{V}_{\text{text}}$ . They can project image and text data onto a latent space that those similar will collocate. For an auxiliary sample  $(\mathbf{x}, y) \in \mathcal{A}$ , we can calculate its similarity with the query  $\mathcal{Q}$ :  $s(\mathbf{x}; \mathcal{Q}) = \mathcal{V}_{\text{image}}(\mathbf{x})^\top \mathcal{V}_{\text{text}}(\mathcal{Q})$ . Based on the similarity score, we propose to train the malicious global model parameters  $\Theta_Q$  with the following routine. At each iteration, we sample a batch of auxiliary data  $\mathcal{B}_{\text{aux}} \subset \mathcal{A}$  and calculate the probability of each auxiliary image  $(\mathbf{x}, y) \in \mathcal{B}_{\text{aux}}$  being aligned with the query, normalized across the batch via a softmax function:

$$\alpha(\mathbf{x}; \mathcal{Q}, \mathcal{B}_{\text{aux}}) = \frac{\exp(s(\mathbf{x}; \mathcal{Q}))}{\sum_{(\mathbf{x}', y') \in \mathcal{B}_{\text{aux}}} \exp(s(\mathbf{x}'; \mathcal{Q}))}. \quad (3)$$

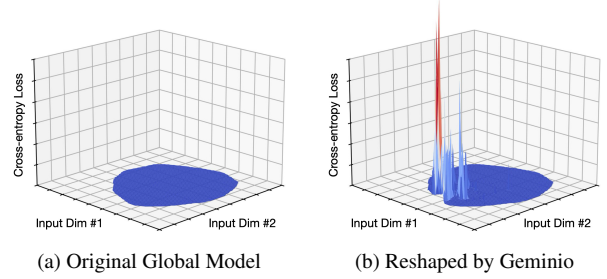


Figure 3. Geminio reshapes the loss landscape of the model such that samples matching the query will have an amplified loss to dominate gradients for targeted reconstruction.

The batch-wise normalization offers information about how one sample is more aligned with the query than another sample. Then, we can train the malicious global model parameters by minimizing

$$\begin{aligned} & \mathcal{L}^{\text{Geminio}}(\mathcal{B}_{\text{aux}}; F_{\Theta_Q}, \mathcal{Q}) \\ &= \frac{\sum_{(\mathbf{x}, y) \in \mathcal{B}_{\text{aux}}} \mathcal{L}(F_{\Theta_Q}(\mathbf{x}); y)(1 - \alpha(\mathbf{x}; \mathcal{Q}, \mathcal{B}_{\text{aux}}))}{|\mathcal{B}_{\text{aux}}| \sum_{(\mathbf{x}', y') \in \mathcal{B}_{\text{aux}}} \mathcal{L}(F_{\Theta_Q}(\mathbf{x}'); y')(1 - \alpha(\mathbf{x}'; \mathcal{Q}, \mathcal{B}_{\text{aux}}))}. \end{aligned} \quad (4)$$

Intuitively, each per-sample loss is associated with a scaling factor (the coefficient). For an auxiliary that has a strong alignment with the query, the corresponding term will be negligible since the coefficient  $(1 - \alpha(\mathbf{x}; \mathcal{Q}, \mathcal{B}_{\text{aux}}))$  is close to zero. In contrast, the term corresponding to an irrelevant auxiliary sample will have its magnitude amplified because of the large coefficient. In order for the malicious model to minimize Equation 4, it must learn to reduce the per-sample loss value of such irrelevant samples. The batch-wise normalization will then increase the per-sample loss value of matched samples. Geminio’s training routine reshapes the loss surface, originally flat (Figure 3a), to have an active response only for those matched samples (Figure 3b).

### 3.2. VLM-Guided Auxiliary Label Generation

The calculation of the per-sample loss value requires the ground-truth label of that input. However, assuming the availability of such a labeled dataset in FL is unreasonable. We propose to misuse the pretrained VLM again to launch Geminio with an unlabeled, possibly off-domain, dataset. In particular, let the class names in a  $K$ -class classification problem be  $[c_1, c_2, \dots, c_K]$ , we can generate a soft label for each auxiliary sample  $\mathbf{x}$  by measuring the similarity of its image features and the text features of each class name. Formally, the soft label  $\mathbf{y} = [y_1, y_2, \dots, y_K]$  is a probability distribution, where  $y_i$  represents the probability of  $\mathbf{x}$  being classified as class  $c_i$  and can be calculated by

$$y_i = \frac{\mathcal{V}_{\text{image}}(\mathbf{x})^\top \mathcal{V}_{\text{text}}(c_i)}{\sum_{j=1}^K \mathcal{V}_{\text{image}}(\mathbf{x})^\top \mathcal{V}_{\text{text}}(c_j)} \quad (5)$$





Figure 4. Geminiio can take task-agnostic queries from the attacker to achieve instance-level targeted reconstruction. While vanilla GIAs cannot recover recognizable images from a large batch, Geminiio narrows down the reconstruction scope and successfully rebuilds high-fidelity images that match the attacker’s queries (e.g., the handgun, rifle, and knife images given the query “Any weapon?” in (a)).

Using soft labels in the cross-entropy loss function, one could launch Geminiio by simply using public datasets or scraping images from the internet.

## 4. Empirical Evaluation

We conduct extensive experiments to analyze Geminiio’s broad applicability to different datasets (ImageNet [4], CIFAR-20 [21], and FER [6]), different neural architectures (ResNet [16], MobileNet [18], EfficientNet [28], and ViT [5]), and different FL scenarios (FedSGD and FedAvg). CIFAR-20 is equivalent to CIFAR-100 but uses 20 super-classes as labels. It provides ground truths to evaluate Geminiio’s task-agnostic targeted retrieval quantitatively.

