OmniFC: Rethinking Federated Clustering via Lossless and Secure Distance Reconstruction

Jie Yan Jing Liu Zhong-Yuan Zhang*

Central University of Finance and Economics jieyan@email.cufe.edu.cn, liu_jing0623@163.com, zhyuanzh@gmail.com

Abstract

Federated clustering (FC) aims to discover global cluster structures across decentralized clients without sharing raw data, making privacy preservation a fundamental requirement. There are two critical challenges: (1) privacy leakage during collaboration, and (2) robustness degradation due to aggregation of proxy information from non-independent and identically distributed (Non-IID) local data, leading to inaccurate or inconsistent global clustering. Existing solutions typically rely on model-specific local proxies, which are sensitive to data heterogeneity and inherit inductive biases from their centralized counterparts, thus limiting robustness and generality. We propose Omni Federated Clustering (OmniFC), a unified and modelagnostic framework. Leveraging Lagrange coded computing, our method enables clients to share only encoded data, allowing exact reconstruction of the global distance matrix—a fundamental representation of sample relationships—without leaking private information, even under client collusion. This construction is naturally resilient to Non-IID data distributions. This approach decouples FC from model-specific proxies, providing a unified extension mechanism applicable to diverse centralized clustering methods. Theoretical analysis confirms both reconstruction fidelity and privacy guarantees, while comprehensive experiments demonstrate OmniFC's superior robustness, effectiveness, and generality across various benchmarks compared to state-of-the-art methods. Code will be released.

1 Introduction

Traditional clustering methods presuppose centralized access to the entire dataset, enabling the construction of global structures such as cluster centroids or kernel matrices. However, in federated settings characterized by data fragmentation across clients and privacy constraints, this assumption breaks down, precluding direct application.

To overcome this, federated clustering (FC) [1] has emerged, enabling clients to collaboratively group data without sharing raw samples, and has found applications in client selection and exploratory data analysis. There are two fundamental challenges: (1) privacy leakage during collaboration, and (2) robustness degradation under non-independent and identically distributed (Non-IID) data. Existing

^{*}Corresponding author.

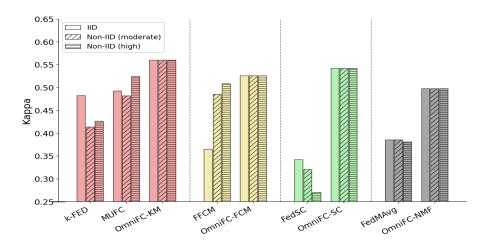


Figure 1: **Robustness to heterogeneity.** We employ COIL-100 and 100 clients to compare the proposed OmniFC with the federated extensions of centralized clustering methods [1, 3, 2, 6, 7]. Compared to existing one-to-one extensions, OmniFC not only unifies the extension of centralized clustering methods but also achieves superior robustness and effectiveness.

FC methods approximate global structures by aggregating model-specific local proxies: federated k-means (KM) and fuzzy c-means (FCM) aggregate local cluster centroids [1, 2, 3, 4, 5], federated spectral clustering (SC) [6] reconstructs the global kernel matrix from local low-rank factors, and federated non-negative matrix factorization (NMF) [7] aggregates local basis matrices. These proxies, however, are computed from biased client-specific datasets, fail to reliably capture global structures, leading to degraded robustness and performance (Fig. 1). Moreover, such methods are tightly bound to specific centralized clustering methods, inheriting restrictive inductive biases—e.g., data compactness in KM [8] and FCM [9], data connectivity in SC [10], and low-rank representation in NMF [11]—thereby confining their performance to assumption-compliant data and limiting their generality.

This work addresses both limitations through a unifying perspective: reconstructing the global pairwise distance matrix, which offers a model-agnostic and fundamental representation of sample relationships, naturally resilient to the Non-IID problem. The key challenge, however, lies in securely computing this matrix without exposing private data. To this end, we propose Omni Federated Clustering (OmniFC), a novel framework that facilitates a unified extension from centralized clustering to FC through lossless and secure distance reconstruction. OmniFC comprises three main steps: local Lagrange-encoded sharing, global distance reconstruction, and cluster assignment. Each client initially encrypts its local data using Lagrange coded computing [12], shares the encoded data with peers for pairwise distance computation, and subsequently transmits the resulting distances to the central server for constructing the global distance matrix. Finally, the global distance matrix can serve as input to centralized clustering methods for performing cluster assignment. Fig. 1 demonstrates the superiority of OmniFC. With respect to distance reconstruction, the proposed OmniFC exhibits two salient features: 1) Efficacy. Both theoretical and empirical analyses consistently demonstrate the capability for lossless reconstruction and robustness to the Non-IID problem. Benefiting from this, the proposed OmniFC achieves lossless federated extensions for pairwise-distance-dependent methods (e.g., SC) and enhances federated extensions for methods (e.g., KM) without explicit dependence on pairwise distances. 2) Security. Theoretical analysis demonstrates that the privacy of local data is preserved during data sharing, as the encoded data prevents the inference of private information even under client collusion. In summary, our contributions are threefold:

- 1) We propose OmniFC, a novel framework that facilitates a unified extension from centralized clustering to FC through lossless and secure distance reconstruction.
- 2) We establish theoretical assurances regarding the efficacy and security of distance reconstruction.
- 3) Experimental results show that our OmniFC outperforms SOTA methods on various benchmarks.

2 Related Work

Centralized Clustering. Traditional centralized clustering aggregates client-held local data on a central server for grouping, with methods making different assumptions—such as compactness [8, 13], connectivity [10, 14], density [15, 16], hierarchy [17, 18], and low-rank representation [11, 19] of the data distribution—to adapt to diverse datasets. However, these methods may become inapplicable due to privacy constraints that prevent the centralization of client data.

Federated Clustering (FC). Unlike centralized clustering, which requires collecting raw client data for model training, FC collects local proxies instead, thus strengthening user privacy protection. To handle this, several recent works have shifted from sharing local private data to exchanging local cluster centroids [1, 2, 4, 5], local basis matrices [7] or synthetic data [20]. Although these methods show promise, these methods either suffer from performance degradation caused by the Non-IID problem or achieve gains at the expense of privacy [20].

Secure FC. Secure FC leverages advanced privacy-preserving techniques—including differential privacy [21], machine unlearning [22], and Lagrange coded computing [12]—to concurrently improve clustering efficacy and fortify data confidentiality. Existing methods typically focus on the effective and secure construction of either global cluster centroids for k-means [23, 3, 24, 25, 26] or a global kernel matrix for spectral clustering [6]. Although promising, these methods remain limited by the Non-IID problem or fail to offer a model-agnostic solution. Moreover, they inherently retain assumptions—such as data compactness [8, 13] and connectivity [10, 14]—from their centralized counterparts, limiting their effectiveness to compliant datasets and thereby reducing their practical applicability.

The most closely related work is SecFC [23], which also leverages Lagrange coded computing to improve clustering accuracy while preserving data confidentiality. In comparison, the proposed OmniFC exhibits the following distinctive characteristics: 1) *Generality*. As a federated variant of k-means (KM), SecFC naturally inherits KM's assumption of data compactness, limiting their effectiveness to KM-friendly datasets. In contrast, the proposed unified and model-agnostic framework OmniFC accommodates diverse datasets by extending beyond KM to encompass alternatives such as spectral clustering and DBSCAN. 2) *One-shot communication scheme*. To obtain more accurate cluster centroids, SecFC requires multiple rounds of communication between the clients and the server. However, multi-round training is unfeasible in some scenarios, like model markets, where users can solely purchase pre-trained models [27]. In contrast, the proposed OmniFC, requiring merely one communication round, exhibits enhanced applicability in such scenarios.

3 Omni Federated Clustering (OmniFC)

This section begins with an overview of the problem definition and the OmniFC framework, followed by a detailed description of OmniFC, and concludes with its privacy analysis and complexity analysis.

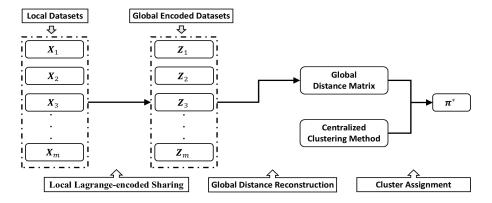


Figure 2: An overview of the proposed OmniFC. The architecture comprises three main steps: 1) Local Lagrange-Encoded Sharing. Each client j ($j \in [m]$) encodes its private data using Lagrange polynomial interpolation and distributes the encoded data to all peers, enabling each client to construct a global encoded dataset while preserving data privacy. 2) Global Distance Reconstruction. Each client j computes pairwise distances within its global encoded dataset and transmits the results to the central server, which leverages them to reconstruct the global distance matrix. 3) Cluster Assignment. A centralized clustering method (e.g., k-means) is applied to the global distance matrix to produce the final clustering result π^* .

3.1 Overview

Problem Definition. Consider a real world dataset $X \in \mathbb{R}^{n \times d}$ comprising n d-dimensional samples $\{x_i\}_{i=1}^n$, which are distributed among m clients, i.e., $X = \bigcup_{j=1}^m X_j$. FC aims to partition X into k clusters with high intra-cluster similarity and low inter-cluster similarity while retaining X_j $(j \in [m] = \{1, 2, \cdots, m\})$ locally. A more detailed summary of notations is presented in Table 3 of the appendix.

