

JANUS: DUAL-SERVER MULTI-ROUND SECURE AGGREGATION WITH VERIFIABILITY FOR FEDERATED LEARNING

Anonymous authors

Paper under double-blind review

ABSTRACT

Secure Aggregation (SA) in federated learning is essential for preserving user privacy by ensuring that model updates are masked or encrypted and remain inaccessible to servers. Although the advanced protocol Flamingo (S&P’23) has made significant strides with its multi-round aggregation and optimized communication, it still faces several critical challenges: (i) *Dynamic User Participation*, where Flamingo struggles with scalability due to the complex setups required when users join or leave the training process; (ii) *Model Inconsistency Attacks (MIA)*, where a malicious server could infer sensitive data, which poses severe privacy risks; and (iii) *Verifiability*, as most schemes lack an efficient mechanism for clients to verify the correctness of server-side aggregation, potentially allowing inaccuracies or malicious actions. We introduce Janus, a generic privacy-enhanced multi-round SA scheme through a dual-server architecture. A new user can participate in training by simply obtaining the servers’ public keys for aggregation, eliminating the need for complex communication graphs. Our dual-server model separates aggregation tasks, ensuring that neither server can successfully launch a MIA without controlling at least $n - 1$ clients. Additionally, we propose a new cryptographic primitive, *Separable Homomorphic Commitment*, integrated with our dual-server approach to ensure the verifiability of aggregation results. Extensive experiments across various models and datasets show that Janus significantly boosts security while enhancing efficiency. It reduces per-client communication and computation overhead from logarithmic to constant scale compared to state-of-the-art methods, with almost no compromise in model accuracy.

1 INTRODUCTION

Traditional machine learning relies on centralized training, where the entire dataset is stored in a single central location and directly accessible by the server. However, users are often reluctant to share data, especially if it involves sensitive information like medical records, photos, or trade secrets. Federated Learning (FL) was proposed to protect user privacy and enable the model training (McMahan et al., 2017). FL is a distributed machine learning framework that uses privacy-preserving cryptographic techniques, which allows participants to collaborate on model training without disclosing their private data. Unfortunately, it has been shown that an adversary can invert a single model update from a target user, thereby revealing a great deal of sensitive information about its local dataset (Hitaj et al., 2017; Nasr et al., 2019; Zhu et al., 2019).

To protect the user gradient information, Secure Aggregation (SA (Bonawitz et al., 2017)) is introduced to enhance the security of FL, which can prevent server access to individual model updates. SA is considered as one of the most robust defenses against gradient inversion and related inference attacks (Huang et al., 2021). Most of the current SA schemes rely on the double-mask, which involves heavy secret sharing, especially as the number of participants grows, requiring two clients to negotiate the key and engage in frequent communication. The advanced SA protocol (BBSA (Bell et al., 2020)) manages the aggregation with thousands of clients and high-dimensional input vectors while tolerating device drops during execution. However, these schemes select a subset of clients and enables aggregation for only one round. Although it is possible to run the protocol multiple times to complete multi-round of aggregation, the *Setup* phase must be re-run for each round to

054 maintain privacy, requiring server interaction with all clients during each step. This results in signif-
055 icant communication overhead and reduced efficiency.

056 Recently, the state-of-the-art Flamingo (Ma et al., 2023) eliminates the need for re-setup in each
057 round, which supports multi-round SA based on the BBSA. It also optimizes the communication
058 graph to improve the system performance, with introducing a set of decryptors to handle part of
059 the computation. While Flamingo marks significant progress, it has limitations in handling dynamic
060 user participation, resisting Model Inconsistency Attacks (MIA) (Pasquini et al., 2022), and ensuring
061 correct server-side aggregation. When users join or leave, the complex setups needed for Flamingo
062 lessen its practicality. The server can still exploit the MIA to infer sensitive information, and clients
063 have no way to verify if the server correctly performed the aggregation or omitted user data.

064 These vulnerabilities stem from the reliance on a single server, which is common in existing schemes
065 due to its simplicity. A single server inherently knows the aggregated results, providing an opportu-
066 nity for a malicious server to compromise the privacy by bypassing the SA protocol (Pasquini et al.,
067 2022). Specifically, the server distributes carefully crafted parameters to non-target users, which
068 can trigger the *dying-ReLU* effect that causes non-target users to generate zero gradients during ag-
069 gregation. As a result, the aggregated gradient effectively reveals the target user’s gradient. This
070 attack affects not only double-mask schemes but all schemes where the server can access the ag-
071 gregation results. While cryptographic signatures could prevent this by allowing users to verify the
072 consistency of received parameters. This approach involves heavy computation and requires users
073 to negotiate the consistency of the received information, which places a large burden on the system.

074 Our research indicates that preventing MIAs necessitates restricting the server’s access to the final
075 aggregation results. To achieve this, we propose a dual-server architecture: one server handles the
076 collection and aggregation of masked gradients, while the other manages the aggregation of masks.
077 If the servers do not collude, neither can access the final aggregated results. This assumption is
078 feasible in many real-world scenarios. For example, banks, financial institutions, and healthcare
079 organizations, despite having different interests, are generally committed to protecting user privacy
080 and complying with regulations. They are motivated to collaborate for the benefit of users and avoid
081 collusion. In the Flamingo scheme, the decryptors can be also considered as one server, with a
082 second server forming the dual-server architecture. This approach ensures security while leveraging
083 the practical willingness of institutions to cooperate for SA. Additionally, numerous studies are
084 relevant to our work, with detailed discussion provided in Appendix A.

085 Another challenge is to ensure the correctness of aggregation, particularly in a dual-server archi-
086 tecture where either server could miscollect or misaggregate masked gradients or masks. An ag-
087 gregation server might prioritize speed over accuracy, performing fast but faulty computations to
088 save resources, which can lead to erroneous results. Since servers are often semi-trusted, they could
089 also deliberately mishandle some gradients or falsify aggregation results, misleading users about the
090 training results (Hahn et al., 2023). Current schemes face difficulties in ensuring efficient verifiabil-
091 ity, typically depending on resource-intensive techniques like homomorphic hashing or signatures.
092 Moreover, errors in aggregation could arise from malicious client submissions, yet current methods
093 fail to enforce strong client-side commitments. To address these challenges, our approach introduces
094 a new cryptographic primitive called separable homomorphic commitment (SHC), which ensures
095 both server-side integrity and client-side data accuracy in the dual-server setting. Homomorphism
096 and separability are two important properties of SHC. The two servers aggregate the different values
097 in the commitment separately. SHC can separates out the part of message and compares them with
the aggregated results, thus enabling the correctness of aggregation.

098 Our main contributions are summarized as follows.

- 099
100 • *Generic construction of dual-server SA with dynamic user participation for FL.* We propose
101 Janus, the first generic construction of SA based on dual-server, which can work well for
102 multiple round of aggregation without re-setup in FL. Our new design avoids heavy com-
103 munication graphs such as complete graphs and k -regular graphs. Additionally, Janus only
104 involves some lightweight components, thus it can avoid the need for time-consuming op-
105 erations such as secret sharing, which in turn dramatically improves the system efficiency.
106 It also enables dynamic user participation with only the servers’ public keys.
- 107 • *A new cryptographic primitive and enhanced privacy with verifiability.* Our primary contri-
bution is the conceptual development of a new cryptographic primitive, termed *Separable*

Homomorphic Commitment (SHC). By analyzing the algebraic properties of current commitment schemes, we identify a common blueprint that can be instantiated to provide novel verification methods for aggregation results. Furthermore, we introduce a dual-server architecture that leverages SHC to enhance both privacy and verifiability. This architecture ensures that aggregation results remain invisible to individual servers, making it impossible for a malicious server to bypass the SA. Consequently, our approach not only enhances resistance to malicious inference attacks but also incorporates verifiability, providing additional security advantages.

- *Implementation and evaluation.* We implemented an instantiation for Janus and evaluated it with similar classical schemes via extensive experiments on different models and datasets. The results show that Janus outperforms in terms of both computation and communication. It reduces per-client overhead from the logarithmic scale of current advanced methods to a constant scale. Table 1 demonstrates that Janus surpasses other state-of-the-art schemes in terms of security, efficiency, and functionality.

Table 1: Comparison of SA Constructions

Scheme	Input Privacy	Multi-round	Verifiability	Dynamic	Versatility	NS*	Efficiency†	MIA
SecAgg (Bonawitz et al., 2017)	✓	✗	✗	✗	✗	1	○	✗
BBSA (Bell et al., 2020)	✓	✗	✗	✗	✗	1	◐	✗
Flamingo (Ma et al., 2023)	✓	✓	✗	✗	✗	2†	◑	✗
Janus	✓	✓	✓	✓	✓	2	●	✓

✓ Support, ✗ No support. Versatility: A generic construction. * Number of servers. † The decryptors of this construction can be abstracted to a server. ‡ More black parts in the circle indicate better efficiency.

2 PRELIMINARIES

2.1 COMMITMENTS

Commitments (Pedersen, 1991) provide the cryptographic cornerstone for integrity and trust in various schemes. It enables participants to commit to values without compromising the confidentiality of the information. Typically, a non-interactive secure commitment scheme consists of the following three algorithms:

1. $CSetup(1^\lambda) \rightarrow pp$. The system initialization algorithm takes as input a security parameter λ , and it outputs the public parameter pp for the commitment scheme.
2. $Commit(pp, v, r) \rightarrow c$. The commitment generation algorithm takes as input a message v from the message space \mathcal{M}_{pp} and a random number (blinder) r in the randomness space \mathcal{R}_{pp} , and it outputs the commitment c in the commitment space \mathcal{C}_{pp} .
3. $Reveal(pp, v, c, r) \rightarrow b$. The revealing commitment algorithm takes as input a message v , a commitment c and a blinder r . If it accepts then the output $b = 1$; otherwise, $b = 0$.

Normally, a secure commitment scheme must satisfy the following three properties.

- **Completeness.** It ensures that if both the committer and the verifier follow the protocol correctly, the verifier will always accept the decommitment (Reveal).

$$\Pr \left(\begin{array}{l} CSetup(1^\lambda) \rightarrow pp; \\ Commit(pp, v, r) \rightarrow c : \\ Reveal(pp, v, c, r) = 1 \end{array} \right) = 1. \tag{1}$$

- **Hiding.** During the commitment phase, the verifier cannot infer the committed value from the commitment. It can ensure that the committed value remains confidential until it is revealed. For any v_1, v_2 of equal length, and any r , the following probability distributions are computationally indistinguishable.

$$\{Commit(pp, v_1, r) \rightarrow c_1\} \stackrel{c}{\approx} \{Commit(pp, v_2, r) \rightarrow c_2\}. \tag{2}$$

- **Binding.** After the commitment is made, the committer cannot change the committed value. It can prevent the committer from cheating by ensuring the immutability of the commitment. There exists a negligible function $\text{negl}(\lambda)$ such that for all non-uniform Probabilistic Polynomial Time (PPT) adversaries \mathcal{A} ,

$$\Pr \left(\begin{array}{l} \text{CSetup}(1^\lambda) \rightarrow pp; \\ \mathcal{A}(pp) \rightarrow (c, r, v_1, v_2) : \\ \text{Reveal}(pp, c, v_1, r) = 1 \wedge \\ \text{Reveal}(pp, c, v_2, r) = 1 \wedge \\ v_1 \neq v_2 \end{array} \right) \leq \text{negl}(\lambda). \quad (3)$$

2.2 MASKING-BASED SECURE AGGREGATION

The One-Time Pad (OTP) is a type of classical encryption which can be perfect secrecy (Katz & Lindell, 2014). Specifically, OTP can encrypt information using either addition or multiplication. Participants can mask their updates to preserve privacy in FL. A formal OTP scheme usually contains the following two algorithms.

1. **Masking** $(x, k) \rightarrow \hat{x}$. The masking algorithm takes as input a secret message x and a private key k , and it outputs the encryption result \hat{x} .
2. **UnMasking** $(\hat{x}, k) \rightarrow x$. The unmasking algorithm takes as input an encrypted message \hat{x} and a private key k , and it outputs the plain message x .