By default, we consider an FL system that trains a ResNet34 model using FedSGD as the protocol with a batch size of 64. For Geminiio, we use the pretrained CLIP [25] to guide the optimization. The gradients are consumed by InvertingGrad [14] to reconstruct private samples. Detailed setup and the source code are provided in the supplementary materials to facilitate further research and reproducibility.

**Outline.** With additional analysis provided in the supplementary material, we would like to deliver three messages via the empirical studies in this section:

- The attacker can freely describe the data valuable to them and “query” the victim’s dataset for targeted reconstruction from a large batch of data. (Section 4.1)
- Geminiio serves as a plugin to existing reconstruction optimization methods and is broadly applicable, even with limited access to auxiliary data. (Section 4.2)
- Geminiio has a high survivability under various FL and defense scenarios. (Section 4.3)

### 4.1. Task-agnostic, Targeted Reconstruction

**Qualitative Analysis.** To showcase Geminiio’s targeted retrieval, Figure 4 provides two example batches of the victim’s private data (top) from different datasets and the corresponding reconstructed images (bottom) for three cases: reconstruction with the vanilla GIA (1st row) and reconstruction with Geminiio given two different queries (2nd and 3rd rows). First, while the vanilla GIA cannot produce recognizable images due to its failure to handle a large batch, Geminiio narrows the reconstruction scope to the data samples that matter most and successfully recovers them with high fidelity. Second, the recovered images match the attacker-provided queries. For instance, a curious attacker may submit a query “Any weapon?” to understand whether the client is, e.g., a weapon enthusiast. Among the 64 images on ImageNet (Figure 4a), only the first three contain weapons and are all successfully reconstructed. Similarly, considering the query “Person with beard and glasses,” while the first five images on FER (Figure 4b) contain a person wearing glasses, only the first two are reconstructed, as the rest do not have a beard. Third, queries can be irrelevant to the ML task. FER classifies facial images into one of the seven emotion expressions (e.g., happy, sad). Even though our example queries describe the appearance of individuals, the targeted reconstructions are successful.

**Comparison with Existing Methods.** Geminiio’s targeted reconstruction is unique and not achievable by existing methods. Figure 5 shows another batch of private data with 32 images of the class “Sombrero.” Imagine that an attacker wants to recover images that contain human faces as a privacy-intrusive example. As shown in the first col-

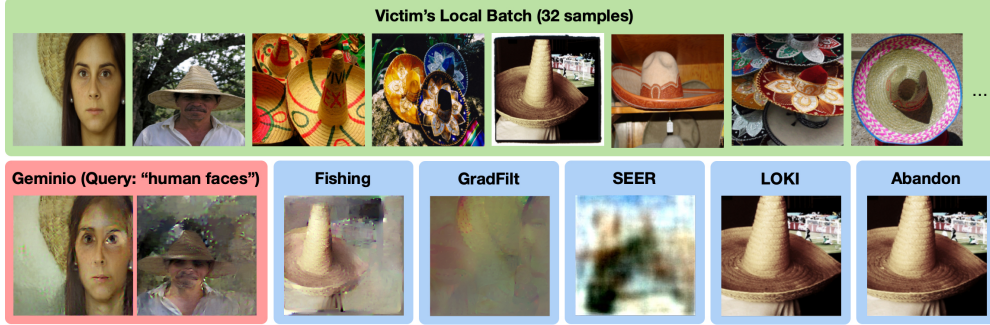


Figure 5. Gemino is the only method that achieves task-agnostic, instance-level targeted reconstruction. Other approaches can only be class-level (Fishing and GradFilt) or can only consider semantic-irrelevant conditions (SEER, LOKI, and Abandon).

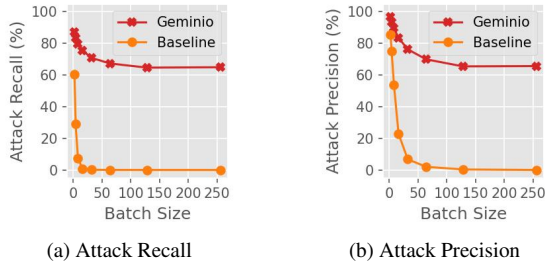


Figure 6. Gemino remains effective even when the batch size used by the victim is large (e.g., 256). In comparison, the baseline method is virtually useless when the batch size is larger than 8.

umn (2nd row), Gemino successfully recovers the first two images in the victim’s private data (1st row). The reconstructed images clearly reveal the facial features of the people with whom the client may interact. In contrast, other methods attempting to narrow the reconstruction scope cannot achieve the same goal. Fishing [35] can only return one random sample of a given class; GradFilt [44] returns all samples of a given class; SEER [12], LOKI [48], and Abandon [1] can only be random or specify semantic-irrelevant conditions (e.g., the brightness level). It is worth emphasizing that Gemino is instance-level. The matched data samples can belong to different classes. It is the only solution that achieves such a fine granularity.

**Quantitative Analysis.** Gemino can pinpoint and reconstruct valuable data samples from a large batch. Figure 6 reports the attack recall and precision on CIFAR-20 over different batch sizes used by the victim. Aligned with evaluating an information retrieval system, the attack recall indicates the percentage of data samples matching the query being retrieved (recovered), while the attack precision refers to the percentage of recovered data samples that indeed match the query. We split the entire training set of CIFAR-20 (50,000 images) into batches. For each of the 100 subclasses in the dataset, we use its name as the query to attack all batches, measure the attack recall and precision, and report their average across 100 subclasses. Following Fishing, we consider a data sample successfully reconstructed if

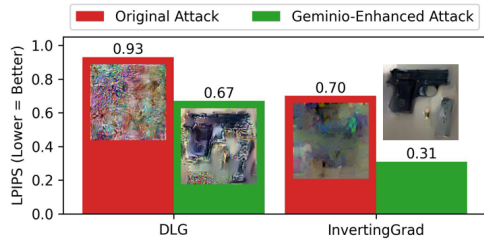


Figure 7. Gemino complements existing reconstruction optimization methods, turns them into targeted attacks, and improves their reconstruction quality.

its output-layer gradients dominate the batch-averaged gradients with a cosine similarity of at least 0.90. Figure 6 shows that Gemino remains effective even when the victim uses a large batch size, such as 256, with an attack recall of 64.96% and precision of 65.67%. Note that the malicious model was trained with a batch size of 64. We also compare Gemino with the baseline approach that uses a VLM to find data samples in the auxiliary dataset that match the query and poison their labels to increase their loss and gradients. As shown in Figure 6 (orange), it cannot provide a meaningful attack unless the batch size is extremely small (e.g., 2). This baseline demonstrates the effectiveness of Gemino in reshaping the loss surface.