Framework Overview. As shown in Fig. 2, OmniFC comprises three main steps: local Lagrange-encoded sharing, global distance reconstruction, and cluster assignment. Each client j ($j \in [m]$) initially encrypts its local data using Lagrange coded computing (LCC) [12], shares the encoded data with peers for pairwise distance computation, and subsequently transmits the resulting distances to the central server for constructing the global distance matrix. Finally, the global distance matrix can serve as input to centralized clustering methods for performing cluster assignment.

3.2 OmniFC

Local Lagrange-encoded Sharing. First, each sample $x_i \in X$ $(i \in [n])$ —regardless of the client to which it is distributed to—is independently transformed from the real domain \mathbb{R}^d to the finite field \mathbb{F}_p^d to ensure numerical stability in secure computation [12], with p denoting a prime. The transformation is defined as:

$$\tilde{\boldsymbol{x}}_i = round(2^q \cdot \boldsymbol{x}_i) + p \cdot \frac{|sign(\boldsymbol{x}_i)| - sign(\boldsymbol{x}_i)}{2},$$
 (1)

where $q \in \mathbb{Z}$ regulates the quantization loss. $round(\cdot)$ and $sign(\cdot)$ represent element-wise rounding and sign functions, respectively. Rounding discretizes continuous values to ensure finite field compatibility, while the sign function facilitates correct mapping of negative values [28]. We denote the transformed form of $X \in \mathbb{R}^{n \times d}$ as $\tilde{X} \in \mathbb{F}_p^{n \times d}$.

Then, each sample $\tilde{x}_i \in \tilde{X}$ $(i \in [n])$ is independently encoded via Lagrange polynomial interpolation by the client (Fig. 3), enabling secret sharing among clients. Specifically, $\tilde{x}_i \in \mathbb{F}_p^d$ is partitioned

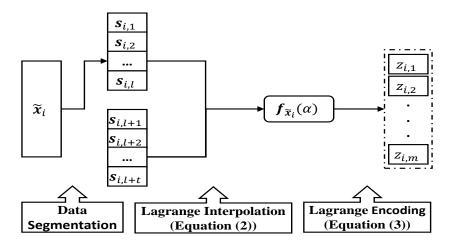


Figure 3: An illustration of the Lagrange encoding. Each sample \tilde{x}_i ($i \in [n]$) is initially divided into l segments $\{s_{i,o}\}_{o=1}^l$. Incorporating t additional noises $\{s_{i,l+o}\}_{o=1}^t$, Lagrange interpolation is then conducted as per Equation (2) to yield $f_{\tilde{x}_i}(\alpha)$. Subsequently, the encoded representations $\{z_{i,j}\}_{j=1}^m$ of \tilde{x}_i are computed according to Equation (3).

into l segments $\{s_{i,o}\}_{o=1}^l$, i.e., $\tilde{\boldsymbol{x}}_i = [s_{i,1}^T, s_{i,2}^T, \cdots, s_{i,l}^T]^T$, and combined with t random noises to construct a polynomial that serves to encode $\tilde{\boldsymbol{x}}_i$. The noises are introduced to ensure privacy protection against potential client collusion [12]. Assuming that d is divisible by l, the client holds data segments $s_{i,o} \in \mathbb{F}_p^{\frac{d}{l}}$ $(o \in [l])$, and samples t additional noises $s_{i,l+o}$ $(o \in [t])$ uniformly from $\mathbb{F}_p^{\frac{d}{l}}$. Based on the segments $\{s_{i,o}\}_{o=1}^{l+t}$, the Lagrange interpolation polynomial $f_{\tilde{\boldsymbol{x}}_i} : \mathbb{F}_p \to \mathbb{F}_p^{\frac{d}{l}}$ of degree l+t-1 can be constructed as follows:

$$f_{\tilde{\boldsymbol{x}}_i}(\alpha) = \sum_{o=1}^{l+t} \boldsymbol{s}_{i,o} \cdot \prod_{o' \neq o} \frac{\alpha - \alpha_{o'}}{\alpha_o - \alpha_{o'}},\tag{2}$$

where $\{\alpha_o\}_{o=1}^{l+t}$ denotes a set of l+t distinct hyperparameters from \mathbb{F}_p , pre-specified through agreement among all clients and the central server. Particularly, each data segment $s_{i,o}$ $(o \in [l])$ can be recovered by setting $\alpha = \alpha_o$, i.e., $f_{\tilde{\boldsymbol{x}}_i}(\alpha_o) = s_{i,o}$. Beyond the $\{\alpha_o\}_{o=1}^{l+t}$ employed in constructing the polynomial $f_{\tilde{\boldsymbol{x}}_i}$, all clients and the central server also pre-select m distinct public hyperparameters $\{\beta_j\}_{j=1}^m$ for encoding, where $\beta_j \in \mathbb{F}_p$ and $\{\alpha_o\}_{o=1}^{l+t} \cap \{\beta_j\}_{j=1}^m = \emptyset$. Based on $\{\beta_j\}_{j=1}^m$, the client encodes its local data $\tilde{\boldsymbol{x}}_i$ into m distinct representations $\{\boldsymbol{z}_{i,j}\}_{j=1}^m$ for secret sharing, with each representation

$$\boldsymbol{z}_{i,j} = \boldsymbol{f}_{\tilde{\boldsymbol{x}}_i}(\beta_j) \tag{3}$$

delivered to the j-th client.

As these operations are defined per sample, they are universally applicable to local data across all clients. Hence, each client j $(j \in [m])$ will possess a global encoded dataset $\boldsymbol{Z}_j \in \mathbb{F}_p^{n \times \frac{d}{l}}$ corresponding to $\tilde{\boldsymbol{X}} \in \mathbb{F}_p^{n \times d}$, where $\boldsymbol{Z}_j = [\boldsymbol{z}_{1,j}, \boldsymbol{z}_{2,j}, \cdots, \boldsymbol{z}_{n,j}]^T = [\boldsymbol{f}_{\tilde{\boldsymbol{x}}_1}(\beta_j), \boldsymbol{f}_{\tilde{\boldsymbol{x}}_2}(\beta_j), ..., \boldsymbol{f}_{\tilde{\boldsymbol{x}}_n}(\beta_j)]^T$.

Global Distance Reconstruction. For each client j ($j \in [m]$), pairwise distances between all encoded representations $z_{i,j}$ and $z_{i',j}$ in Z_j ($i,i' \in [n]$) are calculated and subsequently sent to the central server for constructing the global distance matrix. Specifically, the pairwise distance between $z_{i,j}$ and $z_{i',j}$ can be calculated as:

$$dis(\mathbf{z}_{i,j}, \mathbf{z}_{i',j}) = \|\mathbf{z}_{i,j} - \mathbf{z}_{i',j}\|_{2}^{2}.$$
 (4)

Based on the m distances $\{dis(\boldsymbol{z}_{i,j}, \boldsymbol{z}_{i',j})\}_{j=1}^{m}$ provided by the clients, the server can accurately recover the pairwise distance between the corresponding samples $\tilde{\boldsymbol{x}}_i$ and $\tilde{\boldsymbol{x}}_{i'}$, as demonstrated in Theorem 1.

Theorem 1. Let $f_{\tilde{x}_i,\tilde{x}_{i'}}(\beta): \mathbb{F}_p \to \mathbb{F}_p$ denote the Lagrange interpolation polynomial interpolated from the set $\{(\beta_j, dis(\mathbf{z}_{i,j}, \mathbf{z}_{i',j}))\}_{j=1}^m$:

$$\boldsymbol{f}_{\tilde{\boldsymbol{x}}_{i},\tilde{\boldsymbol{x}}_{i'}}(\beta) = \sum_{j=1}^{m} dis(\boldsymbol{z}_{i,j}, \boldsymbol{z}_{i',j}) \cdot \prod_{j' \neq j} \frac{\beta - \beta_{j'}}{\beta_{j} - \beta_{j'}}, \tag{5}$$

where $z_{i,j}$ and $z_{i',j}$ denote the encoded representations of arbitrary samples \tilde{x}_i and $\tilde{x}_{i'}$ distributed among clients. When $m \geq 2l + 2t - 1$, the distance $dis(\tilde{x}_i, \tilde{x}_{i'})$ can be precisely recovered:

$$dis(\tilde{\boldsymbol{x}}_{i}, \tilde{\boldsymbol{x}}_{i'}) = \sum_{o=1}^{l} \boldsymbol{f}_{\tilde{\boldsymbol{x}}_{i}, \tilde{\boldsymbol{x}}_{i'}}(\alpha_{o}), \tag{6}$$

irrespective of how data is distributed among clients.

Remark 1. The condition $m \geq 2l + 2t - 1$ imposes minimal practical constraint, given that m is predefined by the system while l and t are tunable hyperparameters. This flexibility allows the condition to be met easily, ensuring the theorem's practical applicability and highlighting its relevance to real-world implementations.