Users can apply masking to updates via OTP before uploading to the central servers for aggregation. SA is designed not only to effectively prevent centralized servers from snooping on individual models, but also to defend against attacks from malicious participants and ensure the robustness of the entire FL system. Researchers have proposed several variants of SA to address different threat models and system requirements. We focus on masking-based aggregation schemes. Specifically, there is a set of users \mathcal{U} where $u_i \in \mathcal{U}$ has a private update x_i in FL. In masking-based SA, each u_i adds a pair-wise additive mask to its private update x_i to get the masked vector y_i as follows:

$$y_i = x_i + \sum_{u_j \in \mathcal{U}: i < j} \text{PRG}(s_{i,j}) - \sum_{u_j \in \mathcal{U}: i > j} \text{PRG}(s_{j,i}), \quad (4)$$

where the pseudorandom generator (PRG) can randomly generate a sequence numbers based on the random seed $s_{i,j}$. Note that the masks will be removed when all masked input updates y_i are summed, resulting in

$$\sum_{u_i \in \mathcal{U}} y_i = \sum_{u_i \in \mathcal{U}} \left(x_i + \sum_{i < j} \text{PRG}(s_{i,j}) - \sum_{i > j} \text{PRG}(s_{j,i}) \right) = \sum_{u_i \in \mathcal{U}} x_i. \quad (5)$$

In addition, in order to deal with dropped users during protocol execution, the Shamir secret sharing scheme (Shamir, 1979) is used to share seeds among users. The Diffie-Hellman (DH) key exchange protocol (Diffie & Hellman, 1976) is used to negotiate the seeds $s_{i,j}$ for each pair of users $(u_i, u_j) \in \mathcal{U}$. Note that for large-scale FL applications, the above scheme is not cost-effective. For a n -user FL system, it takes $\mathcal{O}(n^2)$ communication rounds to run the pairwise DH key exchange protocol.

2.3 MODEL INCONSISTENCY ATTACKS

A malicious server \mathcal{AS} intends to obtain private information about the model update of target user U_{tar} . It can distribute elaborately constructed parameters $\theta_{i,t}$ to the non-target users $\{U \setminus U_{tar}\}$ and then send normal parameters $\theta_{tar,t}$ to the target user, where \mathcal{U} denotes the set of all users. This can trigger the *dying-ReLU* (Lu et al., 2019), where the dead layer cannot generate any gradient. Therefore, the non-target user ends up generating tampered model updates $\Delta_{D_{i,t}}^{\theta_{i,t}}$, where the $D_{i,t}$ is the local data of U_i . While the parameters of U_{tar} are real thus generating a right update $\Delta_{D_{tar,t}}^{\theta_{tar,t}}$ on its local data $D_{tar,t}$ in round t . These tampered model updates can enable \mathcal{AS} to obtain the model

updates $\Delta_{\mathcal{D}_{tar,t}}^{\theta_{tar,t}}$ of \hat{U}_{tar} in plaintext. Specifically, the final result of secure aggregation is as follows,

$$\begin{aligned} \mathcal{AS}^{SA}(\Delta_{\mathcal{D}_{1,t}}^{\theta_{1,t}}, \dots, \Delta_{\mathcal{D}_{i-1,t}}^{\theta_{i-1,t}}, \Delta_{\mathcal{D}_{tar,t}}^{\theta_{tar,t}}, \Delta_{\mathcal{D}_{i+1,t}}^{\theta_{i+1,t}}, \dots, \Delta_{\mathcal{D}_{n,t}}^{\theta_{n,t}}) \\ = \mathcal{AS}^{SA}(0, \dots, 0, \Delta_{\mathcal{D}_{tar,t}}^{\theta_{tar,t}}, 0, \dots, 0) = \Delta_{\mathcal{D}_{tar,t}}^{\theta_{tar,t}}. \end{aligned} \quad (6)$$

Once \mathcal{AS} gets the update $\Delta_{\mathcal{D}_{tar,t}}^{\theta_{tar,t}}$, it can get sensitive information about $\mathcal{D}_{tar,t}$ by executing any gradient inversion attack or inference attacks.

3 PROPOSED METHODS

In this section, we design the Janus, a generic privacy-enhanced multi-round SA scheme via a dual-server architecture, where SHC is the core cryptography for verifiability. To facilitate understanding, we first present the new primitive SHC, followed by elaborating on the construction of Janus. Let \odot denote the consecutive operation of \odot . Specifically, $\bigodot_{i=1}^n x_i = x_1 \odot x_2 \dots \odot x_n$, where the \odot indicates addition or multiplication depending on the specific scheme. T is the total number of rounds required for the model to converge and t denotes current round. Let n users participate in FL training, where users are denoted by $\mathcal{U}_t = \{U_i, i \in [1, n]\}$. All users negotiate a model architecture and train the model locally on their private data sets \mathcal{D}_i . There are three types of entities in our system which are aggregation server S_0 , assistant server S_1 , and users. We assume that each user $U_i \in \mathcal{U}_t$ holds a private update x_i of dimension m . For simplicity, we assume that the elements of x_i and $\sum_{U_i \in \mathcal{U}} x_i$ are in \mathbb{Z}_R for R .

3.1 SEPARABLE HOMOMORPHIC COMMITMENT

Definition 1 (*Separable Homomorphic Commitment*). A secure separable homomorphic commitment scheme is a cryptographic protocol that enables secure and flexible commitments. It is comprised of a set of algorithms denoted by the tuple (Setup, Commit, Se, PCommit, Reveal). The formal syntax of each algorithm is described as follows:

- $pp \leftarrow \text{Setup}(1^\lambda)$. A \mathcal{PPT} initialization algorithm takes as input a security parameter λ , and it outputs a public parameters pp .
- $c \leftarrow \text{Commit}(pp, m, r)$. A \mathcal{PPT} commitment algorithm takes as input a public parameter pp , a message m and a random number r , and it outputs a complete commitment c , where $c = (c_m, c_r)$ and c_m is the part associated with the message m and c_r is related to the random number (blinder) r .
- $c_m \leftarrow \text{Se}(pp, c, c_r)$. A Decisional Polynomial Time (\mathcal{DPT}) separation algorithm takes as input a public parameter pp , a complete commitment c and a blinder-related part c_r , and it outputs the message-related commitment c_m .
- $c_m \leftarrow \text{PCommit}(pp, m)$. A \mathcal{DPT} commitment algorithm takes as input a public parameter pp , a message m , and it outputs the message-related commitment c_m .
- $1/0 \leftarrow \text{Reveal}(pp, c, m, r)$. A \mathcal{DPT} revealing commitment algorithm takes as input the public parameter pp , the complete commitment c , the message m and the random blinder r , this algorithm outputs 1 if the m is the valid committed message of c and 0 otherwise.

In addition to the completeness, binding and hiding properties possessed by traditional commitment schemes described in Section 2, the SHC also possess the following two unique properties. The two servers independently aggregate the different values in the commitments. SHC is able to separate part of the message and compare it with the aggregated results, thereby ensuring the correctness of the aggregation.

- **Separability.** The complete commitment c generated by $\text{Commit}(m, r)$ can be divided into two parts $c = (c_m, c_r)$, where c_m is the part associated with the commitment message m and c_r is related to the random blinder r . It can use c_r to extract from the complete commitment c only the parts that are relevant to m . Taking the classic Pedersen commitment (Pedersen, 1991) as an example, the complete commitment is $c = h^r g^m$. Given

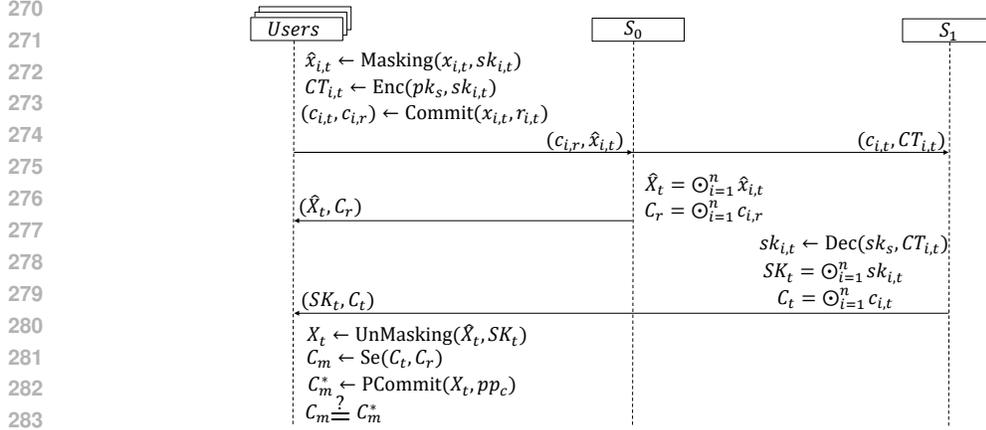


Figure 1: The Workflow of Janus.

$c_r = h^r$, $c_m = g^m$, we can get the c_m from c and c_r via c/c_r . Furthermore, the c_m can be calculated from $\text{PCommit}(m, pp)$.

- **Homomorphism.** Homomorphism facilitates to accomplish secure aggregation. Define the space of message, blinder and commitment as $\mathcal{M}_c, \mathcal{R}_c, \mathcal{C}_c$ respectively.

$$\begin{aligned} \forall (m_0, r_0), (m_1, r_1) \in \mathcal{M}_c \times \mathcal{R}_c : \\ \text{Commit}(m_0 + m_1; r_0 + r_1) = \text{Commit}(m_0; r_0) \cdot \text{Commit}(m_1; r_1). \end{aligned} \quad (7)$$

3.2 THE PROPOSED JANUS

Janus tackles the challenges of dynamic user participation, verifiability, and resistance to model inconsistency attacks that are not addressed in the state-of-the-art Flamingo (S&P’23). Specifically, it has following three key high-level technical ideas:

(1) Dual-server architecture and dynamic user participation. Specifically, the Janus involves two servers, S_0 and S_1 . S_0 is responsible for aggregating the masked updates and S_1 is responsible for aggregating the values associated with the commitments. The dual-server architecture prevents the servers from accessing the final aggregation results, thus effectively avoids attacks such as model reversal and model inconsistency, which are serious privacy leakage in traditional single-server. Furthermore, there is no need to re-establish complex communication diagrams when users join or leave. New users can participate in the new training process by simply generating their own public/private keys and obtaining the servers’ public keys.

(2) Lightweight components and efficient aggregation. Instead of requiring the client to secretly share the mask with all its neighbours as Flamingo and BBSA, Janus does not even require neighbours and avoids the time-consuming process of negotiating keys with each other. It only applies OTP to mask the secret updates and subsequently encrypts the masks via a secure public key encryption. The different messages are then sent to S_0 and S_1 . Thus, no matter how the number of users in the system increases, the operations required by Janus are fixed to the desired constant level.

(3) Verifiability and privacy enhancement. The separability of SHC allows the user to validate the aggregated values locally, thus enabling verifiability. In addition, the binding feature of SHC prevents the client from denying previously sent malicious messages when subsequent misbehavior is detected. This is a feature not available in other advanced schemes. Given the hiding of the SHC and the confidentiality of public key encryption, neither S_0 nor S_1 can access the received secret information. Combined with our dual-server architecture, higher security can be achieved.

