

---

# Federated Spectral Clustering via Secure Similarity Reconstruction

---

Dong Qiao<sup>1,2</sup>   Chris Ding<sup>1</sup>   Jicong Fan<sup>1,2\*</sup>

<sup>1</sup>The Chinese University of Hong Kong, Shenzhen, China

<sup>2</sup>Shenzhen Research Institute of Big Data, Shenzhen, China

dongqiao@link.cuhk.edu.cn   {chrisding,fanjicong}@cuhk.edu.cn

## Abstract

Federated learning has a significant advantage in protecting data and information privacy. Many scholars proposed various secure learning methods within the framework of federated learning but the study on secure federated unsupervised learning especially clustering is limited. We in this work propose a secure kernelized factorization method for federated spectral clustering on distributed data. The method is non-trivial because the kernel or similarity matrix for spectral clustering is computed by data pairs, which violates the principle of privacy protection. Our method implicitly constructs an approximation for the kernel matrix on distributed data such that we can perform spectral clustering under the constraint of privacy protection. We provide a convergence guarantee of the optimization algorithm, reconstruction error bounds of the Gaussian kernel matrix, and the sufficient condition of correct clustering of our method. We also present guarantees of differential privacy. Numerical results on synthetic and real datasets demonstrate that the proposed method is efficient and accurate in comparison to baselines.

## 1 Introduction

In the era of big data, human beings can analyze massive data in various fields due to the improvement of storage and computational capabilities of computing devices [Li *et al.*, 2020b]. Some popular fields such as artificial intelligence, machine learning, internet of things (IoT), and cloud computing have seen explosive development over the past few years. Nevertheless, a side effect of this trend is that individuals and organizations have more and more concerns about potential violation of privacy [Kairouz *et al.*, 2021]. As a result, it has become a challenge to mine valuable information from user data but not to directly access it.

Federated learning [Kairouz *et al.*, 2021; McMahan *et al.*, 2017] can train a global model without retrieving dispersed data [Yang *et al.*, 2018]. This advantage has made it so popular that many scholars have put much effort into the study of federated learning. For example, Yang *et al.* [2019] presented the definitions of horizontal federated learning, vertical federated learning, and federated transfer learning. Some privacy-preserving machine learning models were also presented. For instance, He *et al.* [2020] developed a federated group knowledge transfer algorithm to train small CNNs on edge devices. Chen *et al.* [2018] proposed a protocol to conduct privacy-preserving ridge regression over high-dimensional data. Besides, Kim *et al.* [2018] proposed a block-chained federated learning architecture that enables on-device learning without any central coordination.

Regardless of the great progress of federated learning, it can be found that most of the existing studies are for supervised learning [Li *et al.*, 2020a; Ghosh *et al.*, 2020]. Note that collecting labeled data may deserve very high cost in real situations [Li *et al.*, 2020b] while unlabeled data are abundant. Thus, it

---

\*Corresponding author

is necessary and important to study federated learning for unsupervised learning [Zhang *et al.*, 2020; Tzinis *et al.*, 2021; Zhuang *et al.*, 2021; Dennis *et al.*, 2021] such as clustering [Ng *et al.*, 2001; Fan and Chow, 2017; Fan *et al.*, 2018, 2021; Fan, 2021; Cai *et al.*, 2022; Fan *et al.*, 2022]. For example, Li *et al.* [2021] proposed a federated matrix factorization with a privacy guarantee for recommendation systems. Wang and Chang [2022] proposed two federated matrix factorization algorithms that can be used for federated clustering. Besides, there are some studies on federated spectral clustering. For instance, Wang *et al.* [2020] presented a federated multi-view spectral clustering method under the assumption that the data of each view are in one client. Hernández-Pereira *et al.* [2021] developed a cooperative spectral clustering model to deal with distributed data but the model is linear. However, the study on federated spectral clustering is still very limited and deserves more attention and effort.

In this paper, we propose a federated kernelized factorization method to reconstruct a similarity matrix for secure spectral clustering on distributed data. Our contributions are as follows.

- We propose a federated spectral clustering algorithm and provide convergence guarantee for the optimization.
- We further propose to add noise to the data or the learned factors to enhance the security of clustering and provide guarantees of differential privacy.
- We provide upper bounds for the reconstruction error of the true similarity matrix and theoretical guarantees for correct clustering.

We test our method on both synthetic data and real datasets in comparison to baselines, which verify the effectiveness of our method.

**Notations** We use  $y$ ,  $\mathbf{y}$ , and  $\mathbf{Y}$  to denote scalar, vector, and matrix, respectively. The element of  $\mathbf{Y}$  at row  $i$  and column  $j$  is denoted by  $y_{ij}$ . We use  $\|\cdot\|_2$  to denote the  $\ell_2$  norm of a vector and use  $\text{Tr}(\cdot)$ ,  $\|\cdot\|_F$ , and  $\|\cdot\|_{sp}$  to denote the trace, Frobenius norm, and spectral norm of a matrix respectively. The  $\ell_\infty$  norm and  $\ell_{2,\infty}$  norm of a matrix  $\mathbf{Y}$  are defined as  $\|\mathbf{Y}\|_\infty = \max_{ij} |y_{ij}|$  and  $\|\mathbf{Y}\|_{2,\infty} = \max_j \sqrt{\sum_i y_{ij}^2}$  respectively.  $\mathbf{K}$ ,  $\mathcal{K}$ ,  $\mathcal{K}$ , and  $k$  denote the kernel matrix, kernel function, number of clusters, and the number  $k$  in KNN, respectively.  $\phi$  denotes the feature map induced by  $\mathcal{K}$ .