Then, by converting $dis(\tilde{x}_i, \tilde{x}_{i'})$ from the finite field \mathbb{F}_p back to the real domain \mathbb{R} , the server recovers:

$$dis(\boldsymbol{x}_{i}, \boldsymbol{x}_{i'}) = \begin{cases} \frac{1}{2^{q}} \cdot dis(\tilde{\boldsymbol{x}}_{i}, \tilde{\boldsymbol{x}}_{i'}) & \text{if} \quad 0 \leq dis(\tilde{\boldsymbol{x}}_{i}, \tilde{\boldsymbol{x}}_{i'}) < \frac{p-1}{2} \\ \frac{1}{2^{q}} \cdot (dis(\tilde{\boldsymbol{x}}_{i}, \tilde{\boldsymbol{x}}_{i'}) - p) & \text{if} \quad \frac{p-1}{2} \leq dis(\tilde{\boldsymbol{x}}_{i}, \tilde{\boldsymbol{x}}_{i'}) < p \end{cases}$$
(7)

Based on the recovered distances, we denote the global distance matrix as $D \in \mathbb{R}^{n \times n}$, with each entry defined as $D_{ii'} = dis(\mathbf{x}_i, \mathbf{x}_{i'})$ for $i, i' \in [n]$.

Cluster Assignment. With the recovered global distance matrix $D \in \mathbb{R}^{n \times n}$, the server can directly perform clustering without requiring any modification to existing centralized clustering methods. This characteristic demonstrates the **simplicity** and **flexibility** of the proposed OmniFC framework.

Specifically, pairwise-distance-dependent centralized clustering methods—such as spectral clustering (SC) [10], DBSCAN [15], hierarchical clustering (HC) [17], and k-medoids (KMed) [29]—can seamlessly utilize D for model training, owing to their intrinsic reliance on pairwise sample distances during the clustering process. For methods that do not explicitly depend on pairwise relationships—such as k-means (KM) [8], fuzzy c-means (FCM) [9], and nonnegative matrix factorization (NMF) [11]—the server employs D as a proxy for the raw features $X \in \mathbb{R}^{n \times d}$ to perform clustering, i.e., the distance values become the new features of the samples. This allows these algorithms to operate as if on centralized data, while implicitly leveraging the global structure encoded in D. These federated extensions of centralized methods built upon OmniFC are denoted as OmniFC-SC, OmniFC-DBSCAN, OmniFC-HC, OmniFC-KMed, OmniFC-KM, OmniFC-FCM, and OmniFC-NMF, respectively. Algorithm 1 in the appendix delineates the pseudocode of OmniFC.

3.3 Privacy Analysis

OmniFC adopts LCC encryption to enhance clustering performance while fortifying data privacy. Although LCC enables clients to obtain global awareness via inter-client sharing of Lagrange-encoded data, it also poses emerging privacy threats, as colluding clients may leverage the shared information to infer others' private data [12]. Hence, evaluating OmniFC's resistance to client collusion is

essential for delineating its practical applicability. Theorem 2 provides a formal guarantee that each data point \tilde{x}_i maintains information-theoretic security in the presence of up to t colluding clients, thereby affirming the practical applicability of OmniFC.

Theorem 2. Given the number of noises t, a t-private OmniFC is achievable if $m \ge 2l + 2t - 1$, i.e.,

$$I(\tilde{\boldsymbol{x}}_i; \{\boldsymbol{z}_{i,j}\}_{j \in \boldsymbol{C}}) = 0, \tag{8}$$

where $I(\cdot;\cdot)$ denotes the mutual information function, $\mathcal{C} \subset [m]$ and $|\mathcal{C}| \leq t$.

Remark 2. Mutual information essentially measures how much one piece of information reveals about another—when it equals zero, it means one reveals nothing about the other, thereby preserving privacy. Furthermore, mutual information is not an isolated privacy-preserving metric; it can be compared with other privacy measures (such as differential privacy) within a unified framework [30].

Remark 3. The condition for achieving t-private security in Theorem 2 coincides with that for exact distance reconstruction in Theorem 1, i.e., $m \ge 2l + 2t - 1$. Consequently, by adhering to this constraint, we can increase the number of noises t to strengthen privacy protection without compromising the precision of distance reconstruction.

3.4 Complexity Analysis

Recall that reconstructing the global pairwise distance matrix involves two main stages: local encoding and distance computation on the client side, followed by server-side decoding to recover the global matrix. For each client j, fast polynomial interpolation and evaluation [31] yield an encoding complexity of $\mathcal{O}(\frac{dn_j m \log^2 m}{l})$, while distance computation on the encoded data incurs a complexity of $\mathcal{O}(\frac{dn^2}{l})$, leading to an overall complexity of $\mathcal{O}(\frac{d(n^2+n_j m \log^2 m)}{l})$. n_j is the number of samples held by client j. For the server side, the decoding operation has a complexity of $\mathcal{O}(n^2(l+t)\log^2(l+t))$.

The analysis indicates that, on the client side, increasing the number of data segments l substantially reduces computational complexity, thereby facilitating the calculation of pairwise distances in a $\frac{d}{l}$ -dimensional space, where the feature count $(\frac{d}{l})$ involved in pairwise comparisons decreases as l rises. On the server side, however, a larger l results in greater computational complexity. As for the impact of l on the overall complexity of the reconstruction process, a theoretical analysis is challenging due to the differing hyperparameters involved on both sides.

To handle this, we presented the runtime for reconstructing the global pairwise distance matrix across varying data scales on $10x_73k$. As shown in Table 4, one can observe that: Although the runtime exhibits a marked increase as n grows, enhancing the number of data segments l can considerably boost computational efficiency when n is fixed. For instance, processing 70k samples takes 4352 seconds with l=8, demonstrating that our method remains computationally feasible on large-scale datasets. Importantly, as long as the constraints $m \geq 2(l+t-1)+1$ holds, increasing l have no effect on the reconstruction effectiveness of the global pairwise distance matrix (Theorem 1 and Fig. 5).

4 Experiments

4.1 Experimental Setup

Datasets and Evaluation Criteria. The proposed OmniFC is assessed using seven benchmark datasets across tabular, visual, temporal, and genomic domains, including Iris [32], MNIST [33], Fashion-MNIST [34], COIL-20 [35], COIL-100 [35], Pendigits [36], and 10x_73k [37]. The chosen datasets encompass diverse modalities, dimensionalities, and cluster patterns, facilitating a comprehensive evaluation of the method's generalizability in practical scenarios.

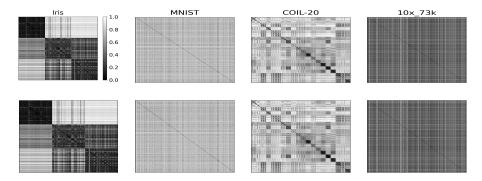


Figure 4: Comparison between the ground-truth (top row) and reconstructed (bottom row) pairwise distance matrices. The visual consistency indicates that the proposed OmniFC faithfully recovers the inter-sample similarity.

The evaluation criteria encompass Normalized Mutual Information (NMI) [38] and Kappa [39], with higher scores indicating improved clustering performance. Despite the widespread use of NMI, increasing evidence suggests it may be misleading, whereas Kappa is more reliable [39, 20, 40]. Hence, our analysis is grounded in Kappa-based results, with NMI-based outcomes relegated to the appendix for reference. Details of datasets and evaluation criteria are provided in Appendix C.1.

Baselines. OmniFC is evaluated in comparison with the federated extensions of several centralized clustering methods, including KM-based (k-FED [1], MUFC [3]), FCM-based (FFCM [2]), SC-based (FedSC [6]), and NMF-based (FedMAvg [7]) methods. To contextualize the performance of federated clustering against its centralized counterpart, we also present results of vanilla KM, FCM, SC, and NMF under centralized settings, referred to as KM_central, FCM_central, SC_central, and NMF_central, respectively.

Federated Settings. Following Ref. [41, 20], we simulate diverse federated settings by partitioning the real-world dataset into k^* subsets—each representing a client—and adjusting the non-IID level p, where k^* denotes the number of true clusters. Specifically, for each client, a fraction p of its data is sampled from a single cluster, while the remaining 1-p portion is drawn uniformly across all clusters. As such, p=0 recovers the IID setting, whereas p=1 induces a maximally skewed distribution, where each client's data is fully concentrated within a single cluster. Since OmniFC is immune to the Non-IID degree, the Non-IID level p is indicated solely during comparisons with the existing FC baselines and omitted elsewhere.

4.2 Experimental Results

Our experiments center on three key aspects: 1) the comparative advantage of OmniFC over existing approaches; 2) the generality of OmniFC in extending centralized clustering methods; and 3) the sensitivity of OmniFC to hyperparameters. Implementation details are provided in Appendix C.2, and supplementary experimental results are presented in Appendix D.

Efficacy Analysis. To comprehensively validate the efficacy of OmniFC, we simulate five scenarios per dataset: IID (p=0), mildly non-IID (p=0.25), moderately non-IID (p=0.5), highly non-IID (p=0.75), and fully non-IID (p=1). As shown in Table 1, the proposed OmniFC enables superior federated extensions for both pairwise-distance-dependent SC and methods that do not explicitly depend on pairwise relationships, such as KM, FCM, and NMF. For SC, our extended results attain centralized-level clustering fidelity while remaining robust to diverse Non-IID conditions, owing to

Table 1: **Kappa of clustering methods in different federated scenarios.** For each comparison, the best result is highlighted in boldface.