Figure 1 shows the workflow of Janus. Subsequently, we provide a detailed description of our Janus, noting that it is a generic construction. Thus, we assume the underlying public key encryption scheme is $\Pi_E = (\text{Setup}, \text{KeyGen}, \text{Enc}, \text{Dec})$, the OTP scheme is $\Pi_O = (\text{Masking}, \text{unMasking})$, and the SHC scheme is $\Pi_S = (\text{Setup}, \text{Commit}, \text{Se}, \text{PCommit}, \text{Reveal})$, in which the setup parts of these schemes are completed in the *Setup* phase of Janus by default. Furthermore, Appendix B

gives the tasks of the different entities in each phase for conciseness and an effective instantiation to demonstrate the practicality. Specifically, Janus consists of the following four phases:

Setup. The objective of this phase is to determine the public parameters pp and specific cryptographic schemes, which ensures that subsequent schemes work properly. In round t , all parties are given the security parameter λ . All public parameters pp of the system are then generated based on λ , e.g., the setup phase and public parameters generation in Π_E, Π_O, Π_S . Each user will generate their private key $sk_{i,t}$ for the OTP. The S_1 will generate its public/private key (pk_s, sk_s) and publish its public key to all participants. Subsequent communications between the users and the servers are encrypted with their respective public keys by default.

Masking and Report. The U_i masks its input updates $x_{i,t}$ via $\text{Masking}(x_{i,t}, sk_{i,t})$ to get the masked updates $\hat{x}_{i,t}$. Subsequently, U_i encrypts the $sk_{i,t}$ using the public key of S_1 via $\text{Enc}(pk_s, sk_{i,t})$ to get the ciphertext $CT_{i,t}$ of $sk_{i,t}$. To achieve subsequent verifiability, U_i makes separable commitment for the input updates $x_{i,t}$ via $\text{Commit}(x_{i,t}, r_{i,t})$ to get the full commitment $c_{i,t}$, where the $r_{i,t}$ is the blinder, the $c_{i,t}$ can be divided into $(c_{i,r}, c_{i,m})$, $c_{i,r}$ is the commitment of blinder and $c_{i,m}$ is the commitment of updates. Then it sends $(\hat{x}_{i,t}, c_{i,r})$ to the aggregation server S_0 and $(c_{i,t}, CT_{i,t})$ to the assistant server S_1 .

Collection and Aggregation. In this phase, the servers will complete the computation secure aggregation and verification for users updates. Specifically, S_0 will aggregate the masked input updates from all users via $X_t = \bigodot_{i=1}^n \hat{x}_{i,t} = \hat{x}_{1,t} \odot \hat{x}_{2,t} \odot \dots \odot \hat{x}_{n,t}$. Then S_0 computes $C_r = \bigodot_{i=1}^n c_{i,r} = c_{1,r} \odot c_{2,r} \odot \dots \odot c_{n,r}$. S_0 sends (\hat{X}_t, C_r) to all users. In fact, \hat{X}_t contains the updated aggregated values for round t and C_r can assist in the validation of aggregated result. For the S_1 , it first decrypts the ciphertext to get the $sk_{i,t}$ via $\text{Dec}(sk_s, CT_{i,t})$. Then it can aggregate the $\bigodot_{i=1}^n sk_{i,t} = sk_{1,t} \odot sk_{2,t} \odot \dots \odot sk_{n,t} = SK_t$. Furthermore, it calculates the aggregation result of the full commitment value for subsequent users to verify the aggregation result completed by S_0 via $\bigodot_{i=1}^n c_{i,t} = c_{1,t} \odot c_{2,t} \odot \dots \odot c_{n,t} = C_t$. Finally, S_1 sends (SK_t, C_t) to all users.

UnMasking and Verification. The users compute the final update results based on the values returned by the two servers and validate the aggregated result. Specifically, U_i gets the final aggregation result via $X_t = \text{UnMasking}(\hat{X}_t, SK_t)$, where the X_t is the updates aggregation result of the round t . To verify the correctness of the aggregation result, U_i extracts the commitment value related to the updates via $C_m = \text{Se}(C_t, C_r)$. The user then calculates the commitment value which is only related to the updates via $C_m^* = \text{PCommit}(X_t, pp_c)$, where the pp_c is the public parameter of the underlying SHC. Finally, U_i compares whether C_m^* and C_m are equal. If they are equal, then the aggregated result is correct; otherwise, it is invalid, and U_i will terminate the subsequent training.

4 EVALUATION

4.1 THEORETICAL ANALYSIS

Janus offers enhanced security compared to state-of-the-art schemes. We give a formal security analysis in appendix C, where Janus can resist MIA and achieve multi-round security. Furthermore, a key advantage of Janus over Flamingo and BBSA is its ability to complete each round with fewer interactions. The two advanced schemes necessitate communication with neighboring nodes to complete the elimination of the mask or decryption process. Assuming that the underlying operations, such as commitments and encryptions, have a complexity of $\mathcal{O}(1)$, Janus demonstrates superior efficiency in terms of interaction count. The remarkable property of Janus is that the system overhead do not grow with the number of users as in previous schemes. The system is designed to be client-friendly, minimizing computational overhead. Clients need only two interactions with the servers to go offline, ensuring there are no issues with aggregation failures or inaccurate results due to user disconnection. We focus on a round of aggregation, with Table 2 presenting the results in comparison to relevant advanced schemes.

Computation Cost. The computation cost of each client consists of: 1) masking the local update by using one-time pad; 2) encrypting the key of one-time pad by public key encryption; 3) committing the local update by using the SHC; 4) unmasking the global aggregation result; 5) separating message-only commitments from the full commitment; 6) calculating the commitment value based on the unmasking result and compare whether it is equal to the separated commitment value to com-

Table 2: Comparison of Performance Analysis

Scheme	Computation		Communication	
	Client	Server	Client	Server
SecAgg	$\mathcal{O}(n^2 + md)$	$\mathcal{O}(dn^2)$	$\mathcal{O}(n + m)$	$\mathcal{O}(n^2 + mn)$
BBSA	$\mathcal{O}(A^2 + lA)$	$\mathcal{O}(n(A^2 + lA))$	$\mathcal{O}(A^2 + l)$	$\mathcal{O}(n(A^2 + l))$
VeriFL	$\mathcal{O}(n)$	$\mathcal{O}(n + l)$	$\mathcal{O}(n)$	$\mathcal{O}(1) + \mathcal{O}(n)$
ELSA	$\mathcal{O}(1 + l)$	$\mathcal{O}(n + nl)$	$\mathcal{O}(1)$	$\mathcal{O}(n)$
Flamingo	Regular Client: $\mathcal{O}(L^2)$ Decryptor: $\mathcal{O}(L^2 + \delta An + (1 - \delta)n + \epsilon n^2)$	$\mathcal{O}(n + L^2)$	Regular Client: $\mathcal{O}(l + A + L^2)$ Decryptors: $\mathcal{O}(L^2 + L + \delta An + (1 - \delta)n)$	$\mathcal{O}(L^3 + n(l + L + A))$
Janus	$\mathcal{O}(1 + l)$	$\mathcal{O}(n + nl)$	$\mathcal{O}(1)$	$\mathcal{O}(n)$

* Let n, L, A denote the total number of clients, the number of decryptors and the upper bound number of neighbors of a client respectively, where $A = \log n$ in BBSA. l denotes the dimension of the update. δ denotes the dropout rate respectively. ϵ is the parameter of graph generation.

plete the verification. All the above operations take only $\mathcal{O}(1)$ time each. Overall, the computational overhead of each client is constant. The computation cost of S_0 mainly consists of aggregating the masking updates from clients and the commitment of random numbers, which both take $\mathcal{O}(n)$. Thus the total computational overhead grows linearly with the number of clients. For S_1 , the computation cost consists of: 1) decrypting the ciphertext of the private key of one-time pad; 2) aggregating the private keys for masking; 3) aggregating the complete commitments for subsequent verification of the aggregation result of S_0 . All these operations mentioned above take $\mathcal{O}(n)$. Overall, the communication overhead of servers grows linearly with the number of clients which takes $\mathcal{O}(n)$.

Communication Cost. Each client needs to send one masked message to S_0 , one encrypted and committed message to S_1 . Overall, the computational overhead of each client is constant. For the servers, S_0 will send the aggregation result of the masking updates to all clients, which takes $\mathcal{O}(n)$. S_1 sends the aggregation result of the key used for one-time pad and the full commitment to all clients, which also takes $\mathcal{O}(n)$. Overall, for servers, their communication overhead grows linearly with the number of clients which takes $\mathcal{O}(n)$.

4.2 MODEL PERFORMANCE

In this section, we carried out various experiments to verify the effectiveness and efficiency of our scheme and to compare it with similar advanced schemes. Our experimental setup includes a 13th Gen Intel(R) Core(TM) i7-13700KF 3.40 GHz processor with 32.0 GB of RAM, a 64-bit Windows 11 operating system, and an RTX 4070Ti GPU display adapter.

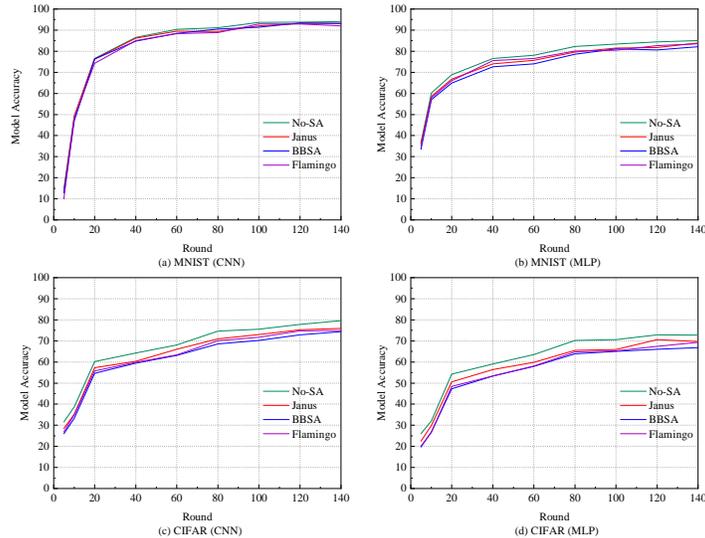


Figure 2: Test accuracy across different datasets and models.

432 **Baselines.** To evaluate the impact of SA on federated learning (e.g., training effectiveness, commu-
 433 nication time), we implemented the original FL framework (No-SA), where the server aggregates
 434 clear updates from users in each training round (McMahan et al., 2017). Bell et al. (2020) optimized
 435 the communication graph of the first mask-based SA scheme (Bonawitz et al., 2017) and proposed
 436 an advanced scheme (BBSA), which we implemented for comparison. For client dropout, we con-
 437 struct the graph with responsive clients, yielding better results than the original. Flamingo (Ma
 438 et al., 2023) introduced multi-round aggregation, and both Flamingo and BBSA involve waiting for
 439 messages from at least t out of n clients.

440 **Datasets and Models.** MNIST consists of 70,000 grayscale handwritten digit images (60,000 for
 441 training, 10,000 for testing), each 28x28 pixels. We use 100 clients, each with 600 training samples.
 442 The global model for MNIST is a fully connected network with layers of size (784, 256, 10). CIFAR-
 443 10 includes 60,000 color images across 10 classes (50,000 training, 10,000 testing), using a CNN
 444 architecture with a batch size of 10, learning rate of 0.001, and 100 training epochs. We employed
 445 SGD as the optimizer, with each client applying SGD once per global epoch (local epoch = 1).

446 To comprehensively evaluate the impact of the security SA in this paper on the model training effec-
 447 tiveness, our experiments are carried out on different datasets and models. We conducted the training
 448 with 100 clients and compare the test accuracy of our Janus with related schemes. Figure 2 shows
 449 the comparison results. The following conclusions can be drawn from the experimental results.
 450 Firstly, the final test accuracy at model convergence is not much different between our scheme and
 451 the compared schemes, in which No-SA has the highest accuracy, and our scheme follows closely.

452 Specifically, for MNIST, the test accuracy of No-SA can reach to 94.1% under the CNN, while the
 453 Janus can also reached about 93.18%. Additionally, the test accuracy of No-SA can reach to 85.04%
 454 under the MLP, while the Janus can also reached about 83.95%. Compared to other schemes, Janus
 455 has considerable accuracy. As for the CIFAR, the test accuracy of No-SA can reach to 77.8% under
 456 the CNN, while the Janus can also reached about 75.94%. Additionally, the test accuracy of No-SA
 457 can reach to 72.8% under the MLP, while the Janus can also reached about 71.6%.