## 2 Federated Spectral Clustering (FedSC)

Suppose we have  $n$  data points of dimension  $m$  distributing in  $P$  clients. For convenience, we denote by  $\mathbf{X} \in \mathbb{R}^{m \times n}$  the matrix composed of all the  $n$  data points and denote by  $\mathbf{X}_p \in \mathbb{R}^{m \times N_p}$  the matrix composed of the  $N_p$  data points in client  $c_p$ , where  $N_p \geq 1$ ,  $p = 1, \dots, P$ , and  $\sum_{p=1}^P N_p = n$ . Without loss of generality, we let  $\mathbf{X} = [\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_P]$ , which means  $\{\mathbf{X}_p\}_{p=1}^P$  are submatrices of  $\mathbf{X}$ . Our goal is to perform spectral clustering on these data to partition them into  $\mathcal{K}$  groups, under the constraint that the data in each client cannot leave the client itself and the privacy of the data should be protected as much as possible, though there could be a central server conducting clustering.

The aforementioned task is non-trivial because in spectral clustering, the first step is constructing an adjacency matrix  $\mathbf{A} \in \mathbb{R}^{n \times n}$ , which has to evaluate the similarity between every data pair  $(\mathbf{x}_i, \mathbf{x}_j)$  using a metric function  $\mathcal{M}(\cdot, \cdot)$  and hence violates the privacy constraint in the task. To solve the problem, we present a federated spectral clustering model in this section.

### 2.1 Similarity Reconstruction via Feature Space Factorization

In spectral clustering, for  $\mathcal{M}(\cdot, \cdot)$ , there are many choices such as  $k$  nearest neighbor similarity and various kernel functions. Let  $\mathcal{K}(\cdot, \cdot)$  be a kernel function and we have

$$\mathcal{K}(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j), \quad (1)$$

where  $\phi : \mathbb{R}^m \rightarrow \mathbb{R}^{m'}$  is a feature map induced by the kernel function<sup>2</sup> and does not need to be carried out explicitly. When it comes to federated spectral clustering, the central server has no access

<sup>2</sup>The most widely-used kernel is the Gaussian kernel  $\mathcal{K}(\mathbf{x}_i, \mathbf{x}_j) = \exp(-\frac{1}{2r^2} \|\mathbf{x}_i - \mathbf{x}_j\|_2^2)$ , of which the feature map  $\phi$  is an infinite-order polynomial feature map and  $r$  is a hyperparameter controlling the smoothness.

to the raw data distributed in clients and hence cannot compute  $\mathcal{K}(\mathbf{x}_i, \mathbf{x}_j)$  using (1). However, if the central server can learn an effective approximation (denoted by  $\widehat{\phi(\mathbf{x}_i)}$ ) for each  $\phi(\mathbf{x}_i)$  without accessing  $\mathbf{x}_i$ ,  $\mathcal{K}(\mathbf{x}_i, \mathbf{x}_j)$  can be estimated, i.e.,

$$\mathcal{K}(\mathbf{x}_i, \mathbf{x}_j) \simeq \widehat{\phi(\mathbf{x}_i)}^T \widehat{\phi(\mathbf{x}_j)}. \quad (2)$$

Thus, inspired by [Fan and Udell, 2019; Fan *et al.*, 2021], we propose to approximate each  $\phi(\mathbf{x}_i)$  by

$$\widehat{\phi(\mathbf{x}_i)} = \phi(\mathbf{Z})\mathbf{c}_i, \quad (3)$$

where  $\mathbf{Z} = [\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_d] \in \mathbb{R}^{m \times d}$ ,  $\phi(\mathbf{Z}) = [\phi(\mathbf{z}_1), \phi(\mathbf{z}_2), \dots, \phi(\mathbf{z}_d)]$ , and  $\mathbf{c}_i \in \mathbb{R}^d$ . Both  $\mathbf{Z}$  and  $\mathbf{c}_i$  are learned from individual columns of  $\mathbf{X}$  and they can be regarded as intermediate variables avoiding the direct access of central server to  $\mathbf{x}_i$  (details of the learning will be introduced later). It follows from (2) and (3) that

$$\mathcal{K}(\mathbf{x}_i, \mathbf{x}_j) \simeq \mathbf{c}_i^T \phi(\mathbf{Z})^T \phi(\mathbf{Z}) \mathbf{c}_j. \quad (4)$$

Thus we can reconstruct the similarity between  $\mathbf{x}_i$  and  $\mathbf{x}_j$  via (4). For convenience, let  $\mathbf{K}_{xx} = \mathcal{K}(\mathbf{X}, \mathbf{X}) = \phi(\mathbf{X})^T \phi(\mathbf{X})$ ,  $\mathbf{K}_{zz} = \mathcal{K}(\mathbf{Z}, \mathbf{Z}) = \phi(\mathbf{Z})^T \phi(\mathbf{Z}) \in \mathbb{R}^{d \times d}$ , and  $\mathbf{C} = [\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_n] \in \mathbb{R}^{d \times n}$ . Then we have

$$\mathbf{K}_{xx} \simeq (\phi(\mathbf{Z})\mathbf{C})^T (\phi(\mathbf{Z})\mathbf{C}) = \mathbf{C}^T \mathbf{K}_{zz} \mathbf{C} \triangleq \hat{\mathbf{K}}_{xx}. \quad (5)$$

Now we use  $\hat{\mathbf{K}}_{xx}$  as a reconstructed similarity matrix for spectral clustering.