## 4.2. Serving as a Plugin with Broad Applicability

**Complementary to Reconstruction Optimization.** Gemino can turn existing reconstruction optimization methods into targeted attacks. In addition to InvertingGrad (the default), we use DLG [49] to reconstruct the victim’s local batch in Figure 4a using the query “Any weapon?”. Figure 7 compares the two reconstruction techniques with and without Gemino’s enhancement. We use the standard metric, LPIPS [42], to understand how well the reconstructed images match the ground truths. A lower score means a higher reconstruction quality. We also provide the reconstructed images closest to the handgun (i.e., the 1st image in the batch) as a visual reference. We can observe that Gemino-

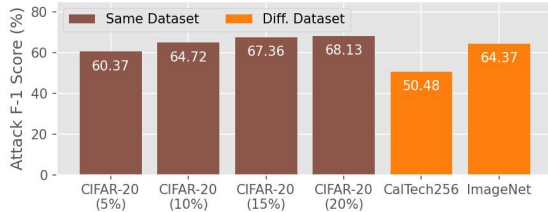


Figure 8. A small number of samples from the same dataset or a different dataset can already drive Geminio.

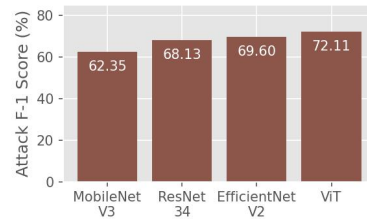
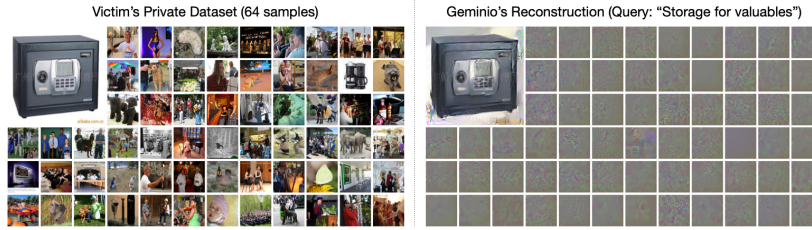
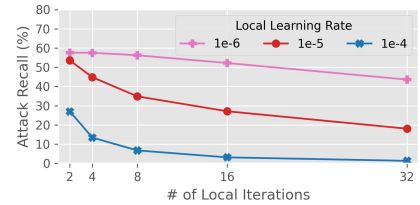


Figure 9. Geminio can attack any neural architectures out of the box without modifying them.



(a) FedAvg - Qualitative



(b) FedAvg - Quantitative

Figure 10. By coordinating the learning rate used by the FL clients, Geminio can be launched under the FedAvg protocol. Even if the client uses a batch size of 8 and runs one entire epoch of training before submitting gradients to the FL server, Geminio can still recover the image of a safe that matches the query.

enhanced attacks are consistently much better than their vanilla counterparts, which cannot recover any recognizable images. An interesting observation is that while DLG is known to be incapable of recovering from large batches and high-resolution images ( $64 \times 64$  as reported in the original paper), it can recover the handgun image well with a resolution of  $224 \times 224$  from a large batch. We conjecture that the gradient amplification in Geminio increases their variance, which will make the gradient matching during the reconstruction easier. We observe this phenomenon even for a batch of just one image.

**Auxiliary Data.** Figure 8 reports the attack F-1 score on CIFAR-20 using different auxiliary datasets. Compared with the default setting with the number of data samples equivalent to 20% of the training dataset, using only 5% of it only leads to a small drop in attack F-1 score, from 68.13% to 60.37%. Alternatively, the attacker can also use a different dataset, such as ImageNet or Caltech256. Even though they are not for the same ML task, the attack F-1 score can still achieve 50.48% and 65.37%, respectively. These datasets are publicly available and can be a practical source of auxiliary data.

**Neural Architectures.** Geminio can attack any neural architecture out of the box. Unlike many targeted attacks that need to inject a malicious module into the architecture, Geminio only modifies the model parameters in a stealthy manner. We conduct experiments to understand how it performs when different architectures are used in the FL system. According to the attack F-1 score reported in Figure 9, we observe that while Geminio works well on different architectures, the effectiveness slightly differs. It is

more effective on ViT and EfficientNetV2 than ResNet34 and MobileNetV3. Interestingly, this particular order reflects the general capability of these models. Hence, we conjecture that for more capable neural architectures, their privacy leakage by Geminio will be more severe.

### 4.3. Resilience to FedAvg and Defenses

While resilience to defenses is not the primary goal for Geminio, we found it to be resistant to popular methods.

**Federated Averaging.** Geminio can survive under FedAvg. Consider a victim having 256 ImageNet images in the private dataset as shown in Figure 10a (left). The victim uses a batch size of 8 and runs one epoch of training before sending the model updates to the server for aggregation. We employ Geminio using a query “Storage for valuables” to simulate a scenario where the attacker wants to know how the client stores the valuables. As shown in Figure 10a (right), it successfully recovers the image of a safe with high fidelity, even detailed enough to identify the specifics of it. The key enabler is to assign a small learning rate to the FL client, which is often the responsibility of the FL server. Figure 10b reports the attack F-1 score with different learning rates assigned to the victim. More local epochs weaken the attack because each iteration modifies the model parameters and may wash out the malicious patterns introduced by Geminio. Setting a small learning rate (e.g.,  $1e-6$ ) can slow down the performance degradation effectively.