458 Figure 3 shows the loss of related schemes during the training process with different datasets and
 459 models. It can be concluded that as the number of training rounds increases, the loss values for
 460 the same dataset with different secure aggregation schemes applied are smoother and eventually all
 461 converge to be almost equal. This shows that our Janus, like advanced schemes, does not result in a
 462 loss of model performance due to the use of secure aggregation. The impact on the model is similar
 463 to that of existing advanced schemes, while protecting users privacy and providing better efficiency.
 464
 465

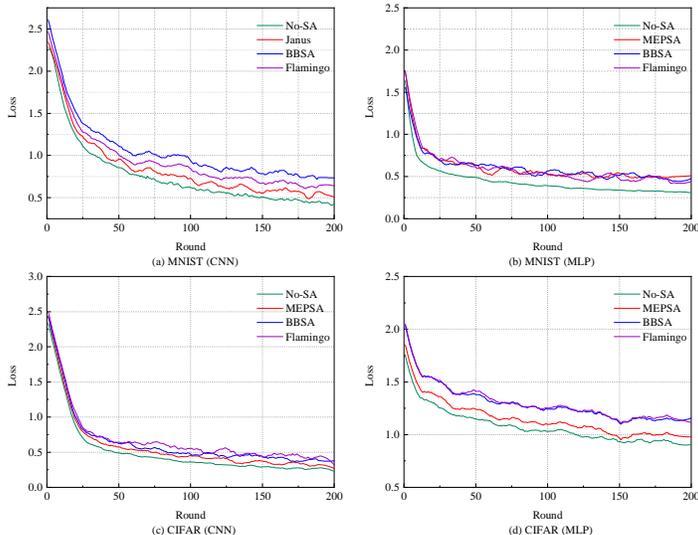


Figure 3: Training loss across different datasets and models.

4.3 COMPUTATION OVERHEAD

Since masking-based schemes are not resistant to user dropouts, we consider this case when implementing BBSA and Flamingo. Specifically, we only consider the case where 10% of users drop out, but it should be noted that in practice the waiting time required to solve the user dropout problem is much longer than that considered in our experiments, due to the complexity and diversity of the real scenarios. Moreover, it is important to note that some of the evaluated schemes inherently support multi-round aggregation, while others do not. We adapted the schemes that lack built-in multi-round aggregation capabilities by running them multiple times to simulate the effect of multi-round aggregation. Although this approach is feasible, it introduces a considerable amount of additional and unnecessary computation overhead. This further highlights the advantages of our proposed scheme, Janus, which is natively designed to support multi-round aggregation without incurring such overhead, thus demonstrating superior efficiency and scalability in practice.

As shown in Figure 4, we present a comparative analysis of the time overhead of various schemes, focusing on the completion time required for a single aggregation. It should be noted that, due to differences in the stages involved across these schemes, only the relevant time-consuming stages were considered for each. From the results, several conclusions can be drawn. First, the computational overhead introduced by SA is within an acceptable range, demonstrating its practicality in real-world applications. More importantly, our proposed scheme exhibits significantly lower overhead, particularly on the client side, which substantially enhances overall efficiency. This improvement can be attributed to the adoption of lightweight cryptographic components, which circumvent time-intensive operations such as secret sharing and DH key negotiation. The absence of these complex operations reduces the computational burden on clients, thereby contributing to the superior performance of our scheme.

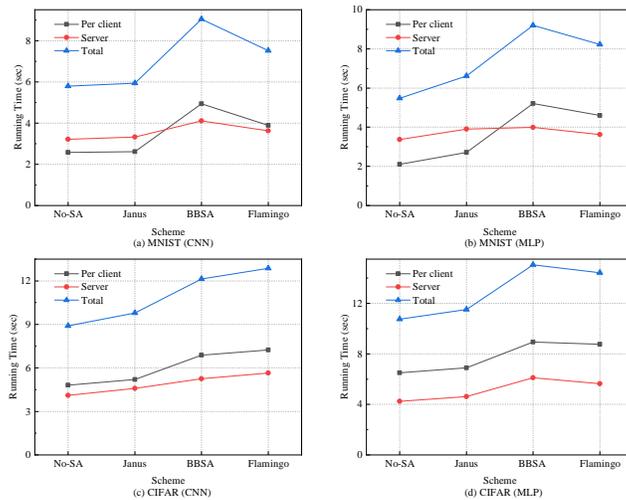


Figure 4: Computation overhead across different datasets and models.

5 CONCLUSION

In this paper, we propose a new cryptographic primitive, i.e., separable homomorphic commitment and design a generic dual-server multi-round SA scheme called Janus for federated learning. Janus addresses the issues of dynamic user participation, verifiability, and resistance to model inconsistency attacks that are not considered in advanced Flamingo (S&P’23). It not only significantly enhances security but also improves system efficiency, which reduces per-client communication and computation overhead from a logarithmic to a constant scale compared to current state-of-the-art methods, with almost no compromise in model accuracy. Finally, we evaluate Janus from both theoretical and experimental perspectives, demonstrating its superior security and performance. Future researches on integrating Janus with various advanced privacy-preserving techniques could further enhance its security. Additionally, secure and effective identification of data poisoning attacks from the users is another worthwhile research direction.

REFERENCES

- 540
541
542 James Henry Bell, Kallista A. Bonawitz, Adrià Gascón, Tancrede Lepoint, and Mariana Raykova.
543 Secure single-server aggregation with (poly)logarithmic overhead. In *CCS*, pp. 1253–1269. ACM,
544 2020.
- 545 Kallista A. Bonawitz, Vladimir Ivanov, Ben Kreuter, Antonio Marcedone, H. Brendan McMahan,
546 Sarvar Patel, Daniel Ramage, Aaron Segal, and Karn Seth. Practical secure aggregation for
547 privacy-preserving machine learning. In *CCS*, pp. 1175–1191. ACM, 2017.
- 548
549 Kallista A. Bonawitz, Fariborz Salehi, Jakub Konečný, Brendan McMahan, and Marco Gruteser.
550 Federated learning with autotuned communication-efficient secure aggregation. In *ACSSC*, pp.
551 1222–1226. IEEE, 2019.
- 552 Carlo Brunetta, Georgia Tsaloli, Bei Liang, Gustavo Banegas, and Aikaterini Mitrokotsa. Non-
553 interactive, secure verifiable aggregation for decentralized, privacy-preserving learning. In
554 *ACISP*, volume 13083 of *Lecture Notes in Computer Science*, pp. 510–528. Springer, 2021.
- 555
556 Huancheng Chen and Haris Vikalo. Recovering labels from local updates in federated learning. In
557 *ICML*. OpenReview.net, 2024.
- 558
559 Whitfield Diffie and Martin E. Hellman. New directions in cryptography. *IEEE Trans. Inf. Theory*,
560 22(6):644–654, 1976.
- 561 Joaquín Delgado Fernández, Sergio Potenciano Menci, Charles Lee, and Gilbert Fridgen. Secure
562 federated learning for residential short term load forecasting. *CoRR*, abs/2111.09248, 2021.
- 563
564 Jiqiang Gao, Boyu Hou, Xiaojie Guo, Zheli Liu, Ying Zhang, Kai Chen, and Jin Li. Secure aggrega-
565 tion is insecure: Category inference attack on federated learning. *IEEE Trans. Dependable Secur.*
566 *Comput.*, 20(1):147–160, 2023.
- 567
568 Kostadin Garov, Dimitar Iliev Dimitrov, Nikola Jovanovic, and Martin T. Vechev. Hiding in plain
569 sight: Disguising data stealing attacks in federated learning. In *ICLR*. OpenReview.net, 2024.
- 570
571 Xiaojie Guo, Zheli Liu, Jin Li, Jiqiang Gao, Boyu Hou, Changyu Dong, and Thar Baker. Verifi-
572 cation-efficient and fast verifiable aggregation for federated learning. *IEEE Trans. Inf.*
Forensics Secur., 16:1736–1751, 2021.
- 573
574 Yue Guo, Antigoni Polychroniadou, Elaine Shi, David Byrd, and Tucker Balch. Microfedml: Pri-
575 vacy preserving federated learning for small weights. *IACR Cryptol. ePrint Arch.*, pp. 714, 2022.
- 576
577 Changhee Hahn, Hodong Kim, Minjae Kim, and Junbeom Hur. Versa: Verifiable secure aggregation
578 for cross-device federated learning. *IEEE Trans. Dependable Secur. Comput.*, 20(1):36–52, 2023.
- 579
580 Briland Hitaj, Giuseppe Ateniese, and Fernando Pérez-Cruz. Deep models under the GAN: infor-
581 mation leakage from collaborative deep learning. In *CCS*, pp. 603–618. ACM, 2017.
- 582
583 Yangsibo Huang, Samyak Gupta, Zhao Song, Kai Li, and Sanjeev Arora. Evaluating gradient inver-
584 sion attacks and defenses in federated learning. In *NeurIPS*, pp. 7232–7241, 2021.
- 584
585 Jonathan Katz and Yehuda Lindell. *Introduction to Modern Cryptography, Second Edition*. CRC
586 Press, 2014.
- 587
588 Ziyao Liu, Jiale Guo, Wenzhuo Yang, Jiani Fan, Kwok-Yan Lam, and Jun Zhao. Privacy-preserving
589 aggregation in federated learning: A survey. *CoRR*, abs/2203.17005, 2022.
- 589
590 Lu Lu, Yeonjong Shin, Yanhui Su, and George E. Karniadakis. Dying relu and initialization: Theory
591 and numerical examples. *CoRR*, abs/1903.06733, 2019.
- 592
593 Yiping Ma, Jess Woods, Sebastian Angel, Antigoni Polychroniadou, and Tal Rabin. Flamingo:
Multi-round single-server secure aggregation with applications to private federated learning. In
SP, pp. 477–496. IEEE, 2023.