Dataset	p	SC-base	ed metho	ds	KM-	based m	ethods		FCM-base	d metho	ods	NMF-bas	MF-based method	
Dataset	P	SC_central	FedSC	Ours	KM_central	k-FED	MUFC	Ours	FCM_central	FFCM	Ours	NMF_central	FedMAvg	Ours
	0.00		0.95	0.95		0.38	0.83	0.95		0.96	0.95		0.50	0.95
	0.25		0.93	0.95		0.95	0.93	0.95		0.49	0.95		0.50	0.95
Iris	0.50	0.95	0.85	0.95	0.95	0.93	0.79	0.95	0.95	0.93	0.95	0.57	0.50	0.95
	0.75		0.93	0.95		0.95	0.81	0.95		0.96	0.95		0.50	0.95
	1.00		0.31	0.95		0.71	0.77	0.95		0.97	0.95		0.50	0.95
	0.00		0.53	0.55		0.43	0.41	0.42		0.48	0.41		0.40	0.38
	0.25		0.54	0.55		0.45	0.50	0.42		0.52	0.41		0.44	0.38
MNIST	0.50	0.55	0.54	0.55	0.47	0.29	0.46	0.42	0.50	0.53	0.41	0.46	0.39	0.38
	0.75		0.58	0.55		0.32	0.47	0.42		0.45	0.41		0.45	0.38
	1.00		0.38	0.55		0.47	0.43	0.42		0.48	0.41		0.46	0.38
	0.00		0.54	0.53		0.46	0.43	0.51		0.51	0.50		0.46	0.49
	0.25		0.52	0.53		0.43	0.40	0.51		0.47	0.50		0.46	0.49
Fashion-MNIST	0.50	0.53	0.54	0.53	0.50	0.48	0.50	0.51	0.53	0.43	0.50	0.51	0.46	0.49
	0.75		0.47	0.53		0.45	0.45	0.51		0.50	0.50		0.46	0.49
	1.00		0.38	0.53		0.32	0.50	0.51		0.46	0.50		0.46	0.49
	0.00		0.68	0.63		0.42	0.58	0.64		0.51	0.59		0.50	0.61
	0.25		0.68	0.63		0.46	0.61	0.64		0.47	0.59		0.51	0.61
COIL-20	0.50	0.61	0.73	0.63	0.64	0.42	0.57	0.64	0.59	0.51	0.59	0.56	0.44	0.61
	0.75		0.54	0.63		0.41	0.58	0.64		0.55	0.59		0.51	0.61
	1.00		0.29	0.63		0.46	0.56	0.64		0.59	0.59		0.52	0.61
	0.00		0.34	0.54		0.48	0.49	0.56		0.37	0.53		0.39	0.50
	0.25		0.32	0.54		0.45	0.50	0.56		0.38	0.53		0.38	0.50
COIL-100	0.50	0.54	0.32	0.54	0.49	0.41	0.48	0.56	0.49	0.49	0.53	0.43	0.39	0.50
	0.75		0.29	0.54		0.41	0.50	0.56		0.48	0.53		0.39	0.50
	1.00		0.27	0.54		0.43	0.52	0.56		0.51	0.53		0.38	0.50
	0.00		0.74	0.72		0.59	0.59	0.62		0.62	0.66		0.33	0.72
	0.25		0.73	0.72		0.46	0.58	0.62		0.61	0.66		0.33	0.72
Pendigits	0.50	0.72	0.72	0.72	0.61	0.48	0.60	0.62	0.66	0.53	0.66	0.45	0.33	0.72
	0.75		0.69	0.72		0.33	0.49	0.62		0.49	0.66		0.33	0.72
	1.00		0.52	0.72		0.53	0.62	0.62		0.70	0.66		0.33	0.72
	0.00		0.52	0.89		0.40	0.63	0.56		0.46	0.55		0.49	0.82
	0.25		0.52	0.89		0.55	0.63	0.56		0.47	0.55		0.49	0.82
10x_73k	0.50	0.89	0.52	0.89	0.85	0.57	0.62	0.56	0.53	0.72	0.55	0.88	0.49	0.82
	0.75		0.54	0.89		0.37	0.65	0.56		0.64	0.55		0.50	0.82
	1.00		0.24	0.89		0.30	0.79	0.56		0.64	0.55		0.50	0.82
count	-	-	8	27	-	2	9	24	-	13	22	-	5	30

lossless pairwise distance reconstruction, which remains unaffected by non-IID severity (see Theorem 1 and Figure 4). For centralized methods not explicitly reliant on pairwise relationships, our extended results generally match—and occasionally exceed—their performance under centralized settings, indicating that the global distance matrix \boldsymbol{D} can serve as an effective surrogate for the raw feature matrix \boldsymbol{X} to perform clustering.

Generality Analysis. To assess OmniFC's generalizability in extending centralized clustering methods, we integrate it with three additional methods (KMed, DBSCAN, and HC) that have been well-studied in centralized contexts but remain underexplored in federated settings. Like SC, all three methods perform clustering based on inter-sample pairwise distances. Hence, by utilizing OmniFC's lossless distance reconstruction, these three methods can be effortlessly integrated into the OmniFC framework to facilitate lossless federated extensions, as shown in Table 2.

Table 2: Kappa of different clustering methods.

Dataset	KMed-bas	ed methods	DBSCAN-l	based methods	HC-based methods		
Dutuset	Central	Ours	Central	Ours	Central	Ours	
Iris	0.94	0.94	0.50	0.50	0.94	0.95	
MNIST	0.31	0.30	0.21	0.21	0.41	0.41	
Fashion-MNIST	0.42	0.42	0.14	0.14	0.45	0.45	
COIL-20	0.41	0.42	0.58	0.58	0.45	0.45	
COIL-100	0.34	0.34	0.37	0.37	0.48	0.48	
Pendigits	0.44	0.45	0.48	0.48	0.57	0.57	
10x_73k	0.30	0.30	0.01	0.01	0.28	0.28	

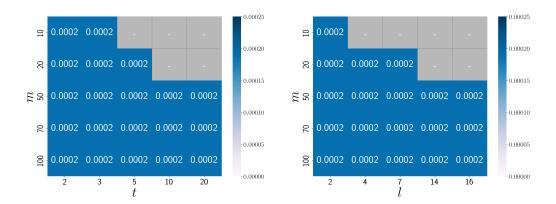


Figure 5: Hyperparameter sensitivity of the global distance matrix reconstruction loss. The gray-highlighted region denotes hyperparameter settings that violate the condition $m \ge 2l + 2t - 1$ in Theorem 1, thus precluding distance reconstruction.

Sensitivity Analysis. To assess the hyperparameter sensitivity of OmniFC, we measure the global distance matrix reconstruction loss—defined as the root-mean-square deviation (RMSE) between the ground-truth and reconstructed pairwise distance matrices—across varying number of clients (m), noises (t), and segments (l). In fact, a theoretical guarantee for this has already been provided in Theorem 1: as long as the condition $m \geq 2l + 2t - 1$ holds, OmniFC is capable of achieving accurate distance reconstruction. This theoretical result is further substantiated by the empirical evidence presented in Fig. 5.

5 Conclusion

This work introduces OmniFC, a unified and model-agnostic framework via lossless and secure distance reconstruction. Unlike existing methods that rely on model-specific proxies and suffer from data heterogeneity, OmniFC adopts a distance-based perspective that is decoupled from specific clustering models. Benefit from this, theoretical and empirical results show that this framework improves robustness under non-IID settings and supports the extension of a wide range of centralized clustering algorithms to FC.

Beyond FC, the proposed framework may open broader opportunities across federated learning. In particular, the reconstructed global distance matrix can naturally function as a global affinity graph, offering new possibilities for advancing federated graph learning and other domains where capturing global sample relationships is fundamental.

Acknowledgements

The authors are grateful to Yong Wang and the three anonymous reviewers for their constructive feedback on an earlier version of this manuscript. Jie Yan is supported by the Strategic Priority Research Program of the Chinese Academy of Sciences, Grant No. XDB1350000.