In the form of federated learning, we expand (3) to

$$\phi(\mathbf{X}_p) \simeq \phi(\mathbf{Z})\mathbf{C}_p, \quad p = 1, \dots, P. \quad (6)$$

It indicates that  $\mathbf{Z}$  is shared for all  $P$  clients and  $\mathbf{C}_p$  is private for client  $c_p$ . Note that (6) is a matrix factorization problem in the feature space induced by a kernel on the data in client  $c_p$ ,  $p = 1, \dots, P$ . Letting  $\mathbf{C} = [\mathbf{C}_1, \dots, \mathbf{C}_P]$ , we solve the following distributed optimization problem<sup>3</sup>

$$\underset{\mathbf{Z}, \mathbf{C}}{\text{minimize}} \quad F(\mathbf{Z}, \mathbf{C}) \triangleq \sum_{p=1}^P \omega_p f_p(\mathbf{Z}, \mathbf{C}_p). \quad (7)$$

In (7),  $f_p$  is a local objective function for client  $c_p$  and  $\omega_1, \dots, \omega_P$  are nonnegative weights for the clients. In this work, we let

$$\begin{aligned} f_p(\mathbf{Z}, \mathbf{C}_p) &= \frac{1}{2} \|\phi(\mathbf{X}_p) - \phi(\mathbf{Z})\mathbf{C}_p\|_F^2 + \frac{\lambda}{2} \|\mathbf{C}_p\|_F^2 \\ &= \frac{1}{2} \text{Tr}(\mathcal{K}(\mathbf{X}_p, \mathbf{X}_p)) - \text{Tr}(\mathbf{C}_p^T \mathcal{K}(\mathbf{Z}, \mathbf{X}_p)) + \frac{1}{2} \text{Tr}(\mathbf{C}_p^T \mathcal{K}(\mathbf{Z}, \mathbf{Z}) \mathbf{C}_p) + \frac{\lambda}{2} \|\mathbf{C}_p\|_F^2, \end{aligned} \quad (8)$$

where  $\lambda \geq 0$  is a penalty parameter. To guarantee the privacy of information, problem (7) shall be solved in the framework of federated learning.

## 2.2 FedSC by Similarity Reconstruction and Model Averaging

In this section, we develop a FedSC algorithm by similarity reconstruction and model averaging. As a classic and popular framework, FederatedAveraging (or FedAvg) is first introduced in [McMahan *et al.*, 2017] for federated learning. In our work, the proposed FedSC is, therefore, built up based on the backbone of FedAvg as in Figure 1. FedSC consists of two stages. The first stage, shown by the left plot of Figure 1, is federated similarity reconstruction, which constructs a similarity matrix in the manner of federated learning. The second stage, shown by the right plot of Figure 1, is using the reconstructed similarity matrix to implement spectral clustering.

### Stage I Federated Similarity Reconstruction

Step ①: As the startup settings for our algorithm, the shared variable  $\mathbf{Z}$  (i.e., the dictionary matrix  $\mathbf{Z}$ ) and each local coefficient matrix  $\mathbf{C}_p$  for  $p = 1, 2, \dots, P$  are initialized randomly in the central server and each client, respectively.

<sup>3</sup>Note that we do not show the data  $\mathbf{X}_p$  in the objective explicitly since it is absorbed into  $f_p$ .

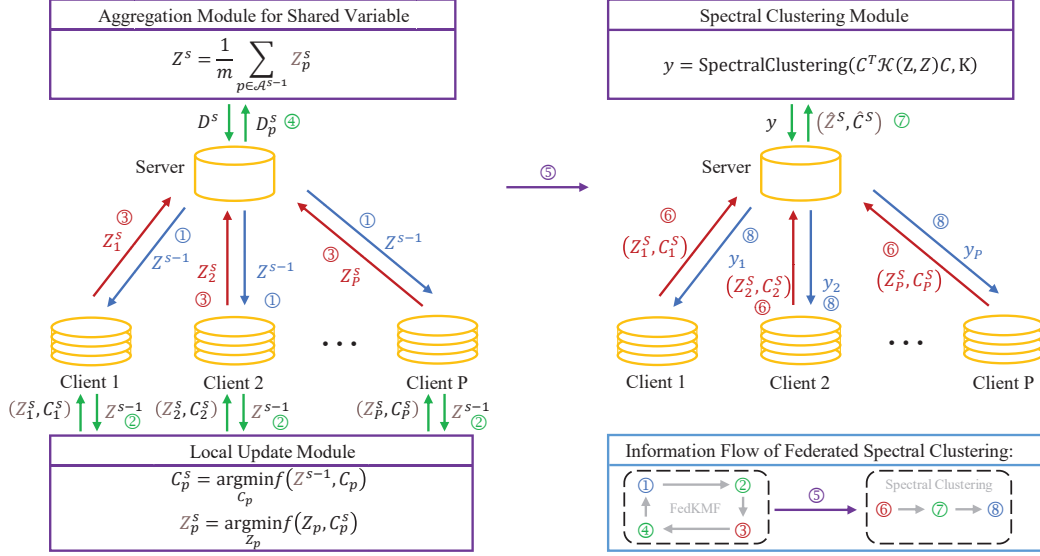


Figure 1: Diagram of the proposed FedSC. Stage I (left plot): Federated Similarity Reconstruction (Steps 1-5). Stage II (right plot): Spectral Clustering (Steps 6-8).

Step ②: For each round  $s$ , where  $1 \leq s \leq S$ , the previous shared variable  $Z$  will firstly be broadcast to each participated client. After that, every client uses this received dictionary matrix  $Z$  to run its own iterative updates of local variables in the Local Update Module (LUM) as:

$$C_p^s = \arg \min_{C_p} f_p(Z^{s-1}, C_p) = \arg \min_{C_p} \frac{1}{2} \|\phi(\mathbf{X}_p) - \phi(Z^{s-1})C_p\|_F^2 + \frac{\lambda}{2} \|C_p\|_F^2 \quad (9)$$

$$Z_p^s = \arg \min_{Z_p} f_p(Z_p, C_p^s) = \arg \min_{Z_p} \frac{1}{2} \|\phi(\mathbf{X}_p) - \phi(Z_p)C_p^s\|_F^2 + \frac{\lambda}{2} \|C_p^s\|_F^2 \quad (10)$$

Step ③: Each client sends back its own dictionary matrix  $Z_p^s$ ,  $p = 1, 2, \dots, P$ , to the central server.