- 594 Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Agüera y Arcas.
595 Communication-efficient learning of deep networks from decentralized data. In *AISTATS*, vol-
596 ume 54 of *Proceedings of Machine Learning Research*, pp. 1273–1282. PMLR, 2017.
597
- 598 Milad Nasr, Reza Shokri, and Amir Houmansadr. Comprehensive privacy analysis of deep learning:
599 Passive and active white-box inference attacks against centralized and federated learning. In *IEEE*
600 *Symposium on Security and Privacy*, pp. 739–753. IEEE, 2019.
- 601 John Nguyen, Kshitiz Malik, Hongyuan Zhan, Ashkan Yousefpour, Mike Rabbat, Mani Malek, and
602 Dzmityr Huba. Federated learning with buffered asynchronous aggregation. In *AISTATS*, volume
603 151 of *Proceedings of Machine Learning Research*, pp. 3581–3607. PMLR, 2022.
- 604 Dario Pasquini, Danilo Francati, and Giuseppe Ateniese. Eluding secure aggregation in federated
605 learning via model inconsistency. In *CCS*, pp. 2429–2443. ACM, 2022.
- 606
- 607 Torben P. Pedersen. Non-interactive and information-theoretic secure verifiable secret sharing. In
608 *CRYPTO*, volume 576 of *Lecture Notes in Computer Science*, pp. 129–140. Springer, 1991.
609
- 610 Pian Qi, Diletta Chiaro, Antonella Guzzo, Michele Ianni, Giancarlo Fortino, and Francesco Piccialli.
611 Model aggregation techniques in federated learning: A comprehensive survey. *Future Gener.*
612 *Comput. Syst.*, 150:272–293, 2024.
- 613 Mayank Rathee, Conghao Shen, Sameer Wagh, and Raluca Ada Popa. ELSA: secure aggregation
614 for federated learning with malicious actors. In *SP*, pp. 1961–1979. IEEE, 2023.
- 615
- 616 Ahmed Salem, Yang Zhang, Mathias Humbert, Pascal Berrang, Mario Fritz, and Michael Backes.
617 MI-leaks: Model and data independent membership inference attacks and defenses on machine
618 learning models. In *NDSS*. The Internet Society, 2019.
- 619 Thomas Sandholm, Sayandev Mukherjee, and Bernardo A. Huberman. SAFE: secure aggregation
620 with failover and encryption. *CoRR*, abs/2108.05475, 2021.
621
- 622 Adi Shamir. How to share a secret. *Commun. ACM*, 22(11):612–613, 1979.
- 623
- 624 Reza Shokri, Marco Stronati, Congzheng Song, and Vitaly Shmatikov. Membership inference at-
625 tacks against machine learning models. In *IEEE Symposium on Security and Privacy*, pp. 3–18.
626 IEEE Computer Society, 2017.
- 627 Jinhyun So, Basak Güler, and Amir Salman Avestimehr. Turbo-aggregate: Breaking the quadratic
628 aggregation barrier in secure federated learning. *IEEE J. Sel. Areas Inf. Theory*, 2(1):479–489,
629 2021.
- 630 Jinhyun So, Corey J. Nolet, Chien-Sheng Yang, Songze Li, Qian Yu, Ramy E. Ali, Basak Guler,
631 and Salman Avestimehr. Lightsecagg: a lightweight and versatile design for secure aggregation
632 in federated learning. In *MLSys*. mlsys.org, 2022.
633
- 634 Jinhyun So, Ramy E. Ali, Basak Güler, Jiantao Jiao, and Amir Salman Avestimehr. Securing secure
635 aggregation: Mitigating multi-round privacy leakage in federated learning. In *AAAI*, pp. 9864–
636 9873. AAAI Press, 2023.
- 637 Timothy Stevens, Christian Skalka, Christelle Vincent, John Ring, Samuel Clark, and Joseph P. Near.
638 Efficient differentially private secure aggregation for federated learning via hardness of learning
639 with errors. In *USENIX Security Symposium*, pp. 1379–1395. USENIX Association, 2022.
640
- 641 Georgia Tsaloli, Bei Liang, Carlo Brunetta, Gustavo Banegas, and Aikaterini Mitrokotsa. sf DEVA:
642 decentralized, verifiable secure aggregation for privacy-preserving learning. In *ISC*, volume
643 13118 of *Lecture Notes in Computer Science*, pp. 296–319. Springer, 2021.
- 644 Haozhao Wang, Haoran Xu, Yichen Li, Yuan Xu, Ruixuan Li, and Tianwei Zhang. Fedcda: Feder-
645 ated learning with cross-rounds divergence-aware aggregation. In *ICLR*. OpenReview.net, 2024.
646
- 647 Di Wu, Jun Bai, Yiliao Song, Junjun Chen, Wei Zhou, Yong Xiang, and Atul Sajjanhar. Fedinverse:
Evaluating privacy leakage in federated learning. In *ICLR*. OpenReview.net, 2024.

- 648 Yueqi Xie, Minghong Fang, and Neil Zhenqiang Gong. Fedredefense: Defending against model
649 poisoning attacks for federated learning using model update reconstruction error. In *ICML*. Open-
650 Review.net, 2024.
- 651
- 652 Guowen Xu, Hongwei Li, Sen Liu, Kan Yang, and Xiaodong Lin. Verifynet: Secure and verifiable
653 federated learning. *IEEE Trans. Inf. Forensics Secur.*, 15:911–926, 2020.
- 654
- 655 Samuel Yeom, Irene Giacomelli, Matt Fredrikson, and Somesh Jha. Privacy risk in machine learn-
656 ing: Analyzing the connection to overfitting. In *CSF*, pp. 268–282. IEEE Computer Society,
657 2018.
- 658 Chengliang Zhang, Suyi Li, Junzhe Xia, Wei Wang, Feng Yan, and Yang Liu. Batchcrypt: Effi-
659 cient homomorphic encryption for cross-silo federated learning. In *USENIX Annual Technical*
660 *Conference*, pp. 493–506. USENIX Association, 2020a.
- 661
- 662 Xianglong Zhang, Anmin Fu, Huaqun Wang, Chunyi Zhou, and Zhenzhu Chen. A privacy-
663 preserving and verifiable federated learning scheme. In *ICC*, pp. 1–6. IEEE, 2020b.
- 664
- 665 Ligeng Zhu, Zhijian Liu, and Song Han. Deep leakage from gradients. In *NeurIPS*, pp. 14747–
666 14756, 2019.
- 667
- 668 Haomin Zhuang, Mingxian Yu, Hao Wang, Yang Hua, Jian Li, and Xu Yuan. Backdoor federated
669 learning by poisoning backdoor-critical layers. In *ICLR*. OpenReview.net, 2024.

670 A RELATED WORKS

671

672

673 The main goal of FL is to protect the privacy of local data while still allowing it to be used to train
674 the public models. A significant amount of research has been conducted around SA. This section
675 reviews the work related to our scheme. Refer to the reference (Qi et al., 2024; Liu et al., 2022) for
676 more extensive survey in this field.

677 **Masking-based SA.** Masking is a classic encryption technique based on one-time pad (Katz & Lin-
678 dell, 2014). Bonawitz et al. (2017) designed the first SA scheme (SecAgg) which using pairwise
679 masks to hide individual inputs for FL. However, their scheme involves a complete communication
680 graph, which incurs heavy computation and communication for each client linear in the number of
681 participants. Subsequently, Bell et al. (2020) replaced that with a k -regular graph of logarithmic
682 degree, which greatly improved the efficiency while maintaining the security. Stevens et al. (2022)
683 replaced the standard mask with learning with errors mask and used verifiable secret sharing to pre-
684 vent malicious users from distributing incorrect shares. Sandholm et al. (2021) arranged the users in
685 the system in the form of a ring chain, the efficiency of the scheme has been significantly improved,
686 and the user drop problem can be effectively solved. Most masking-based schemes require double
687 masking in order to solve the problem of dropped users. Bonawitz et al. (2019) combined the ran-
688 dom rotation technique to actively adjust the quantisation range of the model in order to reduce the
689 model volume. To reduce the communication overhead, TurboAgg (So et al., 2021) divides n users
690 into $n/\log n$ groups and then uses a multi-group loop structure for subsequent aggregation.

691 **Attacks that Bypass SA.** The aggregation results of most existing schemes are visible to both the
692 clients and the aggregation server. However, this can lead to attacks where malicious servers bypass
693 the SA. Pasquini et al. (2022) proposed model inconsistency attacks, where a malicious server can
694 distribute different parameters to targeted and non-targeted users. This can trigger *dying-ReLU*
695 and make the input of non-target users be zero. So et al. (2023) noticed that when the trained
696 model begins to converge, the client model changes little between one training step and the next.
697 A malicious server can infer the updates of a client that participated in the previous round but did
698 not participate in the subsequent round from the aggregation results. Gao et al. (2023) proposed a
699 scheme which can launch a category inference attack even in presence of SA. To avoid this type
700 of attack, when the clients receive the model parameters, they need to verify whether the received
701 parameters are consistent or not, and terminate the training if they are not. But this will increase the
system overhead. Fernández et al. (2021) applied differential privacy on the aggregated model to
hide the aggregation results.

Server-side Attacks and Defenses. Membership inference attacks pose a potential threat from the server side in FL. Specifically, an adversary can determine whether some specific data records are part of the local training dataset of a target user only by accessing the model updates, either through a black-box or white-box approach. Yeom et al. (2018) proposed the first label-based attack, which aims at predicting whether an instance is in the local data of the target user. The attacker exploits the performance disadvantage of the target model in the test dataset to complete this attack. Chen & Vikalo (2024) proposed a general analytical method that allows the FL server to recover client training labels, applicable to various FL algorithms without assumptions on activation functions or batch label composition. Shokri et al. (2017) designed an attack with partial output knowledge in a black box-scenario. Furthermore, Salem et al. (2019) improved a new attack by using the maximum value of the model output confidence. Zhuang et al. (2024) introduced the layer substitution analysis, a new technique that identifies layers critical for backdoor injection, making it well-suited for FL attacks. Leveraging this technique, they developed two layer-wise backdoor attack strategies that successfully implant backdoors into these key layers and evade state-of-the-art defenses without compromising the primary task accuracy.

Meanwhile, Bonawitz et al. (2017) proposed the first SA scheme to compute the sum of model updates hiding personal information. Subsequently, a great deal of research has centred around SA. Techniques such as homomorphic encryption (Zhang et al., 2020a), differential privacy (Stevens et al., 2022), and multi-party computation (Bell et al., 2020) are used to construct SA schemes to protect user privacy from attack by malicious servers. SA based on cryptography aims to prevent attacks by concealing model updates from any potential adversaries. This approach ensures that individual contributions remain private, making it difficult for malicious entities to infer sensitive information from the data.

Recently, Xie et al. (2024) identify a limitation in existing model poisoning attacks defenses: reliance on cross-client or global information, which leads to performance degradation under non-IID data distributions or when there is a large number of malicious clients. Then they establish a crucial distinction between model poisoning attacks and benign model updates by determining whether the update can be approximately reconstructed using distilled local knowledge. Wu et al. (2024) proposed FedInverse, a framework designed to evaluate whether FL models are susceptible to model inversion attacks and quantify the associated data-leakage risks. Garov et al. (2024) showed that all existing malicious server attacks can be identified through systematic checks. Furthermore, they established a set of essential requirements that any practical malicious server attack must meet.

Verifiability. In addition, a malicious server might return incorrect aggregation results to gain an unfair advantage or disrupt the system’s integrity. Such behavior poses significant security threats, as users or clients relying on these results could be misled or manipulated. Therefore verifiable SA is necessary to ensure correct aggregation. Zhang et al. (2020b) verified the aggregation result via homomorphic encryption SA using homomorphic hash function. Additionally, Xu et al. (2020) verified masking-based SA using the same technique. Guo et al. (2021) proposed a verification scheme which focuses on the high dimension inputs. Brunetta et al. (2021) proposed a non-interactive verifiable SA protocol from NIVA, which requires users create a tag for each input shares. In contrast, Tsaloli et al. (2021) proposed a scheme requires only a single tag for each user.

Multi-round Setting and Dynamic Joining. Model convergence in Federated Learning (FL) typically requires multiple rounds of training, with each round contributing incrementally to the overall performance of the global model. However, most existing state-of-the-art SA schemes are designed to support only a single round of aggregation. In addition to protecting user privacy in single rounds of FL training, some studies have looked at privacy issues arising from multiple rounds of FL training. Nguyen et al. (2022) and So et al. (2022) proposed two new schemes support asynchronous aggregation. Guo et al. (2022) designed a multi-round SA protocol for reusable secrets, and their scheme is mainly oriented towards scenarios with small inputs (the input vector with small values).