References

- [1] Don Kurian Dennis, Tian Li, and Virginia Smith. Heterogeneity for the win: One-shot federated clustering. In *International conference on machine learning*, pages 2611–2620. PMLR, 2021.
- [2] Morris Stallmann and Anna Wilbik. Towards federated clustering: A federated fuzzy *c*-means algorithm (ffcm). *arXiv preprint arXiv:2201.07316*, 2022.
- [3] Chao Pan, Jin Sima, Saurav Prakash, Vishal Rana, and Olgica Milenkovic. Machine unlearning of federated clusters. In *The Eleventh International Conference on Learning Representations*, 2023.
- [4] Jinxuan Xu, Hong-You Chen, Wei-Lun Chao, and Yuqian Zhang. Jigsaw game: Federated clustering. *arXiv preprint arXiv:2407.12764*, 2024.
- [5] Kun Yang, Mohammad Mohammadi Amiri, and Sanjeev R Kulkarni. Greedy centroid initialization for federated k-means. *Knowledge and Information Systems*, 66(6):3393–3425, 2024.
- [6] Dong Qiao, Chris Ding, and Jicong Fan. Federated spectral clustering via secure similarity reconstruction. *Advances in Neural Information Processing Systems*, 36:58520–58555, 2023.
- [7] Shuai Wang and Tsung-Hui Chang. Federated matrix factorization: Algorithm design and application to data clustering. *IEEE Transactions on Signal Processing*, 70:1625–1640, 2022.
- [8] James MacQueen. Some methods for classification and analysis of multivariate observations. In *Proceedings of the Fifth Berkeley Symposium on Mathematical Statistics and Probability, Volume 1: Statistics*, volume 5, pages 281–298. University of California press, 1967.
- [9] James C Bezdek, Robert Ehrlich, and William Full. Fcm: The fuzzy c-means clustering algorithm. *Computers & geosciences*, 10(2-3):191–203, 1984.
- [10] Andrew Ng, Michael Jordan, and Yair Weiss. On spectral clustering: Analysis and an algorithm. *Advances in neural information processing systems*, 14, 2001.
- [11] Daniel Lee and H Sebastian Seung. Algorithms for non-negative matrix factorization. *Advances in neural information processing systems*, 13, 2000.
- [12] Qian Yu, Songze Li, Netanel Raviv, Seyed Mohammadreza Mousavi Kalan, Mahdi Soltanolkotabi, and Salman A Avestimehr. Lagrange coded computing: Optimal design for resiliency, security, and privacy. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 1215–1225. PMLR, 2019.
- [13] Abiodun M Ikotun, Absalom E Ezugwu, Laith Abualigah, Belal Abuhaija, and Jia Heming. K-means clustering algorithms: A comprehensive review, variants analysis, and advances in the era of big data. *Information Sciences*, 622:178–210, 2023.
- [14] Ling Ding, Chao Li, Di Jin, and Shifei Ding. Survey of spectral clustering based on graph theory. *Pattern Recognition*, page 110366, 2024.

- [15] Martin Ester, Hans-Peter Kriegel, Jörg Sander, Xiaowei Xu, et al. A density-based algorithm for discovering clusters in large spatial databases with noise. In kdd, volume 96, pages 226–231, 1996.
- [16] HaoChuan Xu and Ninh Pham. Scalable dbscan with random projections. *Advances in Neural Information Processing Systems*, 37:27978–28008, 2024.
- [17] K Chidananda Gowda and GJPR Krishna. Agglomerative clustering using the concept of mutual nearest neighbourhood. *Pattern recognition*, 10(2):105–112, 1978.
- [18] Eduardo Laber and Miguel Batista. On the cohesion and separability of average-link for hierarchical agglomerative clustering. *Advances in Neural Information Processing Systems*, 37:48710–48739, 2024.
- [19] Fangfang Li, Quanxue Gao, Qianqian Wang, Ming Yang, and Cheng Deng. Tensorized soft label learning based on orthogonal nmf. *IEEE Transactions on Neural Networks and Learning Systems*, 2024.
- [20] Jie Yan, Jing Liu, Yi-Zi Ning, and Zhong-Yuan Zhang. Sda-fc: Bridging federated clustering and deep generative model. *Information Sciences*, 681:121203, 2024.
- [21] Cynthia Dwork, Aaron Roth, et al. The algorithmic foundations of differential privacy. *Foundations and Trends® in Theoretical Computer Science*, 9(3–4):211–407, 2014.
- [22] Yinzhi Cao and Junfeng Yang. Towards making systems forget with machine unlearning. In 2015 IEEE symposium on security and privacy, pages 463–480. IEEE, 2015.
- [23] Songze Li, Sizai Hou, Baturalp Buyukates, and Salman Avestimehr. Secure federated clustering. *arXiv preprint arXiv:2205.15564*, 2022.
- [24] Yizhang Wang, Wei Pang, and Witold Pedrycz. One-shot federated clustering based on stable distance relationships. *IEEE Transactions on Industrial Informatics*, 2024.
- [25] Jonathan Scott, Christoph H Lampert, and David Saulpic. Differentially private federated *k*-means clustering with server-side data. *arXiv* preprint arXiv:2506.05408, 2025.
- [26] Abdulrahman Diaa, Thomas Humphries, and Florian Kerschbaum. {FastLloyd}: Federated, accurate, secure, and tunable {k-Means} clustering with differential privacy. In *34th USENIX Security Symposium (USENIX Security 25)*, pages 2733–2752, 2025.
- [27] Jie Zhang, Chen Chen, Bo Li, Lingjuan Lyu, Shuang Wu, Shouhong Ding, Chunhua Shen, and Chao Wu. Dense: Data-free one-shot federated learning. *Advances in Neural Information Processing Systems*, 35:21414–21428, 2022.
- [28] Jiawei Shao, Yuchang Sun, Songze Li, and Jun Zhang. Dres-fl: Dropout-resilient secure federated learning for non-iid clients via secret data sharing. Advances in Neural Information Processing Systems, 35:10533–10545, 2022.
- [29] LKPJ Rdusseeun and P Kaufman. Clustering by means of medoids. In *Proceedings of the statistical data analysis based on the L1 norm conference, neuchatel, switzerland*, volume 31, page 28, 1987.
- [30] Weina Wang, Lei Ying, and Junshan Zhang. On the relation between identifiability, differential privacy, and mutual-information privacy. *IEEE Transactions on Information Theory*, 62(9):5018– 5029, 2016.

- [31] Kiran S Kedlaya and Christopher Umans. Fast polynomial factorization and modular composition. *SIAM Journal on Computing*, 40(6):1767–1802, 2011.
- [32] Ronald A Fisher. The use of multiple measurements in taxonomic problems. *Annals of eugenics*, 7(2):179–188, 1936.
- [33] Li Deng. The mnist database of handwritten digit images for machine learning research [best of the web]. *IEEE signal processing magazine*, 29(6):141–142, 2012.
- [34] Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms. *arXiv* preprint arXiv:1708.07747, 2017.
- [35] Sameer A Nene, Shree K Nayar, Hiroshi Murase, et al. Columbia object image library (coil-20). 1996.
- [36] Fabian Keller, Emmanuel Muller, and Klemens Bohm. Hics: High contrast subspaces for density-based outlier ranking. In 2012 IEEE 28th international conference on data engineering, pages 1037–1048. IEEE, 2012.
- [37] Grace XY Zheng, Jessica M Terry, Phillip Belgrader, Paul Ryvkin, Zachary W Bent, Ryan Wilson, Solongo B Ziraldo, Tobias D Wheeler, Geoff P McDermott, Junjie Zhu, et al. Massively parallel digital transcriptional profiling of single cells. *Nature communications*, 8(1):14049, 2017.
- [38] Alexander Strehl and Joydeep Ghosh. Cluster ensembles—a knowledge reuse framework for combining multiple partitions. *Journal of machine learning research*, 3(Dec):583–617, 2002.
- [39] Xin Liu, Hui-Min Cheng, and Zhong-Yuan Zhang. Evaluation of community detection methods. *IEEE Transactions on Knowledge and Data Engineering*, 32(9):1736–1746, 2019.
- [40] Jie Yan, Xin Liu, Ji Qi, Tao You, and Zhong-Yuan Zhang. The significance of kappa and f-score in clustering ensemble: a comprehensive analysis. *Knowledge and Information Systems*, pages 1–36, 2025.
- [41] Jichan Chung, Kangwook Lee, and Kannan Ramchandran. Federated unsupervised clustering with generative models. In *AAAI 2022 international workshop on trustable, verifiable and auditable federated learning*, volume 4, 2022.
- [42] Jeffrey Humpherys and Tyler J Jarvis. Foundations of Applied Mathematics Volume 2: Algorithms, Approximation, Optimization. SIAM, 2020.
- [43] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.
- [44] Pauli Virtanen, Ralf Gommers, Travis E. Oliphant, Matt Haberland, Tyler Reddy, David Cournapeau, Evgeni Burovski, Pearu Peterson, Warren Weckesser, Jonathan Bright, Stéfan J. van der Walt, Matthew Brett, Joshua Wilson, K. Jarrod Millman, Nikolay Mayorov, Andrew R. J. Nelson, Eric Jones, Robert Kern, Eric Larson, C J Carey, İlhan Polat, Yu Feng, Eric W. Moore, Jake VanderPlas, Denis Laxalde, Josef Perktold, Robert Cimrman, Ian Henriksen, E. A. Quintero, Charles R. Harris, Anne M. Archibald, Antônio H. Ribeiro, Fabian Pedregosa, Paul van Mulbregt, and SciPy 1.0 Contributors. SciPy 1.0: Fundamental Algorithms for Scientific Computing in Python. Nature Methods, 17:261–272, 2020.