Step ④: The central server collects all (or a subset  $\mathcal{A}^{s-1}$  of) these uploaded matrices  $\{Z_p^s\}_{p=1}^P$  and averages them into one new matrix  $Z^s$  in Aggregation Module (AM), i.e.,

$$Z^s = \frac{1}{|\mathcal{A}^{s-1}|} \sum_{p \in \mathcal{A}^{s-1}} Z_p^s \quad (11)$$

where  $|\mathcal{A}^{s-1}|$  is the number of participated clients. In our study, we fix the number of participating clients for each round  $s$ . Therefore, we use the notation  $\bar{P}$  instead of  $|\mathcal{A}^{s-1}|$  in the sequel. This aggregated dictionary matrix  $Z^s$  will then be used to push the next round of federated iteration until the tolerance condition is broken.

Step ⑤: When Stage I comes to an end, the spectral clustering will start.

### Stage II Spectral Clustering

Step ⑥: Each client sends  $(Z_p^s, C_p^s)$  back to the central server for the final aggregation of information.

Step ⑦: The required similarity matrix is then constructed based on the obtained dictionary matrix  $Z^s$  and coefficient matrix  $C^s$  in Spectral Clustering Module (SCM). Based on this approximated similarity matrix, the standard spectral clustering is implemented as usual:

$$y = \text{SpectralClustering}(C^T \mathcal{K}(Z, Z) C, \mathcal{K}). \quad (12)$$

Step ⑧: The central server broadcasts its clustering results to every corresponding client.

### 2.3 Optimization Algorithm for Federated Similarity Reconstruction

As described in the above section, alternate updating of local variables is a key to solving the proposed FedSC problem. In the following two parts, we discuss the optimization for  $\mathbf{Z}$  and  $\mathbf{C}$ , respectively.

For a client  $c_p$ , consider the corresponding local optimization problem

$$\underset{\mathbf{Z}, \mathbf{C}}{\text{minimize}} f_p(\mathbf{Z}, \mathbf{C}) \quad (13)$$

where  $f_p(\mathbf{Z}, \mathbf{C}) = \frac{1}{2} \|\phi(\mathbf{X}_p) - \phi(\mathbf{Z})\mathbf{C}\|_F^2 + \frac{\lambda}{2} \|\mathbf{C}\|_F^2 = \frac{1}{2} \text{Tr}(\mathcal{K}(\mathbf{X}_p, \mathbf{X}_p)) - \text{Tr}(\mathbf{C}^T \mathcal{K}(\mathbf{Z}, \mathbf{X}_p)) + \frac{1}{2} \text{Tr}(\mathbf{C}^T \mathcal{K}(\mathbf{Z}, \mathbf{Z}) \mathbf{C}) + \frac{\lambda}{2} \|\mathbf{C}\|_F^2$ . Let the derivative of  $f_p(\mathbf{Z}, \mathbf{C})$  w.r.t.  $\mathbf{C}$  be zero, we get the following one-step update for  $\mathbf{C}$ :

$$\mathbf{C}_p^s = (\mathcal{K}(\mathbf{Z}^{s-1}, \mathbf{Z}^{s-1}) + \lambda \mathbf{I}_d)^{-1} \mathcal{K}(\mathbf{Z}^{s-1}, \mathbf{X}_p), \quad p = 1, 2, \dots, P. \quad (14)$$

The derivative of  $f_p(\mathbf{Z}, \mathbf{C})$  w.r.t.  $\mathbf{Z}$  is

$$\frac{\partial \mathcal{L}}{\partial \mathbf{Z}} = \frac{1}{\sigma^2} (\mathbf{X}_p \mathbf{W}_Z - \mathbf{Z} \bar{\mathbf{W}}_Z) + \frac{2}{\sigma^2} (\mathbf{Z} \mathbf{Q}_Z - \mathbf{Z} \bar{\mathbf{Q}}_Z), \quad (15)$$

where the intermediate variables are detailed as

$$\begin{aligned} \mathbf{W}_Z &= -\mathbf{C}^T \odot \mathcal{K}(\mathbf{X}_p, \mathbf{Z}) & \bar{\mathbf{W}}_Z &= \text{diag}(\mathbf{1}_n^T \mathbf{W}_Z) \\ \mathbf{Q}_Z &= (0.5 \mathbf{C} \mathbf{C}^T) \odot \mathcal{K}(\mathbf{Z}, \mathbf{Z}) & \bar{\mathbf{Q}}_Z &= \text{diag}(\mathbf{1}_d^T \mathbf{Q}_Z). \end{aligned}$$

Here  $\mathbf{1}_n$  and  $\mathbf{1}_d$  are the column vectors with all elements of 1. Because of the kernel function,  $\mathbf{Z}$  cannot be updated like  $\mathbf{C}_p$ . Here, we use the gradient method to update it. At local iteration  $t$ , by setting  $\mathbf{Z}_p^{s,0} = \mathbf{Z}^{s-1}$ , the update scheme of  $\mathbf{Z}$  is

$$\mathbf{Z}_p^{s,t} = \mathbf{Z}_p^{s,t-1} - \eta_t \frac{\partial f_p}{\partial \mathbf{Z}}(\mathbf{Z}_p^{s,t-1}). \quad (16)$$

where  $\eta_t$  is the step size and can be set as the reverse of the Lipschitz constant of gradient if possible.