Recently, Ma et al. (2023) proposed Flamingo, which has no restrictions on input value. So et al. (2023) mitigated the privacy leakage involved in multi-round aggregation through client selection. Furthermore, the existing schemes do not support dynamic joining. Flamingo assumes that the set of all clients (n) participating in the training is fixed before the training starts and some subset is selected from n in each round t . Therefore, Flamingo does not support the user to dynamically add in the training process. Most current schemes require reconstruction of the communication graph when new users join and require key negotiation with each other user, which imposes huge

756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809

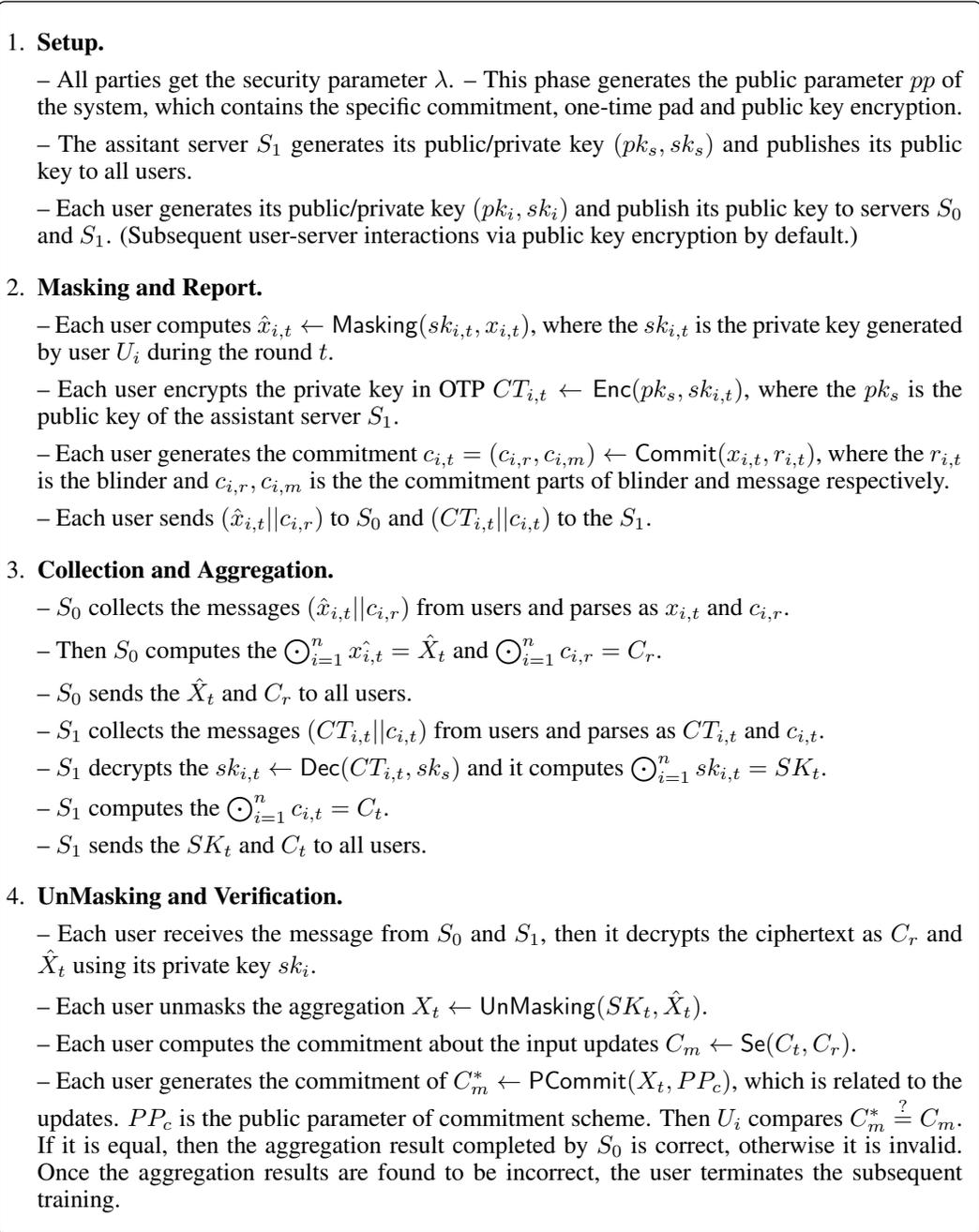


Figure 5: Detailed Construction of Janus.

communication and computation overheads. In addition, Wang et al. (2024) focus on the aggregation of cross-round local models. They proposed FedCDA, a novel cross-round aggregation method that constructs the global model by aggregating local models from multiple rounds based on minimum divergence. To enhance efficiency, FedCDA further introduces an approximation strategy to reduce selection overhead.

B DETAILED JANUS AND ITS INSTANTIATION

In this section, Figure 5 gives the full generic construction of Janus. Furthermore, we give an effective instantiation of our generic construction, where the underlying SHC is Pedersen commitment,

public key encryption is ElGamal, one-time pad is based on normal addition encryption. Specifically, our scheme consists of the following five phases: setup, masking and report, collection and aggregation, collection and aggregation, unmasking and verification.

Setup. This phase determines the public parameters of the system. Firstly, all participants agree on the security parameter λ . The public parameters of the cryptographic primitives are then generated based on the security parameter. Define a triplet (p, q, g, h) , where p is a randomly chosen prime of length $|q| = \lambda + \delta$, the δ is a specified constant, q is a prime order group of \mathbb{Z}_p^* , and g, h are random generators of group of q order, $q = (p - 1)/\gamma$ is prime and the γ is a specified small integer. U_i generates public/private key $(sk_i, pk_i) = (sk_i, g^{sk_i} \pmod{p})$, where the $sk_i \in \mathbb{Z}_p^*$. S_1 generates public/private key $(sk_s, pk_s) = (sk_s, g^{sk_s} \pmod{p})$ where the $sk_s \in \mathbb{Z}_p^*$. Then S_1 and U_i publish their public keys to all entities while store their private keys secretly.

Masking and Report. Each user U_i trains local data \mathcal{D}_i to get the updates $x_{i,t}$ for round t . U_i masks the vector by $\hat{x}_{i,t} = \text{Masking}(x_{i,t}, sk_{i,t}) = x_{i,t} + sk_{i,t} \pmod{p}$. Then U_i encrypts the $sk_{i,t}$ by $\text{Enc}(pk_s, sk_{i,t}) = CT_{i,t} = (g^{sk_{i,t}} \pmod{p}, sk_{i,t} pk_s^{k_{i,t}} \pmod{p})$. Furthermore, U_i commits the $x_{i,t}$ by $c_{i,t} = \text{Commit}(x_{i,t}, r_{i,t}) = g^{x_{i,t}} h^{r_{i,t}} \pmod{p}$, where the $r_{i,t} \in \mathbb{Z}_p^*$ and $c_{i,t} = (c_{i,r}, c_{i,m}) = (h^{r_{i,t}} \pmod{p}, g^{x_{i,t}} \pmod{p})$. Finally, U_i sends $c_{i,r}$ and $\hat{x}_{i,t}$ to S_0 , $c_{i,t}$ and $CT_{i,t}$ to S_1 .

Collection and Aggregation. Subsequently, S_0 receives the message from U_i . Then it computes $\bigodot_{i=1}^n \hat{x}_{i,t} = \hat{x}_{1,t} + \hat{x}_{2,t} + \dots + \hat{x}_{n,t} = \hat{X}_t$ and $\bigodot_{i=1}^n c_{i,r} = h^{r_{1,t}} h^{r_{2,t}} \dots h^{r_{n,t}} \pmod{p} = C_r$. Then S_0 sends C_r and \hat{X}_t to all users. When the S_1 receives the message from U_i . It first decrypts $sk_{i,t} = \text{Dec}(sk_s, CT_{i,t}) = sk_{i,t} pk_s^{k_{i,t}} (g^{sk_{i,t}})^{-1} \pmod{p}$. Subsequently, it computes $\bigodot_{i=1}^n c_{i,t} = c_{1,t} c_{2,t} \dots c_{n,t} \pmod{p} = C_t$. Then it computes $\bigodot_{i=1}^n sk_{i,t} = sk_{1,t} + sk_{2,t} + \dots + sk_{n,t} = SK_t$. Finally, S_1 sends the C_t and SK_t to all users.

Unmasking and Verification. When U_i receives the message from S_0 and S_1 . Firstly, U_i computes the $X_t = \text{Unmasking}(\hat{X}_t, SK_t) = \hat{X}_t - SK_t$ to get the aggregation result X_t . To verify the validity of the aggregation results, U_i separates the parts of the commitments that are only relevant to the input updates by $\text{Se}(C_t, C_r) = C_m$. Then U_i makes a commitment to the aggregation result from S_0 through $\text{PCommit}(X_t, pp_c) = g^{\hat{X}_t} \pmod{p} = C_m^*$, where pp_c is the public parameters of the underlying SHC. Eventually U_i compares whether $C_m^* \stackrel{?}{=} C_m$ holds, if it does it indicates that the aggregation result \hat{X}_t from S_0 is correct, otherwise the aggregation result is not valid. U_i will refuse to accept the results of the aggregation and aborted the subsequent training.

Correctness. The correctness of this instantiation requires each user will obtain the correct aggregation result and the valid verification as long as each entities run the protocol honestly. It is not hard to prove this due to the correctness of the underlying public key encryption, one-time pad and SHC. Specifically, we assume that the aggregation server S_0 receives all masked-input and performs Janus correctly, the following condition holds.

$$\begin{aligned}
\bigodot_{i=1}^n \hat{x}_{i,t} &= \hat{x}_{1,t} + \hat{x}_{2,t} + \dots + \hat{x}_{n,t} \\
&= x_{1,t} + sk_{i,t} + x_{2,t} + sk_{2,t} + \dots + x_{n,t} + sk_{n,t} \\
&= \bigodot_{i=1}^n x_{i,t} + \bigodot_{i=1}^n sk_{i,t} \\
&= X_t + SK_t,
\end{aligned} \tag{8}$$

where the $\bigodot_{i=1}^n sk_{i,t}$ is computed by S_1 . The final aggregation result is $\bigodot_{i=1}^n x_{i,t} = \bigodot_{i=1}^n \hat{x}_{i,t} - \bigodot_{i=1}^n sk_{i,t} = X_t$. If the validation passes, the following condition holds.

$$\begin{aligned}
C_t &= g^{x_{1,t}} h^{r_{1,t}} g^{x_{2,t}} h^{r_{2,t}} \dots g^{x_{n,t}} h^{r_{n,t}} \\
&= g^{x_{1,t} + x_{2,t} + \dots + x_{n,t}} h^{r_{1,t} + r_{2,t} + \dots + r_{n,t}}, \\
C_r &= h^{r_{1,t} + r_{2,t} + \dots + r_{n,t}}, \\
C_m &= C_t / C_r = g^{x_{1,t} + x_{2,t} + \dots + x_{n,t}}, \\
C_m^* &= g^{X_t}.
\end{aligned} \tag{9}$$

If the aggregation result X_t from S_0 is correct, then the $C_m^* = C_m$ will always hold.

C SECURITY ANALYSIS

In this section, we intend to demonstrate the security of our generic construction. We first give the threat model and prove the Janus can protect the privacy of users’ local updates and the aggregated updates. Finally, we give the security proof of single round and multi-round.

C.1 THREAT MODEL

All users agree to publish the final results of model aggregation only to each user, but not to the servers to resist MIA. These users have a common interest in soundness (i.e., getting the correct global model aggregation updates from untrusted servers) and privacy (i.e., hiding local model updates from each other and the server). The specific assumptions in our paper are as follows: The two servers will not collude but may perform incorrect aggregation. The scheme also allows for up to $n - 2$ clients to collude. Specifically, even if the server aggregates incorrect results, our scheme provides verifiability, which enables us to detect such behavior and mitigate the associated risks. If the server colludes with up to $n - 2$ clients, it can only obtain the additive result of the remaining two uncolluding clients. This result is an aggregation of two encrypted or obfuscated values, making it impossible to recover each uncolluding user’s specific gradient information. This ensures that the colluding entities cannot initiate a MIA or access the private information of the remaining two non-colluding clients. When $n - 2$ clients collude, this assumption is even weaker, as the absence of server involvement further limits the accessible information, making it even harder to extract useful data. If only a single server is corrupted, this does not compromise individual user privacy. For instance, with server S_0 , as long as the underlying encryption algorithm is secure, the server cannot access the user-submitted private data without the user’s private key. Similarly, for server S_1 , the security of the underlying SHC ensures that its hiding properties prevent S_1 from obtaining any private information. In conclusion, the assumptions of our scheme are reasonable and well-supported. We will incorporate these clarifications in the revised version to better highlight the theoretical advantages of our approach. In addition, we assume the channel between each user and servers are secure, which allows each entities to authenticate the incoming messages and prevent outsiders from injecting their responses. Furthermore, we assume that there is no collusion between all entities in the system. Our security proofs are based on this threat model.