**Gradient Pruning.** A popular defense to prune gradients of small magnitudes. Figure 11 reports the reconstruction quality on CIFAR-20 with varying pruning ratios. We reconstruct 100 batches and measure the average LPIPS. We



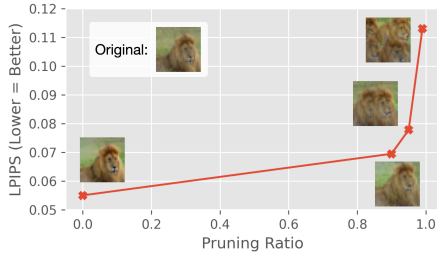


Figure 11. Gradient pruning cannot mitigate Geminio unless the pruning ratio is high, which can hinder the regular learning of FL.

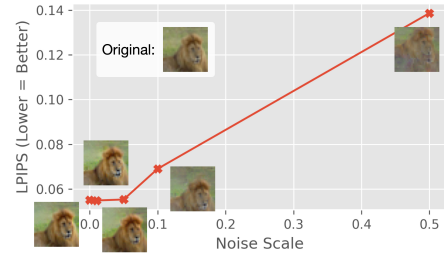


Figure 12. Adding Laplacian noise cannot prevent Geminio from retaining gradients of targeted samples, unless the degree of noise is significant, which can hinder the regular learning of FL.

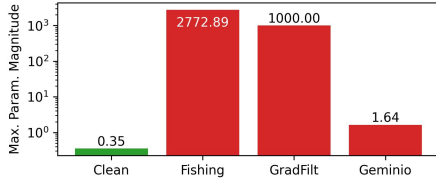


Figure 13. Fishing and GradFilt can be detected easily by model parameter inspection. Geminio does not need to set some model parameters to a large value, making it comparable with the clean model.

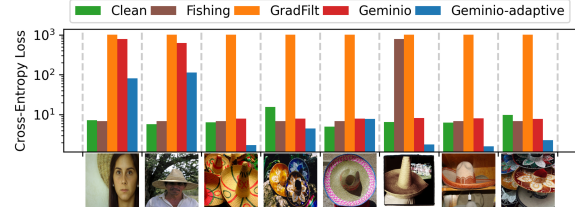


Figure 14. Targeted GIAs rely on amplifying the loss value of certain samples. An FL client who has access to it may detect such attacks.

observe that even with 95% of small gradients being set to zero, the perceptual quality of reconstructed images is still comparable to no defense. The reconstructed images become barely perceptible when 99% of the gradients are zeroed. In practice, such a setting is prohibited because it also removes the useful learning signals for training the ML model. Hence, gradient pruning cannot mitigate Geminio.

**Laplacian Noise.** Another popular defense is to add Laplacian noise to gradients. Figure 12 reports the reconstruction quality on CIFAR-20 with varying scales of Laplacian noise. Following [39], we use a per-layer noise injection. At each layer, we obtain its maximum gradient and scale it by a factor to be the standard deviation of the Laplacian noise with a zero mean for injection. A scale of 0.10 is already considered significant, but it barely affects the perceptual quality of reconstructed images. The reconstruction becomes severely affected when the noise scale is increased to 0.50, but it also washes out useful learning signals. Hence, injecting noise is not a viable defense against Geminio.

**Model Parameter Inspection.** As the client regularly receives model parameters from the server, it is natural to inspect whether they contain anomalies as a detection mechanism. Figure 13 reports the maximum magnitude of model parameters of a clean model and three poisoned models by Fishing, GradFilt, and Geminio. We observe that Fishing and GradFilt send a model with parameters deviating significantly from the clean one (2772.89 and 1000, respectively). In contrast, Geminio is only 1.64, close to the clean model (i.e., 0.35). Hence, setting a threshold may be able to detect Fishing and GradFilt, but not Geminio.

**Per-sample Loss Inspection.** The FL client may analyze the loss value per sample at each local training iteration.

We use the batch in Figure 5 and show the loss magnitude for each of the first 8 samples. All three attacks introduce a high loss value to the targeted samples. For Fishing, it successfully isolates the 6th sample in the batch, causing its loss to be significantly higher than the others. For GradFilt, since all samples in this batch are of the same target class (i.e., “sombbrero”), all samples have a magnified loss equal to 1000. For Geminio, the first two samples matching the attacker’s query (i.e., “human faces”) have amplified loss while the rest remains small. These are expected behaviors because targeted GIAs use the same principle: magnifying the gradients of desired samples to make them dominate the average gradients. While loss inspection seems promising, an advanced adversary could conduct an adaptive attack to suppress the loss values when training the malicious model (see Geminio-adaptive in Figure 5). Hence, more robust defenses need to be developed as future work.

## 5. Conclusions

We have introduced Geminio, a gradient inversion attack that harnesses the image-text association capabilities of pre-trained VLMs to enable language-guided targeted reconstructions. Our extensive experiments have yielded three key insights. First, Geminio enables the attacker to provide a natural language query to describe the data of value and reconstructs those matched samples from large data batches. Second, it serves as a plugin to enhance existing reconstruction optimization methods, broadly applicable to different neural architectures, auxiliary datasets, and FL protocols. Third, existing defenses are insufficient to mitigate Geminio. An advanced attacker can adapt Gemi-



nio to harden loss inspection. We believe that Geminio will inspire further research into the new threats posed by recent advancements in natural language processing, as they can be exploited as a “communication” interface for the adversary to express their goals and launch more flexible attacks.