We summarize the optimization details in Algorithm 1 (shown in Appendix A).

### 2.4 Convergence Analysis of The Proposed Algorithm

First of all, it is obvious that all local objective functions  $f_p(\cdot, \cdot)$  for  $p = 1, \dots, P$  are lower bounded. To analyze the convergence of Algorithm 1, we make two assumptions. The first one is the Lipschitz continuity of the gradient of the local objective functions.

**Assumption 2.1.** The gradients of all local objective functions  $f_p(\cdot, \cdot)$  for  $p = 1, \dots, P$  are  $L_{Z_p}^s$ -Lipschitz continuous in  $\mathbf{Z}$ , that is

$$\|\nabla_{\mathbf{Z}} f_p(\mathbf{Z}^{s,t}, \mathbf{C}_p^s) - \nabla_{\mathbf{Z}} f_p(\mathbf{Z}^{s,t-1}, \mathbf{C}_p^s)\|_F \leq L_{Z_p}^s \|\mathbf{Z}^{s,t} - \mathbf{Z}^{s,t-1}\|_F. \quad (17)$$

In addition, there exist some lower and upper bounds for  $L_{Z_p}^s$ , i.e.,  $0 < \underline{L}_Z \leq L_{Z_p}^s \leq \bar{L}_Z$  hold for all  $p = 1, \dots, P$  and  $s = 1, \dots, S$ .

The second assumption, similar to [Li *et al.*, 2019; Lian *et al.*, 2017], is as follows.

**Assumption 2.2.** The difference between the local gradient and the global gradient is bounded as

$$\|\nabla_{\mathbf{Z}} f_p(\mathbf{Z}, \mathbf{C}_p) - \nabla_{\mathbf{Z}} F(\mathbf{Z}, \mathbf{C})\|_F \leq \zeta, \quad \forall p = 1, \dots, P. \quad (18)$$

To build the convergence condition, we define the following iterative terms of  $\mathbf{Z}^{s,t}$  and  $\mathbf{C}^s$  for all  $t = 1, \dots, Q$  and  $s = 1, \dots, S$ :

$$T_C(\mathbf{Z}^{s,0}, \mathbf{C}^s) = \sum_{p=1}^P \omega_p \|\mathbf{C}_p^s - \mathbf{C}_p^{s-1}\|_F^2 \quad (19)$$

$$T_Z(\mathbf{Z}^{s,t}, \mathbf{C}^s) = \|\mathbf{Z}^{s,t} - \mathbf{Z}^{s,t-1}\|_F^2$$

where the instantaneous average  $\mathbf{Z}^{s,t}$  is defined as

$$\mathbf{Z}^{s,t} = \frac{1}{P} \sum_{p \in \mathcal{A}^s} \mathbf{Z}_p^{s,t}. \quad (20)$$

Based on the above assumptions, we provide the following convergence guarantee for Algorithm 1.

**Theorem 2.3** (Convergence of Algorithm 1). *Suppose Assumption 2.1 and Assumption 2.2 hold. Let  $T = S(1 + Q)$  be the total number of global and local rounds. Then the sequence  $\{(\mathbf{Z}^{s,t}, \mathbf{C}^s)\}$  generated by Algorithm 1 with stepsize  $1/L_{Z_p}^s$  and  $\omega_p = \frac{N_p}{n}$  satisfies*

$$\frac{1}{T} \left[ \sum_{s=1}^S T_C(\mathbf{Z}^{s,0}, \mathbf{C}^s) + \sum_{s=1}^S \sum_{t=1}^Q T_Z(\mathbf{Z}^{s,t}, \mathbf{C}^s) \right] \leq \frac{D}{T} [F(\mathbf{Z}^{1,0}, \mathbf{C}^0) - f] + \frac{16\zeta^2\psi D}{P\underline{L}_Z} \quad (21)$$

where  $\psi = 1 + \frac{(P+8)(Q-1)(2Q-1)}{P-4(Q-1)^2(1+\underline{L}_Z^2/\underline{L}_Z^2)}$  and  $D = \frac{2}{\underline{\gamma}_{\min} + \lambda} + \frac{4}{\underline{L}_Z}$ .

The proof can be found in Appendix E. We see that when  $T \rightarrow \infty$ , the algorithm converges to a finite value, which is small if  $\zeta$  is small and  $\underline{L}_Z$  is close to  $\underline{L}_Z$ .

### 3 Security-Enhanced FedSC

In order to enhance the security of FedSC, we present two noise-augmented variants of the proposed algorithm in this section.

#### 3.1 FedSC with Perturbed Data

We add random noise to the data in each client and then perform Algorithm 1 to reconstruct a similarity matrix, which further improves the privacy of data. Specifically, the data  $\mathbf{X} \in \mathbb{R}^{m \times n}$  is perturbed by a noise matrix  $\mathbf{E} \in \mathbb{R}^{m \times n}$  to form the noisy data matrix

$$\tilde{\mathbf{X}} = \mathbf{X} + \mathbf{E}, \quad (22)$$

where  $E_{ij} \sim \mathcal{N}(0, \sigma^2)$ . We then perform Algorithm 1 with a Gaussian kernel of parameter  $r$  on  $\tilde{\mathbf{X}}$  and obtain  $\mathbf{Z}, \mathbf{C} = [\mathbf{C}_1, \mathbf{C}_2, \dots, \mathbf{C}_P]$ , and

$$\hat{\mathbf{K}}_{\tilde{x}\tilde{x}} = \mathbf{C}^T \mathcal{K}(\mathbf{Z}, \mathbf{Z}) \mathbf{C}. \quad (23)$$

We have the following reconstruction (for the true similarity matrix  $\mathbf{K}_{xx} = \mathcal{K}(\mathbf{X}, \mathbf{X})$ ) error bound<sup>4</sup>.