C.2 PRIVACY FROM USERS

In the “honest but curious” setting, each client will honestly adhere to the protocol but attempt to infer the local gradients of clients and the aggregated gradients. Therefore, we can use the standard simulation proof for multi-party computation protocols to demonstrate the privacy of our generic construction. We first consider privacy protection against honest-but-curious clients who hold their own local gradients and have access to the global gradients. Specifically, let Π denote the proposed Janus involving n users C_1, C_2, \dots, C_n and two servers S_0 and S_1 . Each user holds a local update gradient x_i , Janus securely computes the aggregated global update X . All participants may attempt to infer more additional information, the Π satisfies the following privacy guarantee:

- For each honest-but-curious client C_i , the client learns nothing beyond its own local gradient x_i and the final global aggregated gradient X . Formally, for each C_i , there exists a \mathcal{PPT} simulator \mathcal{S}_i such that:

$$\{\mathbf{View}_{\Pi}(C_i)\} \approx \{\mathcal{S}_i(x_i, X)\}, \quad (10)$$

where $\mathbf{View}_{\Pi}(C_i)$ denotes the view of C_i during the real execution of Π , x_i is the C_i ’s local updates and X is the final global aggregated result.

- For S_0 and S_1 , they learns nothing beyond the masked aggregated results and the aggregated results of masks. This can ensure they will learn nothing about the final global aggregated gradient X , thus resisting the MIA. Formally, for S_0 and S_1 , there exists a \mathcal{PPT} simulator \mathcal{S}_{server} such that:

$$\{\mathbf{View}_{\Pi}(S_0, S_1)\} \approx \{\mathcal{S}_{server}(\hat{X}, CT)\}, \quad (11)$$

where $\mathbf{View}_{\Pi}(S_0, S_1)$ denotes the view of two servers during the real execution of Π , \hat{X} is the masked aggregation result and CT is the ciphertext of masks.

Given any subset $\mathcal{U} \subseteq \mathcal{C}$ of the users, where the \mathcal{C} is the set of all users in the system ($|\mathcal{C}| = m$). Let the $\mathbf{REAL}_{\mathcal{U}}^{\mathcal{C}, \lambda}(\{\{\hat{x}_{i,t}\}_{i \in \mathcal{C}}, (c_{1,r}, c_{2,r} \dots c_{m,r})\})$ be a random variable representing the joint view of the users in \mathcal{U} . This suggests that all these honest but curious clients learned was the aggregation of gradients of all clients and their own gradients.

Functionality \mathcal{F}_s
Parties: users $1, \dots, N$ from \mathcal{S}_t and two servers S_0 and S_1 .
Parameters: corrupted rate η , number of participating training clients per-round n .
<ul style="list-style-type: none"> • \mathcal{F}_s receives a set of corrupted parties \mathcal{C} from the adversary \mathcal{A}, where the $\mathcal{C} / \mathcal{S}_t \leq \eta$. • For each round t: <ol style="list-style-type: none"> 1. \mathcal{F}_s receives a set of N clients \mathcal{S}_t and updates $x_{i,t}$ from clients $i \in \{\mathcal{S}_t \setminus \mathcal{C}\}$. 2. \mathcal{F}_s sends \mathcal{S}_t to \mathcal{A} and requests a set M_t. \mathcal{F}_s computes the $X_t = \bigodot_{i \in \{M_t \setminus \mathcal{C}\}} x_{i,t}$ if $M_t \subseteq \mathcal{S}_t$ and continues; otherwise \mathcal{F}_s sends abort to all honest participants. 3. There are two scenarios based on whether the servers are corrupted by \mathcal{A} as follows. <ul style="list-style-type: none"> – Corrupted: \mathcal{F}_s outputs X_t to all the participants corrupted by \mathcal{A}. – Not corrupted: \mathcal{F}_s requests a mask SK_t from \mathcal{A} and outputs $X_t \odot SK_t$ to S_1.

Figure 6: Ideal functionality for Janus.

C.3 SINGLE-ROUND SECURITY

Theorem 1 (Security of Janus) *Let the security parameter be λ and n be the number of users for aggregation in each round. Assuming the existence of secure underlying one-time pad, SHC, and public key encryption. Our generic construction can securely realize the ideal functionality \mathcal{F}_s under the presence of a static adversary controlling η fraction of n users (and the server S_1) as shown in Figure 6.*

$$\mathbf{REAL}_{\Pi, \mathcal{A}}^{\mathcal{F}_s, \mathcal{F}_{sum}^t}(\lambda, n, x_{\mathcal{S}_t}) \approx \mathbf{IDEAL}_{\mathcal{F}_s, \mathcal{S}}^{\mathcal{F}_{sum}^t}(\lambda, n, x_{M_t}). \quad (12)$$

Proof. We first prove the security of a single round aggregation. Our generic scheme (denoted as Π) securely realizes the ideal functionality \mathcal{F}_{sum}^t (Figure 7) in the random oracle model. We can find from the ideal function \mathcal{F}_{sum}^t that it is the M_t sent by the adversary \mathcal{A} that determines the actual result. We assume the \mathcal{A} controls a set of clients and denote the set of corrupted clients as \mathcal{C} .

Event 1. We start with the servers not being corrupted by the \mathcal{A} . Now, we first build a simulator \mathcal{S} in the ideal world, running \mathcal{A} as a subroutine. Specifically, the simulation for round t is as follows.

1. \mathcal{S} receives a set M_t from the adversary \mathcal{A} .
2. \mathcal{S} acquires Z_t from the \mathcal{F}_{sum}^t .
3. *Masking and Report.* \mathcal{S} interacts with \mathcal{A} as in the masking and report phase and acts as honest users in $i \in \{M_t \setminus \mathcal{C}\}$ with the masked updates $x'_{i,t}$ such that the $Z_t = \bigodot_{i \in \{M_t \setminus \mathcal{C}\}} x'_{i,t}$. Here the input update $x'_{i,t}$ and the mask $sk_{i,t}$ are generated by \mathcal{S} itself.
4. *Collection and Aggregation.* In this phase, \mathcal{S} interacts with \mathcal{A} , where \mathcal{A} performs as a honest participant as in the collection and aggregation of Π .
5. *UnMasking and Verification.* \mathcal{S} interacts with \mathcal{A} as honest participants in the unmasking and verification phase.
6. In the above steps, if all honest participants would abort in the protocol in this round of aggregation, then \mathcal{S} sends **abort** to \mathcal{F}_{sum}^t . Finally, \mathcal{A} outputs the value at random and terminates this aggregation.

972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025

Functionality \mathcal{F}_{sum}^t
Parties: users from \mathcal{S}_t and two servers.
Parameters: corrupted rate η .
<ul style="list-style-type: none"> • \mathcal{F}_{sum}^t receives a set of corrupted participants \mathcal{C} from the adversary \mathcal{A} and $x_{i,t}$ from clients $i \in \{\mathcal{S}_t \setminus \mathcal{C}\}$. • \mathcal{F}_s sends \mathcal{S}_t to \mathcal{A} and requests a set M_t. \mathcal{F}_s computes the $Z_t = \bigodot_{i \in \{M_t \setminus \mathcal{C}\}} x_{i,t}$ if $M_t \subseteq \mathcal{S}_t$ and continues; otherwise \mathcal{F}_s sends abort to the all honest participants. • For each round t: <ul style="list-style-type: none"> 1. \mathcal{F}_s receives a set of N clients \mathcal{S}_t and updates $x_{i,t}$ from clients $i \in \{\mathcal{S}_t \setminus \mathcal{C}\}$. 2. \mathcal{F}_s sends \mathcal{S}_t to \mathcal{A} and requests a set M_t. \mathcal{F}_s computes the $Z_t = \bigodot_{i \in \{M_t \setminus \mathcal{C}\}} x_{i,t}$ if $M_t \subseteq \mathcal{S}_t$ and continues; otherwise \mathcal{F}_s sends abort to the all honest participants. 3. There are two scenarios based on whether the servers are corrupted by \mathcal{A} as follows. <ul style="list-style-type: none"> – Corrupted: \mathcal{F}_s outputs X_t to all the participants corrupted by \mathcal{A}. – Not corrupted: \mathcal{F}_s requests a mask SK_t from \mathcal{A} and outputs $Z_t \odot SK_t$ to S_1.

Figure 7: Ideal functionality for Report and Collection in Round t .

We construct a series of hybrid execution programs from the real world to the ideal world.

Hybrid 1. The view of \mathcal{A} in the real-world execution is the same as the ideal world, when \mathcal{S} has actual inputs from honest participants $\{x_{i,t}\}, i \in \mathcal{S}_t \setminus \mathcal{C}$, the individual masks $sk_{i,t}$ and the SK_t .

Hybrid 2. \mathcal{S} does not use the actual masks in one-time pad between honest participants. It generates a random mask $sk'_{i,t}$ from the $\{0, 1\}^\lambda$, then it computes the corresponding one-time pad ciphertext as $(\hat{x}'_{i,t})$. We argue the view of \mathcal{A} in this hybrid is the computationally indistinguishable from the previous hybrid 1 as follows.

Firstly, the mask $sk_{i,t}$ is computed from the \mathcal{R}_C of the one-time pad, and the mask $sk'_{i,t}$ is randomly sampled in the ideal world. Let the M_t denotes the set of users chosen by \mathcal{A} in the ideal world. \mathcal{A} in the ideal and real world can observe $\text{Masking}(x_{i,t}, sk_{i,t})$ between a user $i \notin M_t$ and a client $i \in M_t$. This indistinguishability stems from the selection of random masks in the specific underlying one-time pad. Secondly, \mathcal{A} can observe the ciphertexts generated from the $sk'_{i,t}$. The distribution of the ciphertexts is computationally indistinguishable from the \mathcal{A} observed from the real world, which is depend on the security of the underlying OTP.

Hybrid 3. In this hybrid, instead of using one-time pad with actual personal mask $sk_{i,t}$ randomly selected from the space \mathcal{R}_C , \mathcal{S} uses masks randomly sampled from $\{0, 1\}^\lambda$. Before the proof, we model the generation of masks as a random oracle \mathcal{O}_R (see more details in the prior work (Bonawitz et al., 2017)). For $\forall i \in \{M_t \setminus \mathcal{C}\}$, the \mathcal{S} samples $sk'_{i,t}$ randomly and programs \mathcal{O}_R as $sk'_{i,t} = \hat{x}_{i,t} \odot x_{i,t}$, where the $\hat{x}_{i,t}$ is observed in the real world and the \odot denotes the inverse operation of \bigodot . From the perspective of \mathcal{A} , the distributions of $\hat{x}_{i,t}$ in this hybrid and the previous one are statistically indistinguishable.

Additionally, \mathcal{A} learns the $sk'_{i,t}$ in the clear for $i \in M_t$ in the real and ideal world. The distributions of $sk'_{i,t}$ are identical. However, \mathcal{A} learn nothing about $sk'_{i,t}$ for $i \notin M_t$ in both worlds because of the semantic security of the underlying one-time pad. From the view of \mathcal{A} , this hybrid is computationally indistinguishable from the previous hybrid.

Hybrid 4. In this hybrid, instead of control the random oracle as in the previous hybrid, \mathcal{S} will program the random oracle \mathcal{O}_R as $sk'_{i,t} = \hat{x}'_{i,t} \odot x'_{i,t}$. Specifically, the $x'_{i,t}$ s are chosen such that $\bigodot_{i \in \{M_t \setminus \mathcal{C}\}} x_{i,t} = \bigodot_{i \in \{M_t \setminus \mathcal{C}\}} x'_{i,t}$. From the view of \mathcal{A} this hybrid is the same as the previous one, which can be derived from Lemma 6.1 of the prior work (Bonawitz et al., 2017).

Hybrid 5. Similar to the previous operation, this hybrid replaces the mask of honest participants with the result from the random oracle. \mathcal{S} will abort if the \mathcal{A} would cheat by sending invalid masked updates to \mathcal{S} . In the phase of unmasking and verification, the \mathcal{A} would cheat by sending different M_t

1026 to \mathcal{F}_{sum}^t . \mathcal{S} will simulate the following protocol (see as Lemma 1) and output whatever the protocol
 1027 outputs. It is identical to the previous hybrid by doing this.

1028 The final hybrid precisely represents the execution of the ideal world. The aforementioned events
 1029 indicate that our system is secure in the ideal world with a single round process.