## References

- [1] Franziska Boenisch, Adam Dziedzic, Roei Schuster, Ali Shahin Shamsabadi, Ilia Shumailov, and Nicolas Papernot. When the curious abandon honesty: Federated learning is not private. In *2023 IEEE 8th European Symposium on Security and Privacy (EuroS&P)*, pages 175–199. IEEE, 2023. 3, 6
- [2] Franziska Boenisch, Adam Dziedzic, Roei Schuster, Ali Shahin Shamsabadi, Ilia Shumailov, and Nicolas Papernot. Reconstructing individual data points in federated learning hardened with differential privacy and secure aggregation. In *2023 IEEE 8th European Symposium on Security and Privacy (EuroS&P)*, pages 241–257. IEEE, 2023. 3
- [3] Ka-Ho Chow, Wenqi Wei, and Lei Yu. Imperio: Language-guided backdoor attacks for arbitrary model control. In *International Joint Conference on Artificial Intelligence*, 2024. 3
- [4] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009. 3, 5, 12
- [5] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2021. 5
- [6] Dumitru, Ian Goodfellow, Will Cukierski, and Yoshua Bengio. Challenges in representation learning: Facial expression recognition challenge. <https://kaggle.com/competitions/challenges-in-representation-learning-facial-expression-recognition-challenge>, 2013. Kaggle. 5
- [7] Will Cukierski Yoshua Bengio Dumitru, Ian Goodfellow. Challenges in representation learning: Facial expression recognition challenge, 2013. 12
- [8] Haokun Fang and Quan Qian. Privacy preserving machine learning with homomorphic encryption and federated learning. *Future Internet*, 13(4):94, 2021. 3
- [9] Hao Fang, Bin Chen, Xuan Wang, Zhi Wang, and Shu-Tao Xia. Gifd: A generative gradient inversion method with feature domain optimization. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4967–4976, 2023. 2
- [10] Hao Fang, Yixiang Qiu, Hongyao Yu, Wenbo Yu, Jiawei Kong, Baoli Chong, Bin Chen, Xuan Wang, Shu-Tao Xia, and Ke Xu. Privacy leakage on dnns: A survey of model inversion attacks and defenses. *arXiv preprint arXiv:2402.04013*, 2024. 1, 2
- [11] Liam Fowl, Jonas Geiping, Wojtek Czaja, Micah Goldblum, and Tom Goldstein. Robbing the fed: Directly obtaining private data in federated learning with modified models. *arXiv preprint arXiv:2110.13057*, 2021. 2, 3
- [12] Kostadin Garov, Dimitar I Dimitrov, Nikola Jovanović, and Martin Vechev. Hiding in plain sight: Disguising data stealing attacks in federated learning. *arXiv preprint arXiv:2306.03013*, 2023. 3, 6, 13
- [13] Gartner. What’s new in the 2023 gartner hype cycle for emerging technologies. <https://www.gartner.com/en/articles/what-s-new-in-the-2023-gartner-hype-cycle-for-emerging-technologies>. 1
- [14] Jonas Geiping, Hartmut Bauermeister, Hannah Dröge, and Michael Moeller. Inverting gradients-how easy is it to break privacy in federated learning? *Advances in neural information processing systems*, 33:16937–16947, 2020. 1, 2, 3, 5, 13
- [15] Ali Hatamizadeh, Hongxu Yin, Holger R. Roth, Wenqi Li, Jan Kautz, Daguang Xu, and Pavlo Molchanov. Gradvit: Gradient inversion of vision transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10021–10030, 2022. 3
- [16] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 5, 13
- [17] Geoffrey E Hinton, Nitish Srivastava, Alex Krizhevsky, Ilya Sutskever, and Ruslan R Salakhutdinov. Improving neural networks by preventing co-adaptation of feature detectors. *arXiv preprint arXiv:1207.0580*, 2012. 12
- [18] Andrew G. Howard, Mark Sandler, Grace Chu, Liang-Chieh Chen, Bo Chen, Mingxing Tan, Weijun Wang, Yukun Zhu, Ruoming Pang, Vijay Vasudevan, Quoc V. Le, and Hartwig Adam. Searching for mobilenetv3. *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 1314–1324, 2019. 5
- [19] Yangsibo Huang, Samyak Gupta, Zhao Song, Kai Li, and Sanjeev Arora. Evaluating gradient inversion attacks and defenses in federated learning. *Advances in neural information processing systems*, 34:7232–7241, 2021. 1
- [20] Yuxin Wen Jonas Geiping, Liam Fowl. Breaching - a framework for attacks against privacy in federated learning (<https://github.com/JonasGeiping/breaching>), 2022. 11, 13
- [21] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009. 5
- [22] Quande Liu, Cheng Chen, Jing Qin, Qi Dou, and Pheng-Ann Heng. Feddg: Federated domain generalization on medical image segmentation via episodic learning in continuous frequency space. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1013–1023, 2021. 1
- [23] Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas. Communication-efficient learning of deep networks from decentralized data. In *Artificial intelligence and statistics*, pages 1273–1282. PMLR, 2017. 1, 2, 3