- [45] Madson Luiz Dantas Dias. fuzzy-c-means: An implementation of fuzzy *c*-means clustering algorithm., May 2019.
- [46] Tian Li, Anit Kumar Sahu, Ameet Talwalkar, and Virginia Smith. Federated learning: Challenges, methods, and future directions. *IEEE signal processing magazine*, 37(3):50–60, 2020.
- [47] Bo Yang, Xiao Fu, Nicholas D Sidiropoulos, and Mingyi Hong. Towards k-means-friendly spaces: Simultaneous deep learning and clustering. In *international conference on machine learning*, pages 3861–3870. PMLR, 2017.
- [48] Lukas Miklautz, Timo Klein, Kevin Sidak, Collin Leiber, Thomas Lang, Andrii Shkabrii, Sebastian Tschiatschek, and Claudia Plant. Breaking the reclustering barrier in centroid-based deep clustering. arXiv preprint arXiv:2411.02275, 2024.
- [49] Sheng Zhou, Hongjia Xu, Zhuonan Zheng, Jiawei Chen, Zhao Li, Jiajun Bu, Jia Wu, Xin Wang, Wenwu Zhu, and Martin Ester. A comprehensive survey on deep clustering: Taxonomy, challenges, and future directions. *ACM Computing Surveys*, 57(3):1–38, 2024.
- [50] Jie Yan, Jing Liu, Ji Qi, and Zhong-Yuan Zhang. Privacy-preserving federated deep clustering based on gan. *arXiv preprint arXiv:2211.16965*, 2022.
- [51] Jie Yan, Jing Liu, and Zhong-Yuan Zhang. Ccfc: Bridging federated clustering and contrastive learning. *arXiv preprint arXiv:2401.06634*, 2024.
- [52] Jie Yan, Jing Liu, Yi-Zi Ning, and Zhong-Yuan Zhang. Ccfc++: Enhancing federated clustering through feature decorrelation. *arXiv preprint arXiv:2402.12852*, 2024.

NeurIPS Paper Checklist

1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [Yes]

Justification: Please refer to the abstract and introduction.

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: Please refer to Section E.

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best

judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [Yes]

Justification: Please refer to Section B.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: Please refer to Sections 4.1 and C.2. .

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived
 well by the reviewers: Making the paper reproducible is important, regardless of
 whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example

- (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
- (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
- (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
- (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [No]

Justification: We have provided publicly available dataset information in Section C.1. The code for the proposed OmniFC will be released upon the paper's formal publication.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: Please refer to Sections 4.1 and C.2.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [No]

Justification: The deterministic nature of the global distance matrix reconstruction in the proposed OmniFC ensures consistent quality evaluations across multiple runs.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: Please refer to Section C.

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.

- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines?

Answer: [Yes]

Justification: The research conducted in the paper conform with the NeurIPS Code of Ethics.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. Broader impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [Yes]

Justification: This paper mainly targets privacy-preserving clustering in federated scenarios. By addressing this challenge, we can further promote the practical deployment of clustering in sensitive domains, such as healthcare and finance, while safeguarding data security and

user privacy.

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks,

mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: This paper poses no such risks.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with
 necessary safeguards to allow for controlled use of the model, for example by requiring
 that users adhere to usage guidelines or restrictions to access the model or implementing
 safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do
 not require this, but we encourage authors to take this into account and make a best
 faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: The original sources of the datasets are cited in Section C.1.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. New assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [NA]

Justification: This paper does not release new assets.

Guidelines:

- The answer NA means that the paper does not release new assets.
- · Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. Crowdsourcing and research with human subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: This paper does not involve crowdsourcing nor research with human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. Institutional review board (IRB) approvals or equivalent for research with human subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: This paper does not involve crowdsourcing nor research with human subjects.

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.

• For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

16. Declaration of LLM usage

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does not impact the core methodology, scientific rigorousness, or originality of the research, declaration is not required.

Answer: [NA]

Justification: The core method development in this research does not involve LLMs as any important, original, or non-standard components.

- The answer NA means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy (https://neurips.cc/Conferences/2025/LLM) for what should or should not be described.

A Pseudocode of the Proposed OmniFC

The procedure of OmniFC is formally presented in Algorithm 1. On the client side, each sample \tilde{x}_i is independently encoded into $z_{i,j}$ based on Equation (3), and then transmitted to the j-th client, where $i \in [n]$ and $j \in [m]$. Then, each client j computes pairwise distances between all encoded representations $z_{i,j}$ and $z_{i',j}$ $(i,i' \in [n])$ using Equation (4), and transmits the results to the central server. On the server side, the global distance matrix is reconstructed based on Equations (6) and (7), and subsequently utilized by a centralized clustering algorithm to derive the final clustering outcome π^* .

Algorithm 1: OmniFC

Input: Local datasets $\{X_j\}_{j=1}^m$, prime number p, the number of segments l, the number of noises t, pre-specified hyperparameters $\{\alpha_o\}_{o=1}^{l+t}$ and $\{\beta_j\}_{j=1}^m$.

Output: The final partition π^* .

- 1 Clients execute:
- **2 Local Lagrange Encoding and Secret Sharing:**
- Each sample \tilde{x}_i is encoded via Equation (3), i.e., $z_{i,j} = f_{\tilde{x}_i}(\beta_j)$, and subsequently
- transmitted to the j-th client, where $i \in [n]$ and $j \in [m]$.
- **5 Global Distance Reconstruction:**
- Each client j computes pairwise distances between all encoded representations $z_{i,j}$ and
- $z_{i',j}$ $(i,i' \in [n])$ using Equation (4), and transmits the results to the central server.
- 8 Server executes:
- 9 Global Distance Reconstruction:
- The server reconstructs the global distance matrix according to Equations (6) and (7).
- 11 Cluster assignment:
- The global distance matrix is fed into a centralized clustering method to obtain π^* .

B Proofs of Theorems

Before proving the theorems, we first summarize some notations used throughout the main text and this appendix, and introduce two lemmas from Ref. [42] and [12]. Refer to Table 3 for the notions, with the lemmas delineated below.

Lemma 1. [42] Given n distinct points $\{(\boldsymbol{x}_i, \boldsymbol{y}_i)\}_{i=1}^n$ with mutually different x_i , there exists a unique polynomial $f(\boldsymbol{x})$ of degree no greater than n-1 that interpolates the data, i.e., $f(\boldsymbol{x}_i) = \boldsymbol{y}_i$.

Lemma 2. [12] Given the number of noises t, and a polynomial f used to compute $f(\tilde{X})$, and the degree of f is denoted as deg(f). When $m \ge deg(f)(l+t-1)+1$, a t-private LCC encryption is achievable, e.g.,

$$I(\tilde{\boldsymbol{x}}_i; \{\boldsymbol{z}_{i,j}\}_{j \in \mathcal{C}}) = 0, \tag{9}$$

where $I(\cdot;\cdot)$ denotes the mutual information function, $\mathcal{C} \subset [m]$ and $|\mathcal{C}| \leq t$.

Proof of Theorem 1. With Lemma 1, we prove Theorem 1 as follows.

Proof. The server possesses only the pre-defined public hyperparameters $\{\alpha_o\}_{o=1}^{l+t}$, $\{\beta_j\}_{j=1}^m$ and the distance $\{dis(\boldsymbol{z}_{i,j},\boldsymbol{z}_{i',j})\}_{j=1}^m$. For each distance $dis(\boldsymbol{z}_{i,j},\boldsymbol{z}_{i',j})$ $(j \in [m])$, it can be further formulated as:

$$dis(\mathbf{z}_{i,j}, \mathbf{z}_{i',j}) = \|\mathbf{z}_{i,j} - \mathbf{z}_{i',j}\|_{2}^{2} = \|\mathbf{f}_{\tilde{\mathbf{x}}_{i}}(\beta_{j}) - \mathbf{f}_{\tilde{\mathbf{x}}_{i'}}(\beta_{j})\|_{2}^{2},$$
(10)

Table 3: Notations.

Notation	Explanation
m	Number of clients.
$oldsymbol{X} \in \mathbb{R}^{n imes d}$	The centralized dataset $X \in \mathbb{R}^{n \times d}$ consists of n d -dimensional samples $\{x_i\}_{i=1}^n$, which are distributed among m clients, i.e., $X = \bigcup_{j=1}^m X_j$.
$ ilde{m{X}} \in \mathbb{F}_p^{n imes d}$	\tilde{X} denotes the representation of X over the finite field $\mathbb{F}_p^{n \times d}$, consisting of n d -dimensional samples $\{\tilde{x}_i\}_{i=1}^n$, where each \tilde{x}_i corresponds to the transformed version of x_i in \mathbb{F}_p^d .
l	Number of data segments.
t	Number of noises.
$oldsymbol{s}_{i,o} \in \mathbb{F}_p^{rac{d}{l}}$	$\tilde{\boldsymbol{x}}_i = [\boldsymbol{s}_{i,1}^T, \boldsymbol{s}_{i,2}^T, \cdots, \boldsymbol{s}_{i,l}^T]^T$, where $\boldsymbol{s}_{i,o}$ denotes the o -th segment of $\tilde{\boldsymbol{x}}_i$ for $o \in [l] = \{1, 2, \cdots, l\}$. For $l < o \le l+t$, $\boldsymbol{s}_{i,o}$ corresponds to the o -th noise uniformly sampled from \mathbb{F}_{n}^{d} .
$\{\alpha_o\}_{o=1}^{l+t}$	A collection of $l+t$ distinct hyperparameters from \mathbb{F}_p , predetermined by consensus between all clients and the central server, serves to construct the Lagrange interpolation polynomial.
$\{\beta_j\}_{j=1}^m$	A collection of m distinct hyperparameters from \mathbb{F}_p , predetermined by consensus between all clients and the central server, serves to encode the local data into m distinct representations.
$oldsymbol{Z}_j \in \mathbb{F}_p^{n imes rac{d}{l}}$	The global encoded dataset possessed by client j ($j \in [m]$). $\mathbf{Z}_j = [\mathbf{z}_{1,j}, \mathbf{z}_{2,j}, \cdots, \mathbf{z}_{n,j}]^T$, where $\mathbf{z}_{i,j}$ ($i \in [n]$) is the encoded representation of $\tilde{\mathbf{x}}_i$ at client j .