**Theorem 3.1** (Error bound of similarity matrix reconstruction). *Suppose  $\|\mathbf{X}\|_{2,\infty} = \theta$ ,  $\|\mathbf{C}\|_{2,\infty} = \tau_C$ , and  $\left\| \phi(\mathbf{Z})\mathbf{C} - \phi(\tilde{\mathbf{X}}) \right\|_{2,\infty} \leq \gamma$ , where  $\theta$ ,  $\tau_C$ , and  $\gamma$  are some nonnegative constants. Then with the probability at least  $1 - n(n-1)e^{-t}$ , the reconstructed similarity matrix  $\hat{\mathbf{K}}_{\tilde{x}\tilde{x}}$  satisfies*

$$\left\| \hat{\mathbf{K}}_{\tilde{x}\tilde{x}} - \mathbf{K}_{xx} \right\|_{\infty} \leq \frac{1}{r^2} \left[ (\sigma\xi + \sqrt{2}\theta)^2 - 2\theta^2 \right] + (\sqrt{d}\tau_C + 1)\gamma \quad (24)$$

where  $\xi = \sqrt{(m + 2\sqrt{mt} + 2t)}$ .

Note that  $r$  is the hyperparameter of the Gaussian kernel. In our experiment,  $r$  was automatically estimated as the mean of all pairwise distances between data points, i.e.,  $r = \frac{1}{n^2} \sum_{i,j} \|\mathbf{x}_i - \mathbf{x}_j\|_2$ . Assume  $|x_{ik} - x_{jk}| = \mathcal{O}(\varepsilon)$  for all  $i, j \in [n], k \in [m]$ , then  $\|\mathbf{x}_i - \mathbf{x}_j\| = \mathcal{O}(\sqrt{m}\varepsilon)$ , which means  $r^2$  is linear with  $m\varepsilon^2$ . Thus, the reconstruction error  $\|\hat{\mathbf{K}}_{\tilde{x}\tilde{x}} - \mathbf{K}_{xx}\|_{\infty}$  is upper bounded by  $\mathcal{O}(\sigma^2/\varepsilon^2 + \sqrt{d}\gamma\tau_C)$ , where  $\varepsilon^2/\sigma^2$  can be regarded as a signal-noise ratio. Therefore, the bound is useful. In general, Theorem 3.1 indicates that when the added noise is small and the optimization makes  $\gamma$  small, the reconstruction error for the true similarity matrix is less than a small constant with high probability. This verified the effectiveness of our similarity reconstruction method.

It should be pointed out that  $\hat{\mathbf{K}}_{\tilde{x}\tilde{x}}$  is not guaranteed to be a sparse matrix and hence the corresponding graph may not contain multiple connected components. We therefore use an extra KNN-based operation to get a sparse similarity matrix, which may also reduce the computational cost of eigenvalue decomposition when  $n$  is very large. Specifically, we let

$$\hat{\mathbf{K}}_{\tilde{x}\tilde{x}} = \text{getSparseMatrixbyKNN}(\hat{\mathbf{K}}_{\tilde{x}\tilde{x}}, k) \quad (25)$$

<sup>4</sup>We defer the proof for all theoretical results to the supplementary material.

which only keeps the largest  $k$  connections from each point to other points. Finally, we perform spectral clustering using  $\hat{\mathbf{K}}_{\bar{x}\bar{x}}$ . The central server broadcasts the clustering results to each participating client.

As mentioned before, one can choose to inject some noise into its raw data to avoid privacy leakage. However, a question is how much noise we can add to the data to the largest extent for the guarantee of correct clustering. We first present the following definitions.

**Definition 3.2** (Local neighbor set). Suppose  $\mathbf{x}_i$  and  $\mathbf{x}_j$  are data points of  $\mathbf{X} \in \mathbb{R}^{m \times n}$  with class labels  $L_i$  and  $L_j$  respectively, and let  $\text{KNN}(\mathbf{x}_i)$  be the set of the  $k$ -nearest neighbors of  $\mathbf{x}_i$ . We define

$$\mathcal{N}_i^{k, \text{intra}} := \{\mathbf{x}_j \in \mathbf{X} | L_i = L_j \text{ and } \mathbf{x}_j \in \text{KNN}(\mathbf{x}_i)\}. \quad (26)$$

**Definition 3.3** (Global neighbor set). Suppose  $\mathbf{x}_i$  and  $\mathbf{x}_j$  are data points of  $\mathbf{X} \in \mathbb{R}^{m \times n}$  with class labels  $L_i$  and  $L_j$  respectively, and let  $\text{KNN}(\mathbf{x}_i)$  be the point set of  $k$ -nearest neighbors of  $\mathbf{x}_i$ . We define

$$\mathcal{N}_i^{k, \text{global}} := \{\mathbf{x}_j \in \mathbf{X} | \mathbf{x}_j \in \text{KNN}(\mathbf{x}_i)\}. \quad (27)$$

If we call the local neighbor of data point  $\mathbf{x}_i$  the *intra-class neighbor* of  $\mathbf{x}_i$ , another definition called *inter-class neighbor* of  $\mathbf{x}_i$  can be further introduced as follows.

**Definition 3.4** (Inter-class neighbor set). Suppose  $\mathbf{x}_i$  and  $\mathbf{x}_j$  are data points of data matrix  $\mathbf{X} \in \mathbb{R}^{m \times n}$  with class labels,  $L_i$  and  $L_j$ , respectively, and let  $\text{KNN}(\mathbf{x}_i)$  be the point set of  $k$ -nearest neighbors of  $\mathbf{x}_i$ . We define

$$\mathcal{N}_i^{k, \text{inter}} := \{\mathbf{x}_j \in \mathbf{X} | L_i \neq L_j \text{ and } \mathbf{x}_j \in \text{KNN}(\mathbf{x}_i)\}. \quad (28)$$

Based on the above definitions, the following definition is presented to determine whether a data point can be correctly clustered or not.