- [24] Pretom Roy Ovi, Emon Dey, Nirmalya Roy, and Aryya Gangopadhyay. Mixed quantization enabled federated learning to tackle gradient inversion attacks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5046–5054, 2023. 2
- [25] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021. 2, 4, 5, 13
- [26] Yichuan Shi, Olivera Kotevska, Viktor Reshniak, Abhishek Singh, and Ramesh Raskar. Dealing doubt: Unveiling threat models in gradient inversion attacks under federated learning, a survey and taxonomy. *arXiv preprint arXiv:2405.10376*, 2024. 2
- [27] H Takahashi, J Liu, Y Liu, and Y Liu. Breaching fedmd: Image recovery via paired-logits inversion attack. In *The IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023. 2
- [28] Mingxing Tan and Quoc V. Le. Efficientnetv2: Smaller models and faster training. In *International Conference on Machine Learning*, 2021. 5
- [29] Aidmar Wainakh, Fabrizio Ventola, Till Müßig, Jens Keim, Carlos Garcia Cordero, Ephraim Zimmer, Tim Grube, Kristian Kersting, and Max Mühlhäuser. User label leakage from gradients in federated learning. In *Privacy Enhancing Technologies Symposium*, 2022. 11
- [30] Kang Wei, Jun Li, Ming Ding, Chuan Ma, Howard H Yang, Farhad Farokhi, Shi Jin, Tony QS Quek, and H Vincent Poor. Federated learning with differential privacy: Algorithms and performance analysis. *IEEE transactions on information forensics and security*, 15:3454–3469, 2020. 3
- [31] Wenqi Wei, Ling Liu, Margaret Loper, Ka-Ho Chow, Mehmet Emre Gursoy, Stacey Truex, and Yanzhao Wu. A framework for evaluating client privacy leakages in federated learning. In *Computer Security—ESORICS 2020: 25th European Symposium on Research in Computer Security, ESORICS 2020, Guildford, UK, September 14–18, 2020, Proceedings, Part I 25*, pages 545–566. Springer, 2020. 1
- [32] Wenqi Wei, Ling Liu, Yanzhao Wu, Gong Su, and Arun Iyengar. Gradient-leakage resilient federated learning. In *2021 IEEE 41st International Conference on Distributed Computing Systems (ICDCS)*, pages 797–807. IEEE, 2021. 3
- [33] Wenqi Wei, Ka-Ho Chow, Fatih Ilhan, Yanzhao Wu, and Ling Liu. Model cloaking against gradient leakage. In *2023 IEEE International Conference on Data Mining (ICDM)*, pages 1403–1408. IEEE, 2023. 3
- [34] Jie Wen, Zhixia Zhang, Yang Lan, Zhihua Cui, Jianghui Cai, and Wensheng Zhang. A survey on federated learning: challenges and applications. *International Journal of Machine Learning and Cybernetics*, 14(2):513–535, 2023. 1
- [35] Yuxin Wen, Jonas A Geiping, Liam Fowl, Micah Goldblum, and Tom Goldstein. Fishing for user data in large-batch federated learning via gradient magnification. In *International Conference on Machine Learning*, pages 23668–23684. PMLR, 2022. 2, 3, 6, 12, 13
- [36] Ruihan Wu, Xiangyu Chen, Chuan Guo, and Kilian Q Weinberger. Learning to invert: Simple adaptive attacks for gradient inversion in federated learning, 2023. 2
- [37] Xiaoyu Wu, Yang Hua, Chumeng Liang, Jiaru Zhang, Hao Wang, Tao Song, and Haibing Guan. Cgi-dm: Digital copyright authentication for diffusion models via contrasting gradient inversion. *2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10812–10821, 2024.
- [38] Ran He Yanbo Wang, Jian Liang. Towards eliminating hard label constraints in gradient inversion attacks. In *The Twelfth International Conference on Learning Representations*, 2024. 2
- [39] Zipeng Ye, Wenjian Luo, Qi Zhou, and Yubo Tang. High-fidelity gradient inversion in distributed learning. In *AAAI Conference on Artificial Intelligence (AAAI)*, 2024. 2, 8
- [40] Hongxu Yin, Arun Mallya, Arash Vahdat, Jose M Alvarez, Jan Kautz, and Pavlo Molchanov. See through gradients: Image batch recovery via gradinversion. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 2, 11, 13
- [41] Chengliang Zhang, Suyi Li, Junzhe Xia, Wei Wang, Feng Yan, and Yang Liu. {BatchCrypt}: Efficient homomorphic encryption for {Cross-Silo} federated learning. In *2020 USENIX annual technical conference (USENIX ATC 20)*, pages 493–506, 2020. 3
- [42] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *CVPR*, 2018. 6
- [43] Rui Zhang, Song Guo, Junxiao Wang, Xin Xie, and Dacheng Tao. A survey on gradient inversion: Attacks, defenses and future directions. *arXiv preprint arXiv:2206.07284*, 2022. 1
- [44] Rui Zhang, Song Guo, and Ping Li. Gradfilt: Class-wise targeted data reconstruction from gradients in federated learning. In *Companion Proceedings of the ACM on Web Conference 2024*, pages 698–701, 2024. 3, 6, 13
- [45] Tuo Zhang, Lei Gao, Chaoyang He, Mi Zhang, Bhaskar Krishnamachari, and A Salman Avestimehr. Federated learning for the internet of things: Applications, challenges, and opportunities. *IEEE Internet of Things Magazine*, 5(1):24–29, 2022. 1
- [46] Zhiqiu Zhang, Zhu Tianqing, Wei Ren, Ping Xiong, and Kim-Kwang Raymond Choo. Preserving data privacy in federated learning through large gradient pruning. *Computers & Security*, 125:103039, 2023. 3
- [47] Bo Zhao, Konda Reddy Mopuri, and Hakan Bilen. idlg: Improved deep leakage from gradients. *arXiv preprint arXiv:2001.02610*, 2020. 2, 11
- [48] Joshua C Zhao, Atul Sharma, Ahmed Roushdy Elkordy, Yahya H Ezzeldin, Salman Avestimehr, and Saurabh Bagchi. Loki: Large-scale data reconstruction attack against federated learning through model manipulation. In *2024 IEEE Symposium on Security and Privacy (SP)*, pages 1287–1305. IEEE, 2024. 3, 6, 13
- [49] Ligeng Zhu, Zhijian Liu, and Song Han. Deep leakage from gradients. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2019. 2, 6, 13

# Geminio: Language-Guided Gradient Inversion Attacks in Federated Learning

## Supplementary Material

### Outline

This document provides additional details to support our main paper. It is organized as follows:

- Section A: Geminio Strengthens Label Inference Attacks
- Section B: Geminio Works Under Homomorphic Encryption
- Section C: Geminio Supports Different Local Batch Sizes
- Section D: Experiment Setup
- Section E: Additional Visual Examples

### A. Geminio Strengthens Label Inference Attacks

Label inference is a prerequisite for gradient inversion, with various attack methods being proposed [29, 40, 47]. Surprisingly, Geminio is not just compatible with them but also boosts their accuracy. We use five label inference attacks provided by the `breaching` library [20] and compare the original attack with the Geminio-enhanced one. Since our problem setting focuses on targeted reconstructions, we only need to make sure the class labels with matched samples in the local batch are inferred. The success or failure of inferring other class labels is unimportant because their gradients are small and negligible in the gradient matching (reconstruction optimization) process. Figure 15 reports the results measured on CIFAR-20. This dataset provides ground truths for conducting such quantitative studies. When gradients submitted by the victim are generated based on the Geminio-poisoned malicious model, all label inference attacks are consistently improved. This phenomenon can be explained by our observation in Figure 16 that the class labels containing matched samples in the local batch have their gradients amplified. Since those attacks share the same principle to examine the gradient magnitude of different classes, Geminio facilitates this label inference process.