implying that it corresponds to the evaluation of a degree-2(l+t-1) polynomial at β_j . According to Lemma 1, the polynomial can be uniquely interpolated from 2(l+t-1)+1 distinct points. That is, when $m \geq 2l+2t-1$, the polynomial can be interpolated from the set $\{(\beta_j, dis(\boldsymbol{z}_{i,j}, \boldsymbol{z}_{i',j}))\}_{j=1}^m$, and it is exactly $\boldsymbol{f}_{\tilde{\boldsymbol{x}}_i, \tilde{\boldsymbol{x}}_{i'}}(\beta)$, i.e.,

$$\boldsymbol{f}_{\tilde{\boldsymbol{x}}_{i},\tilde{\boldsymbol{x}}_{i'}}(\beta) = \left\| \boldsymbol{f}_{\tilde{\boldsymbol{x}}_{i}}(\beta) - \boldsymbol{f}_{\tilde{\boldsymbol{x}}_{i'}}(\beta) \right\|_{2}^{2}.$$
 (11)

Particularly, by assigning $\beta = \alpha_o$ ($o \in [l]$), the distance between the o-th data segments of \tilde{x}_i and $\tilde{x}_{i'}$ can be accurately recovered:

$$f_{\tilde{x}_{i},\tilde{x}_{i'}}(\alpha_{o}) = \left\| f_{\tilde{x}_{i}}(\alpha_{o}) - f_{\tilde{x}_{i'}}(\alpha_{o}) \right\|_{2}^{2} = \left\| s_{i,o} - s_{i',o} \right\|_{2}^{2}.$$
 (12)

Consequently, the distance between \tilde{x}_i and $\tilde{x}_{i'}$ can be precisely reconstructed:

$$\sum_{o=1}^{l} f_{\tilde{x}_{i}, \tilde{x}_{i'}}(\alpha_{o}) = \sum_{o=1}^{l} \|s_{i,o} - s_{i',o}\|_{2}^{2} = dis(\tilde{x}_{i}, \tilde{x}_{i'})$$
(13)

Note that since the above proof does not impose any constraints on the distribution of \tilde{x}_i and $\tilde{x}_{i'}$ across clients, Equation (13) holds irrespective of how data is distributed among clients.

Proof of Theorem 2. With Lemma 2, we prove Theorem 2 as follows.

Table 4: Reconstruction time (seconds) of the global pairwise distance matrix across diverse data scales under varying l. Computational efficiency can be substantially improved by increasing the number of data segments l.

\overline{n}	1k	2k	5k	10k	20k	40k	70k
l=2	4	15	78	285	1053	4634	13576
l=4	3	9	55	154	620	2037	7094
l = 8	2	6	32	124	435	1471	4352

Proof. We prove Theorem 2 by instantiating Lemma 2 with the specific polynomial structure used in the OmniFC framework.

Recall that Lemma 2 states that a t-private LCC encryption is achievable when

$$m > \deg(\mathbf{f})(l+t-1)+1,$$

where f is the polynomial used in the encoding scheme, and l is the number of data segments.

In the OmniFC setting, the polynomial f is a quadratic distance-based function of degree 2, i.e., deg(f) = 2. Plugging this into the general LCC bound yields:

$$m > 2l + 2t - 1$$
.

Therefore, under this condition, the mutual information between any private input \tilde{x}_i and the encoded messages observed by up to t colluding clients satisfies:

$$I(\tilde{\boldsymbol{x}}_i; \{\boldsymbol{z}_{i,i}\}_{i \in \mathcal{C}}) = 0,$$

where $\mathcal{C} \subset [m]$ and $|\mathcal{C}| \leq t$.

This guarantees t-privacy in the OmniFC framework, thus completing the proof.

C Experimental Details

All experiments are implemented in Python and executed on a system equipped with an Intel Core i7-12650H CPU, 16GB of RAM, and an NVIDIA GeForce RTX 4060 GPU.

C.1 Datasets and Evaluation Criteria

Datasets. As shown in Table 5, we select seven benchmark datasets across tabular, visual, temporal, and genomic domains, including Iris [32], MNIST [33], Fashion-MNIST [34], COIL-20 [35], COIL-100 [35], Pendigits [36], and 10x_73k [37]. The chosen datasets encompass diverse modalities, dimensionalities, and cluster patterns, facilitating a comprehensive evaluation of the method's generalizability in practical scenarios.

Fig. 6 exemplifies, through the Iris dataset, our simulation of federated scenarios under different Non-IID conditions. We simulate diverse federated settings by evenly partitioning the Iris dataset into 3 (the number of true clusters) subsets—each representing a client—and adjusting the non-IID level p. For a client with 50 datapoints, the first $p \cdot 50$ datapoints are sampled from a single cluster, and the remaining $(1-p) \cdot 50$ ones are randomly sampled from any cluster. As such, p=0 recovers the IID setting, whereas p=1 induces a maximally skewed distribution, where each client's data is fully concentrated within a single cluster.

Table 5: Description of datasets.

Dataset	Type	Size	Image size/Features	Class
Iris	tabular	150	4	3
MNIST	image	70000	28×28	10
Fashion-MNIST	image	70000	28×28	10
COIL-20	image	1440	128×128	20
COIL-100	image	7200	128×128	100
Pendigits	time series	10992	16	10
10x_73k	gene	73233	720	8

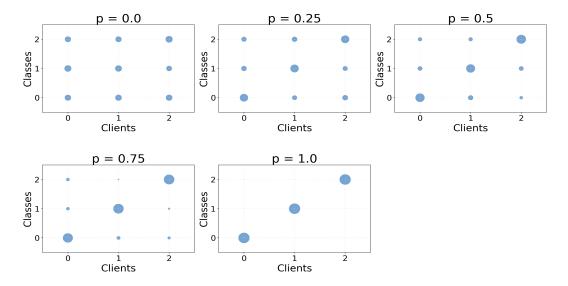


Figure 6: Data partition visualization on Iris.

Evaluation Criteria. Evaluation is based on two metrics—normalized mutual information (NMI) [38] and Kappa [39]—where elevated scores denote superior clustering quality. Despite being widely adopted, NMI has been shown to have limitations, such as the finite size effect, and fails to account for the importance of small clusters [39, 20, 40]. In contrast, Kappa addresses these concerns, making it a more reliable alternative for clustering evaluation. Hence, our analysis is grounded in Kappa-based results, with NMI-based outcomes serving only as supplementary references.

C.2 Implementation Details

All centralized clustering methods are implemented by leveraging existing open-source Python libraries: KM, KMed, SC, NMF, and DBSCAN utilize the sklearn library [43], HC employs the scipy library [44], and FCM adopts an individual open-source implementation [45]. For OmniFC, $\{\alpha_o\}_{o=1}^{l+t}$ is set as a sequence of l+t consecutive odd integers starting from 1, while $\{\beta_j\}_{j=1}^m$ is set as a sequence of m consecutive even integers starting from 0. The default values of l and t are set to 2.

We acknowledge the importance of the constraint $m \ge 2(l+t-1)+1$ for successful distance reconstruction. Among these hyperparameters, m (number of clients) is typically predetermined

by the federated scenario, while l (segments) and t (noises) are tunable hyperparameters that can be adjusted to satisfy this constraint. In practice, violating this condition is extremely rare because the number of clients m in typical federated learning scenarios ranges from tens to millions [46], providing substantial flexibility for hyperparameter selection. In practice, m is fixed by the system, while the choice of l and t depends on the relative emphasis placed on communication and computational efficiency versus privacy preservation. Note that as long as $m \geq 2(l+t-1)+1$ holds, the specific values of l and t have no effect on the reconstruction effectiveness of the global pairwise distance matrix (Theorem 1 and Fig. 5). The hyperparameter l governs the dimensionality $(\frac{d}{l})$ of both the shared encoded data across clients and the features involved in pairwise distance computation. A larger l reduces the volume of shared data and the number of features compared, thus improving communication and computational efficiency (Table 4). The hyperparameter t governs the information-theoretic security, with higher values enhancing resilience against client collusion (Theorem 2). However, the constraint $m \geq 2(l+t-1)+1$ precludes simultaneous increases in l and t, indicating that their selection hinges on the trade-off between communication/computational efficiency and privacy preservation.