**Definition 3.5** (Correct clustering). Suppose  $\mathbf{x}_i \in \mathbb{R}^m$  and  $\mathbf{x}_j \in \mathbb{R}^m$  are data points of data matrix  $\mathbf{X} \in \mathbb{R}^{m \times n}$ ,  $\mathbf{x}_i$  is said to be correctly clustered with a tolerance of  $\epsilon$  if

- a.  $\hat{\mathbf{K}}_{ij} \geq \max_k(\hat{\mathbf{K}}_{ik}^{\text{inter}}) - \epsilon$  for any of  $\mathbf{x}_j \in \mathcal{N}_i^{k, \text{intra}}$ ,
- b.  $\hat{\mathbf{K}}_{ij} \leq \min_k(\hat{\mathbf{K}}_{ik}^{\text{intra}}) + \epsilon$  for any of  $\mathbf{x}_j \in \mathcal{N}_i^{k, \text{inter}}$ .

Based on Definition 3.5, the following theorem gives the guarantee of our security-enhanced FedSC.

**Theorem 3.6** (Guarantee of noisy spectral clustering). Let  $B(\sigma) = \frac{1}{r^2} [(\sigma\xi + \sqrt{2}\theta)^2 - 2\theta^2] + (\sqrt{d}\tau_C + 1)\gamma$ . Then with the probability of at least  $1 - n(n-1)e^{-t}$ , performing spectral clustering using  $\hat{\mathbf{K}}_{\bar{x}\bar{x}}$  yields correct clustering results if

$$B(\sigma) \leq \frac{\epsilon}{2} - \max_i \frac{1}{4} \left[ \max_k(\mathbf{K}_{ik}^{\text{inter}}) - \min_k(\mathbf{K}_{ik}^{\text{intra}}) \right] \quad (29)$$

where  $\mathbf{K}_{ik}^{\text{inter}} = (\mathbf{K}_{xx})_{ik}^{\text{inter}}$  and  $\mathbf{K}_{ik}^{\text{intra}} = (\mathbf{K}_{xx})_{ik}^{\text{intra}}$ .

Based on Theorem 3.1 and Theorem 3.6, we can get a bound on the variance of noise for FedSC with perturbed data:

$$\sigma \leq \frac{1}{\xi} \left[ \sqrt{r^2(B_1 - B_2) + 2\theta^2} - \sqrt{2}\theta \right] \quad (30)$$

where  $B_1 = \frac{\epsilon}{2} - \max_i \frac{1}{4} [\max_k(\mathbf{K}_{ik}^{\text{inter}}) - \min_k(\mathbf{K}_{ik}^{\text{intra}})]$  and  $B_2 = (\sqrt{d}\tau_C + 1)\gamma$ . This bound indicates that the intensity of noise should not be too strong otherwise it may seriously affect the performance of federated spectral clustering. But, at least under this bound, one can choose to inject as much noise as possible into the raw data to ensure data security and privacy.

Using Theorem 3.22 in [Dwork *et al.*, 2014] and the post-processing property of differential privacy, we have the following privacy guarantee for this enhanced FedSC algorithm.

**Proposition 3.7.** FedSC with perturbed data given by (22) is  $(\epsilon, \delta)$ -differentially private if  $\sigma \geq 2c\tau_X/\epsilon$ , where  $c^2 > 2 \ln(1.25/\delta)$ .

Based on this proposition and (30), we obtain the following privacy-utility trade-off:

$$2\sqrt{2 \ln 1.25/\delta} \tau_X / \varepsilon < \sigma \leq \frac{1}{\xi} \left[ \sqrt{r^2(B_1 - B_2) + 2\theta^2} - \sqrt{2\theta} \right]. \quad (31)$$

This ensures both clustering performance and  $(\varepsilon, \delta)$ -differential privacy. In particular, if we substitute  $\sigma$  with the upper bound, we can get a strong level of privacy but the worst utility. By the way,  $B_1 - B_2$  is related to the property of the data. A larger  $B_1 - B_2$  means a better property for clustering, which further provides a larger upper bound for the noise level  $\sigma$ , yielding a stronger privacy guarantee.

### 3.2 FedSC with Perturbed Factors

In FedSC with perturbed factor, we added Gaussian noise to  $\mathbf{Z}$  in every round of the optimization but added Gaussian noise to  $\mathbf{C}$  in the last round of the optimization. To be more specific,  $\tilde{\mathbf{C}} = \mathbf{C} + \mathbf{E}_C$ , and  $\tilde{\mathbf{Z}} = \mathbf{Z} + \mathbf{E}_Z$ , where the entries of  $\mathbf{E}_C$  and  $\mathbf{E}_Z$  are drawn from  $\mathcal{N}(0, \sigma_C^2)$  and  $\mathcal{N}(0, \sigma_Z^2)$  respectively. Then we perform spectral clustering using the following reconstructed kernel matrix:

$$\tilde{\mathbf{K}}_{xx} = \tilde{\mathbf{C}}^T \mathcal{K}(\tilde{\mathbf{Z}}, \tilde{\mathbf{Z}}) \tilde{\mathbf{C}}. \quad (32)$$

The following theorem shows the reconstruction error bound for the ground truth kernel matrix  $\mathbf{K}_{xx}$ .