### B. Geminio Works Under Homomorphic Encryption

Our threat model considers an active attacker who is the FL server. The attacker can execute Geminio under homomorphic encryption by controlling only one client. As the malicious client can obtain the victim’s gradients in plain text, Geminio can be run on the client side and perform identically to FL without homomorphic encryption. Figure 17 provides reconstruction results with “luxury watches” as the attacker’s query. The two watches can be retrieved from the victim’s gradients, leading to a high-quality reconstruction

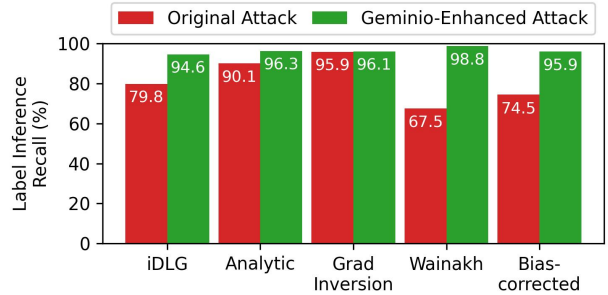


Figure 15. Geminio consistently improves five label inference attacks. Given an attacker’s query, it leads to a high success rate in inferring class labels containing matched samples in the local batch.

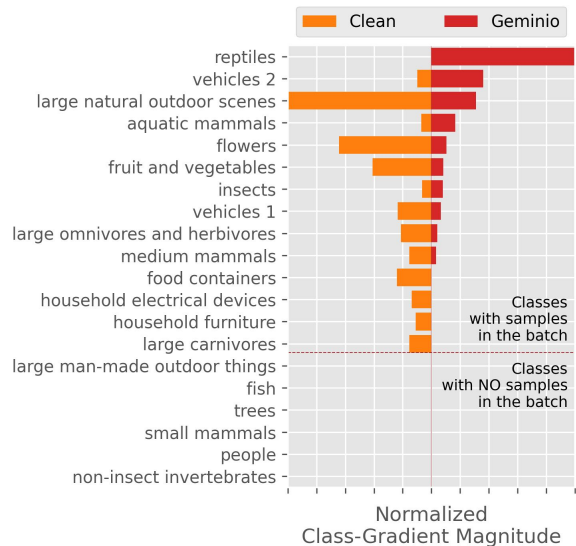


Figure 16. Label inference attacks examine the per-class gradient magnitude. Compared with a clean model, Geminio, with the query “dinosaur,” will amplify the gradients of the class(es) to which the matched samples belong (the class “reptiles” in this example). This facilitates the label inference process.

where we can even read the brand for the first image to be Rolex.

### C. Geminio Supports Different Local Batch Sizes

During Geminio’s optimization, minibatch training needs to be conducted but this training batch size does not need to match the local batch size used by the client. Figure 18 re-

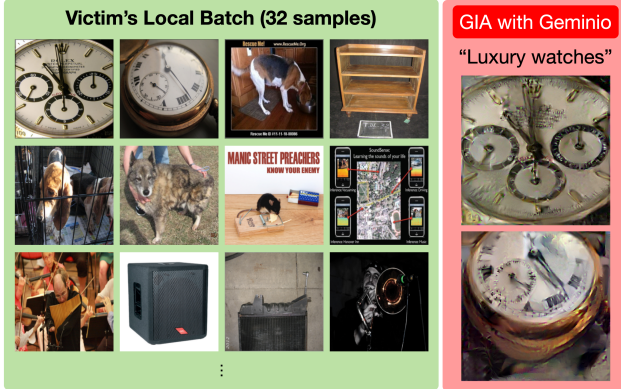


Figure 17. By controlling one FL client, Geminio can retrieve targeted private samples under FL with homomorphic encryption.

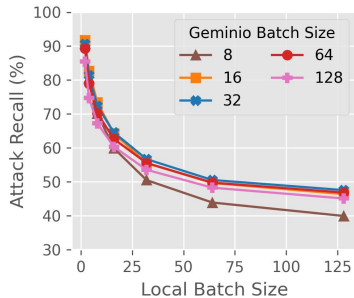


Figure 18. The training batch size used by Geminio to poison the model is irrelevant to the local batch size to be used by the victim.

ports the attack recall with varying local batch sizes. We repeat the experiment using different training batch sizes for Geminio to optimize the malicious global model. We observe that their targeted retrieval performances are similar, with the smallest batch size of 8 being slightly worse. For instance, when Geminio uses a batch size of 64 for its optimization, the malicious global model can be sent to clients with any local batch size, which may or may not be controlled by the server (e.g., depending on the computing resources of the client device).

## D. Experiment Setup

Our experiments cover a wide range of datasets, ML models, and FL scenarios to analyze Geminio’s properties. Here, we describe the default experiment setup.

### D.1. Datasets

We conduct experiments on three datasets: ImageNet [4], CIFAR-20 [17], and Facial Expression Recognition (FER) [7]. By default, visual examples are based on ImageNet.

The scenario of fine-grained targeted retrieval by Gem-

Table 1. The superclasses and their subclasses in CIFAR-100. We create a benchmark dataset, CIFAR-20, that uses the 20 superclasses for the classification problem and the 100 subclass names as queries. This design gives us ground truths for the instance-level retrieval.