Additionally, several clustering methods evaluated in our experiments demand full $n \times n$ pairwise distance matrix computations, imposing substantial computational and memory burdens on large-scale datasets. To facilitate the execution of comprehensive experiments, we implement a subsampling strategy whereby 1000 samples are randomly drawn from datasets exceeding 5000 entries to form the experimental subset. This approach balances computational efficiency with the preservation of the original data distribution, enabling fair and meaningful comparisons across methods. The sensitivity of the proposed OmniFC with respect to the number of samples is presented in Appendix D.2.

D Supplementary Experimental Results

D.1 NMI-based Evaluation Results

To supplement the Kappa-based evaluation results and to enable broader comparability with existing FC works, we additionally provide NMI-based evaluation results in Tables 6 and 7. Similar to the Kappa-based evaluation results, the numerical results based on NMI also confirm the effectiveness and generalizability of OmniFC.

D.2 Sensitivity Analysis

To assess the sensitivity of the proposed OmniFC concerning the number of samples, we evaluate the global distance matrix reconstruction loss—defined as the root-mean-square deviation (RMSE) between the ground-truth and the reconstructed pairwise distance matrices—across different sample sizes. As shown in Table 8, OmniFC exhibits favorable scalability concerning sample size.

D.3 Non-IID Partitioning based on Dirichlet Distributions

To test more diverse Non-IID scenarios, we used Dirichlet distributions with varying concentrations α to model label heterogeneity on Pendigits in Table 9, and other experimental settings are consistent with those in Table 1. Experimental results indicate that the OmniFC method exhibits superior efficacy and robustness.

D.4 The Feasibility of Extending Deep Clustering

To further showcase the flexibility of the proposed framework, we validated the feasibility of employing OmniFC to extend deep clustering. As a foundational deep k-means (KM) method, DCN [47]

Table 6: **NMI of clustering methods in different federated scenarios.** For each comparison, the best result is highlighted in boldface.

Dataset	m	SC-base	d metho	ds	KM-	based m	ethods		FCM-base	d metho	ds	NMF-bas	sed method	s
Dataset	\boldsymbol{p}	SC_central	FedSC	Ours	KM_central	k-FED	MUFC	Ours	FCM_central	FFCM	Ours	NMF_central	FedMAvg	Ours
	0.00		0.90	0.90		0.66	0.76	0.90		0.91	0.90		0.73	0.90
	0.25		0.85	0.90		0.90	0.85	0.90		0.72	0.90		0.73	0.90
Iris	0.50	0.90	0.75	0.90	0.90	0.87	0.70	0.90	0.90	0.87	0.90	0.56	0.73	0.90
	0.75		0.85	0.90		0.90	0.74	0.90		0.91	0.90		0.73	0.90
	1.00		0.29	0.90		0.70	0.70	0.90		0.93	0.90		0.73	0.90
	0.00		0.59	0.58		0.51	0.48	0.46		0.53	0.43		0.48	0.47
	0.25		0.60	0.58		0.49	0.52	0.46		0.53	0.43		0.45	0.47
MNIST	0.50	0.58	0.59	0.58	0.54	0.39	0.50	0.46	0.55	0.52	0.43	0.47	0.43	0.47
	0.75		0.59	0.58		0.46	0.52	0.46		0.52	0.43		0.47	0.47
	1.00		0.45	0.58		0.51	0.55	0.46		0.57	0.43		0.47	0.47
	0.00		0.61	0.61		0.56	0.56	0.52		0.61	0.53		0.53	0.55
	0.25		0.60	0.61		0.54	0.54	0.52		0.59	0.53		0.53	0.55
Fashion-MNIST	0.50	0.61	0.61	0.61	0.62	0.57	0.60	0.52	0.61	0.58	0.53	0.60	0.53	0.55
	0.75		0.55	0.61		0.55	0.54	0.52		0.61	0.53		0.53	0.55
	1.00		0.39	0.61		0.48	0.59	0.52		0.58	0.53		0.53	0.55
	0.00		0.80	0.75		0.65	0.74	0.74		0.71	0.72		0.62	0.75
	0.25		0.78	0.75		0.70	0.73	0.74		0.69	0.72		0.62	0.75
COIL-20	0.50	0.75	0.80	0.75	0.74	0.66	0.72	0.74	0.75	0.72	0.72	0.70	0.62	0.75
	0.75		0.69	0.75		0.67	0.73	0.74		0.74	0.72		0.63	0.75
	1.00		0.46	0.75		0.69	0.72	0.74		0.75	0.72		0.63	0.75
	0.00		0.67	0.79		0.76	0.76	0.79		0.69	0.79		0.70	0.76
	0.25		0.66	0.79		0.75	0.76	0.79		0.71	0.79		0.70	0.76
COIL-100	0.50	0.79	0.66	0.79	0.77	0.75	0.76	0.79	0.79	0.77	0.79	0.72	0.70	0.76
	0.75		0.64	0.79		0.75	0.76	0.79		0.77	0.79		0.70	0.76
	1.00		0.61	0.79		0.75	0.79	0.79		0.81	0.79		0.70	0.76
	0.00		0.77	0.72		0.67	0.67	0.67		0.68	0.70		0.42	0.71
	0.25		0.76	0.72		0.62	0.66	0.67		0.68	0.70		0.42	0.71
Pendigits	0.50	0.72	0.74	0.72	0.69	0.63	0.67	0.67	0.69	0.67	0.70	0.55	0.42	0.71
	0.75		0.75	0.72		0.50	0.64	0.67		0.65	0.70		0.42	0.71
	1.00		0.62	0.72		0.64	0.71	0.67		0.69	0.70		0.42	0.71
	0.00		0.71	0.85		0.68	0.65	0.58		0.69	0.58		0.66	0.78
	0.25		0.71	0.85		0.70	0.68	0.58		0.70	0.58		0.66	0.78
10X_73k	0.50	0.85	0.70	0.85	0.82	0.73	0.72	0.58	0.68	0.79	0.58	0.83	0.66	0.78
	0.75		0.59	0.85		0.65	0.73	0.58		0.83	0.58		0.66	0.78
	1.00		0.19	0.85		0.49	0.80	0.58		0.82	0.58		0.66	0.78
count	-	-	13	22	-	5	12	18	-	22	13	-	1	34

Table 7: **NMI of different clustering methods.**

Dataset	KMed-bas	ed methods	DBSCAN-	based methods	HC-based methods		
Butuset	Central	Ours	Central	Ours	Central	Ours	
Iris	0.86	0.86	0.73	0.73	0.89	0.90	
MNIST	0.38	0.38	0.56	0.56	0.49	0.49	
Fashion-MNIST	0.49	0.49	0.53	0.53	0.54	0.54	
COIL-20	0.60	0.61	0.86	0.86	0.70	0.69	
COIL-100	0.69	0.69	0.85	0.85	0.78	0.78	
Pendigits	0.56	0.55	0.74	0.74	0.69	0.69	
10x_73k	0.31	0.31	0.45	0.45	0.51	0.50	

Table 8: Sample-size sensitivity of the global distance matrix reconstruction loss on MNIST.

\overline{n}	1000	2000	3000	5000	7000	10000
RMSE	0.0002	0.0002	0.0002	0.0002	0.0002	0.0002

Table 9: **Kappa of clustering methods in different federated scenarios.** For each comparison, the best result is highlighted in boldface.

α	SC-based methods		KM-based methods				FCM-based methods			NMF-based methods				
	SC_c	central	FedSC	OmniFC-SC	KM_central	k-FED	MUFC	OmniFC-KM	FCM_central	FFCM	OmniFC-FCM	NMF_central	FedMAvg	OmniFC-NMF
1000			0.71	0.72		0.58	0.61	0.62		0.60	0.66		0.33	0.72
5	0	.72	0.73	0.72	0.62	0.48	0.61	0.62	0.65	0.52	0.66	0.70	0.33	0.72
0.001			0.59	0.72		0.51	0.60	0.62		0.54	0.66		0.33	0.72
count		-	1	2	-	0	0	3	-	0	3	-	0	3

Table 10: **Kappa of clustering methods in different simulated federated datasets.** OmniFC-DCN achieves superior clustering performance relative to its shallow counterpart, OmniFC-KM.

Method	MNIST	Pendigits	10x_73k 0.56 0.60	
OmniFC-KM	0.42	0.62		
OmniFC-DCN	0.43	0.64		

has spurred a range of advanced clustering techniques and found widespread use in privacy-sensitive domains such as medicine and finance [48]. Therefore, extending DCN could significantly impact multiple research fields and accelerate progress on downstream applications.

To this end, we extend DCN analogously to KM by using the reconstructed distance matrix in place of raw features as the model input. This extended variant is denoted as OmniFC-DCN. As evidenced in Table 10, OmniFC-DCN outperforms its shallow counterpart, OmniFC-KM, in clustering performance, further highlighting the flexibility of OmniFC.

E Limitation

This work primarily focuses on extending shallow centralized clustering methods and may be less effective for high-dimensional or intrinsically complex data. A promising future direction is to explore how the reconstructed global distance matrix can substantially support the federated extension of deep centralized clustering methods, thereby enabling more powerful representation learning under complex data distributions [49, 50, 51, 52].