1030 **Event 2.** In this event, the server is not corrupted by \mathcal{A} , the whole simulation is the same as Event 1,
 1031 except that the \mathcal{S} will program the masks added by the \mathcal{A} in each step.

1032 We complete the proof that for any single round t , the protocol Π always securely realizes the ideal
 1033 functionality \mathcal{F}_{sum}^t in the presence of a static malicious adversary.

1034 **Lemma 1** *Assume there exists a PKI and a secure signature scheme, there are 3ζ participants with
 1035 at most ζ colluding malicious participants. Specifically, each party has an input bit of 0 or 1 from
 1036 a server. There exists a one-round protocol enabling each honest participant to determine whether
 1037 the server sent the same value to all honest participants.*

1038 **Proof.** When an honest participant receives at least 2ζ messages with the same value, it indicates that
 1039 the server has sent the same value to all honest participants. This is because, in the given system,
 1040 the threshold of 2ζ identical messages can only be met if a large majority of honest participants
 1041 have received the same value. Specifically, let the total number of participants in the system be
 1042 $n = t_h + t_m$, where t_h denotes the number of honest participants and t_m denotes the number of
 1043 malicious participants. For security and consistency in distributed protocols, the parameter ζ is set
 1044 such that $t_h > 2\zeta$. When an honest participant receives no fewer than 2ζ identical messages, it can
 1045 confidently conclude that at least $\zeta + 1$ of these messages were sent by distinct honest participants,
 1046 ensuring consistency of the message content. Hence, it can be inferred that the server has broadcast
 1047 the same value to all honest participants.

1048 Conversely, if an honest participant receives fewer than 2ζ messages with the same value, this sug-
 1049 gests that the server may have sent different messages to different participants during the commu-
 1050 nication process. Since the number of identical messages received by honest participants falls short
 1051 of forming a consensus of 2ζ , it implies that the server may have engaged in malicious behavior by
 1052 sending inconsistent messages to various honest participants. To ensure the security and consistency
 1053 of the protocol, in such a scenario, the honest participant will abort the protocol execution. This
 1054 abort mechanism effectively prevents potential security threats and data integrity issues that could
 1055 arise due to inconsistent messages from the server.

1056 C.4 MULTI-ROUND SECURITY

1057 Our threat model assumes the corrupted rate is η , which means that \mathcal{A} controls ηn clients throughout
 1058 total T rounds. In order to prove the security of the multi-round scheme on the basis of the above
 1059 single-round security proof. The mask $sk_{i,t}$ computed from \mathcal{O}_R of the underlying SHC Π_C . Let
 1060 the Δ_t be distribution of the view of \mathcal{A} in the single round t and the total number of rounds needed
 1061 for the model to converge is T . If there exists an adversary \mathcal{B} , and two rounds of aggregation
 1062 $t_1, t_2 \in [0, T]$, where \mathcal{B} can distinguish between Δ_{t_1} and Δ_{t_2} , then we can construct an adversary
 1063 \mathcal{A} breaks the security of the underlying Π_C . We call the challenger in the security game of Π_C as \mathcal{S} .
 1064 Specifically, there exists two worlds ($b = 0$ or 1) for the \mathcal{O}_R game. The \mathcal{S} uses a random function if
 1065 the $b = 0$. When $b = 1$, \mathcal{S} actual Π_C . Then we build the \mathcal{A} as follows. On input t_1, t_2 from \mathcal{B} , the \mathcal{A}
 1066 asks for $sk_{i,t}$ for all honest participants in the round t_1 and t_2 . Then \mathcal{A} could computes the masked
 1067 updates from the Π prescribed. It generates two views $\Delta_{t_1}, \Delta_{t_2}$ and sends them to the \mathcal{B} . Finally, \mathcal{A}
 1068 outputs whatever the \mathcal{B} outputs as the answer.

1069 C.5 RESISTING MIA

1070 The MIA is effective primarily because the server is aware of the final aggregated result. If the server
 1071 can manipulate the parameters sent to different clients, it can introduce inconsistencies that influ-
 1072 ence the model training process. The key to resisting this attack is to ensure that all clients start with
 1073 the same initial model parameters. This uniformity can be achieved through two main approaches:
 1074 using a public bulletin board where the initial parameters are posted for everyone to see, or through
 1075 mutual agreement among clients to verify that the parameters they receive are indeed consistent

across the network. The public bulletin board approach suffers from centralized dependency, information leakage, and scalability issues, while the mutual agreement method has high communication complexity, scalability limitations, and is vulnerable to Sybil attacks. Both methods face challenges in maintaining consistency and security as the number of clients increases.

A significant advancement in Janus, which brought forward a novel concept: making the aggregation results in the system visible only to the clients. In this approach, the computation of the final aggregation result is performed locally by each client, rather than centrally by the server. This means that even if the server, denoted as S_0 , disseminates inconsistent model parameters to different clients, it remains unaware of the actual final aggregated model. This paradigm shift ensures that the server cannot gain insight into the final result, thus preventing it from introducing systematic biases.

Additionally, we assume that S_0 and the clients are not colluding. In other words, the server cannot conspire with any client to manipulate the aggregation process. By decentralizing the aggregation computation and keeping the final result private among the clients, the Janus effectively mitigates the risk of a successful MIA. This approach not only enhances the security of the federated learning framework but also reinforces the privacy and trustworthiness of the system by limiting the server’s influence over the final model.

D APPENDIX FOR REBUTTAL REVISION

D.1 MORE COMPARISON METHODS

Table 3 demonstrates that Janus surpasses other state-of-the-art schemes in terms of security, efficiency, and functionality. Specifically, our scheme achieves optimal efficiency while providing enhanced security and functionality. Our scheme makes weaker assumptions compared to ELSA Rathee et al. (2023), resulting in higher security while supporting multi-round aggregation with a significant performance improvement. Compared to VeriFL Guo et al. (2021), Janus does not require constructing complex communication graphs or performing time-consuming secret sharing operations, which leads to substantial performance gains.

Table 3: Comparison of SA Constructions

Scheme	Input Privacy	Multi-round	Verifiability	Dynamic	Versatility	NS*	Efficiency†	MIA
SecAgg (Bonawitz et al., 2017)	✓	✗	✗	✗	✗	1	○	✗
BBSA (Bell et al., 2020)	✓	✗	✗	✗	✗	1	◐	✗
VeriFL Guo et al. (2021)	✓	✗	✓	✗	✗	1	◑	✗
ELSA Rathee et al. (2023)	✓	✗	✗	✗	✗	2	◒	✗
Flamingo (Ma et al., 2023)	✓	✓	✗	✗	✗	2†	◓	✗
Janus	✓	✓	✓	✓	✓	2	◔	✓

✓ Support, ✗ No support. Versatility: A generic construction. * Number of servers. † The decryptors of this construction can be abstracted to a server. ‡ More black parts in the circle indicate better efficiency, and the theoretical support comes from the computation efficiency analysis in Table 2.

D.2 ADDITIONAL EXPERIMENTS

Our scheme is not limited to specific models or datasets. To better support this conclusion, we have added more experimental results for the CIFAR-100 dataset on different models. Specifically, the CIFAR-100 dataset is a challenging benchmark dataset widely used in machine learning and computer vision research. It contains 60,000 color images, each of size 32x32 pixels, distributed across 100 distinct classes. Each class is organized hierarchically, with 20 superclasses grouping the 100 fine-grained categories, adding a layer of complexity. This fine-grained nature, combined with the small image resolution, makes classification tasks on CIFAR-100 particularly difficult, as it demands models to capture subtle features and patterns. The dataset is balanced, with each class containing 500 training images and 100 test images, ensuring uniform representation while amplifying the difficulty of distinguishing between visually similar categories. The specific experimental results are shown in Figure 8 and 9, which show that Janus is comparable to most of the existing schemes in terms of performance, but Janus has an obvious advantage in terms of computational overhead.

1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187

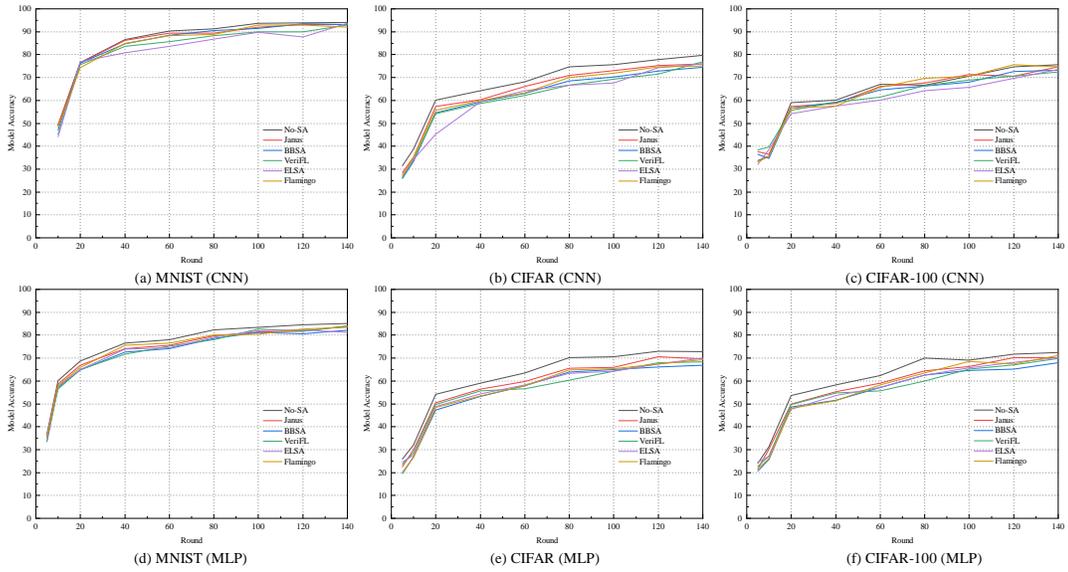


Figure 8: Test accuracy across different datasets and models.

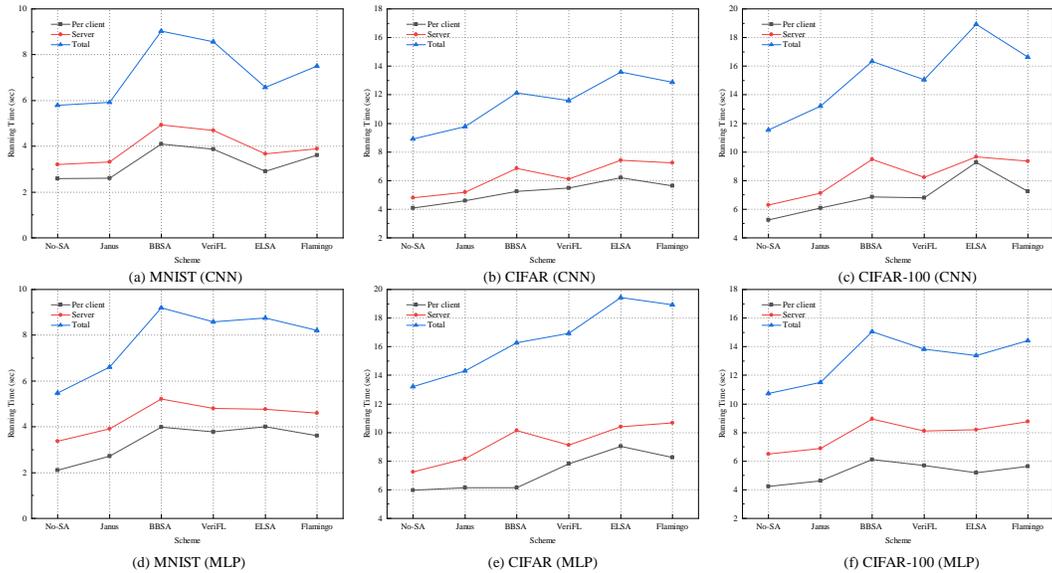


Figure 9: Computation overhead across different datasets and models.