Superclass (20)	Subclasses (100)
aquatic mammals	beaver, dolphin, otter, seal, whale
fish	aquarium fish, flatfish, ray, shark, trout
flowers	orchids, poppies, roses, sunflowers, tulips
food containers	bottles, bowls, cans, cups, plates
... (16 more rows) ...	

inio can be imagined as an attacker writing a “query” to search for relevant records in the victim’s private database. Quantitative evaluation requires two ingredients: (i) a benchmark dataset with ground truths and (ii) a set of indicative performance metrics.

**Benchmark: CIFAR-20** The benchmark dataset should include a set of queries, each is a textual description and associated with a list of relevant images. Then, we can randomly sample a local batch from the dataset, use Geminio to reconstruct images given different queries, and measure how many relevant images are successfully reconstructed. This process repeats for a number of random local batches until, e.g., all training images are processed. To showcase instance-level retrieval better, the queries should not be the class names of the classification problem. Based on these requirements, we created a variant of CIFAR-100 and named it CIFAR-20. Each image in CIFAR-100 is associated with two official labels, a subclass and a superclass (see Table 1 for four superclasses and their subclasses). We use the 20 superclasses for the classification problem and the 100 subclasses as queries. With this design, we can easily obtain images in the local batch that should be retrieved for a given query (i.e., a subclass name).

**Metrics: Attack Recall and Precision** We follow Fishing’s approach [35] to determine whether an image in a local batch dominates and will be reconstructed. In particular, if the gradients produced by an image have a cosine similarity with the average gradients over a threshold, it is considered a reconstructed sample. While Fishing uses 0.95 as the threshold, we found that this is overly restrictive. Instead, we use 0.90. Note that we observe multiple examples where targeted reconstruction succeeds even if the cosine similarity is below 0.90. Our choice (i.e., 0.90) is still conservative. A more principled approach is considered as our future work. Based on this thresholding, we can measure the percentage of targeted images being reconstructed (i.e., Attack Recall) and, among all reconstructed images, the percentage of them being the actual targeted images (i.e., Attack



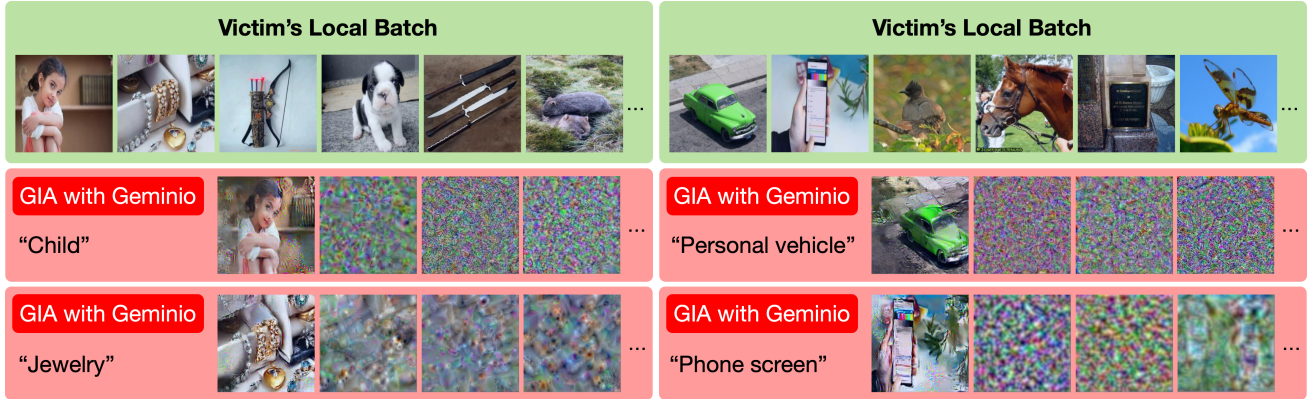


Figure 19. Geminio considers the attacker-specified query to pinpoint and reconstruct the matched samples in the private local batch.

Precision).

## D.2. FL Configuration

The FL system aims to train a ResNet34 [16] model. Following existing works [12, 14, 35, 40, 44, 48, 49], we use FedSGD to be the default protocol. The FL client receives a model from the server, updates it with a batch of private samples, and returns the gradients to the server, which is malicious, and attempts to reconstruct private samples from it.

## D.3. Attack Configuration

For Geminio, we use CLIP [25] with the ViT-L/14 Transformer architecture as the pretrained VLM<sup>2</sup> to process auxiliary data, which comes from the respective validation set. Geminio poisons the model with a training batch size 64 using Adam as the optimizer. For gradient inversion, we use InvertingGrad [14].

## D.4. Computing Environment

All experiments are conducted on a server with Intel® Xeon® Gold 6526Y CPU, 64GB RAM, and two NVIDIA RTX 5880 Ada Lovelace GPUs.

## D.5. Implementation

Geminio is written in PyTorch and can be easily integrated into existing GIAs. Our implementation uses breaching [20], a collection of GIAs, to demonstrate such a plug-and-play feature. We first extracted image features from auxiliary data, which took about 7 minutes for ImageNet. Given a query from the attacker, Geminio can use those pre-generated image features to poison the model in less than 8 minutes.

<sup>2</sup><https://huggingface.co/openai/clip-vit-large-patch14>

## E. Additional Visual Examples

We provide additional visual examples in Figure 19. Geminio can prioritize reconstruction to recover those samples that match the attacker-provided queries. For the first local batch (left), the query “child” leads to the reconstruction of the 1st image, while the query “jewelry” to the same batch recovers the necklace (i.e., the 2nd image). Similarly, for the second local batch (right), the green car (i.e., the 1st image) will be recovered if the attacker provides “personal vehicle” as the query. However, if the query is “phone screen” instead, the same reconstruction optimization will recover the 2nd image in the batch automatically.