**Theorem 3.8.** *Assume  $\|\phi(\mathbf{Z})\mathbf{C} - \phi(\mathbf{X})\|_{2,\infty} \leq \gamma$ ,  $\|\mathbf{C}\|_{2,\infty} \leq \tau_C$ . Then with probability at least  $1 - (n + d)e^{-t}$ , it holds that*

$$\left\| \tilde{\mathbf{K}}_{xx} - \mathbf{K}_{xx} \right\|_{\infty} \leq \gamma_{zc}(\gamma_{zc} + 2) \quad (33)$$

where  $\gamma_{zc} = \gamma + \sqrt{d} \left( \sigma_C \xi_d + \tau_C \sqrt{2 \left( 1 - \exp\left(-\frac{\sigma_Z^2 \xi_d^2}{2r^2}\right)\right)} \right)$  and  $\xi_d^2 = d + 2\sqrt{dt} + 2t$ .

We see that, given a fixed  $\gamma$ , the reconstruction error becomes smaller if  $\sigma_Z$  and  $\sigma_C$  are smaller. Based on Theorem 3.8 and Definitions 3.2-3.5, we can obtain a bound similar to that in Theorem 3.6 to guarantee correct clustering, which will not be detailed here.

**Theorem 3.9.** *In FedSC, assume  $\max_{(p,j)} \{\|\mathbf{x}_{p_j}\|, \|\mathbf{x}'_{p_j}\|\} \leq \tau_X$ ,  $\max_{(i,j)} \|\mathbf{z}_i - \mathbf{x}_j\|_{\infty} = \Upsilon$ ,  $\|\mathbf{Z}_p^s\|_{sp} \leq \tau_Z \forall s$ , and  $\|\mathbf{C}^S\|_{2,\infty} \leq \tau_C$ , we perturb  $\{\mathbf{Z}_p^s\}_{p=1}^P$ ,  $\forall s = 1, 2, \dots, S$  with noise drawn from  $\mathcal{N}(0, \sigma_Z^2)$  with the parameter  $\sigma_Z \geq \sqrt{(8S\Delta^2(g_Z) \log(e + (\varepsilon_Z/\delta_Z)))/\varepsilon_Z^2}$  where  $\Delta(g_Z) = \frac{2\sqrt{d}\tau_C\tau_X\eta_k}{r^2} \left\{ 1 + (\tau_X + \tau_Z) \frac{(\tau_X + \Upsilon)}{r^2} \right\}$  and perturb  $\{\mathbf{C}_p^S\}_{p=1}^P$  with noise drawn from  $\mathcal{N}(0, \sigma_C^2)$  with the parameter  $\sigma_C \geq 2c\lambda^{-1}\sqrt{d}\tau_X(\tau_X + \Upsilon)/(r^2\varepsilon_C)$  for  $c^2 > 2\ln(1.25/\delta_C)$ . Then, FedSC is  $(\varepsilon_C + \varepsilon_Z, \delta_C + \delta_Z)$ -differentially private.*

Theorem 3.9 shows that our FedSC with perturbed factors can protect the data privacy provided that the noises added to  $\mathbf{C}$  and  $\mathbf{Z}$  are sufficiently large. Similarly to (31), we can also get a privacy-utility trade-off using Theorem 3.8 and Theorem 3.9, which is detailed in Appendix B.

## 4 Related Work

It should be pointed out that the study on federated spectral clustering in literature is very limited. Besides our work, the only work that aims to address the problem is [Hernández-Pereira *et al.*, 2021]. More introduction and discussion about the related work (federated matrix factorization/clustering [Yang *et al.*, 2021; Ghosh *et al.*, 2020; Dennis *et al.*, 2021; Wang and Chang, 2022] and spectral clustering [Von Luxburg, 2007; Hernández-Pereira *et al.*, 2021]) are in the supplementary material.

## 5 Experiments

### 5.1 Performance on similarity reconstruction

Taking the COIL20 dataset [Nene *et al.*, 1996] as an example, we first obtain the similarity matrix from vanilla spectral clustering based on the same kernel function. Then, we use the proposed method to derive the estimated similarity matrix  $\tilde{\mathbf{K}}_{\tilde{x}\tilde{x}}$  which is actually an approximation of ground truth. To



make it clearer, we also give the corresponding sparse similarity matrices by KNN sparsification (25). Figure 2 shows the similarity matrices constructed by different methods. We see that the proposed method can be able to successfully reconstruct the similarity matrix in the federated scenarios. The reconstruction errors on synthetic data, iris, banknote authentication, and COIL20 are in the supplementary material.

## 5.2 Clustering performance of FedSC

In this subsection, we check the clustering performance of the proposed security-enhanced FedSC method on both synthetic and real-world datasets. The synthetic dataset is generated from concentric circles. The details are in the supplementary. This synthetic dataset is visualized in Figure 3(a). Here, we continue to adopt the aforementioned COIL20 as an example of real-world datasets.

Taking the synthetic dataset as an example, the first group of cases helps illustrate the effectiveness of the proposed FedSC method. we first apply the vanilla spectral clustering method to the clean data. The predictive result is shown in Figure 3(b). It is clear that the vanilla spectral clustering method correctly clusters the data points lying in concentric circles. We then use the proposed FedSC method to cluster the data. One can find in Figure 3(c) that almost all of the data points also have been grouped correctly. However, when we inject some volume of noise

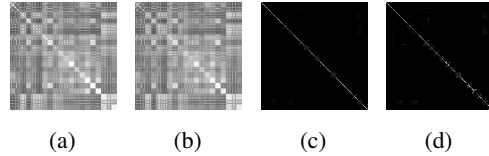


Figure 2: Visualization of similarity matrices: (a) similarity matrix of vanilla spectral clustering; (b) approximated similarity matrix of the proposed method; (c)(d) the corresponding sparse similarity matrices generated by KNN sparsification.

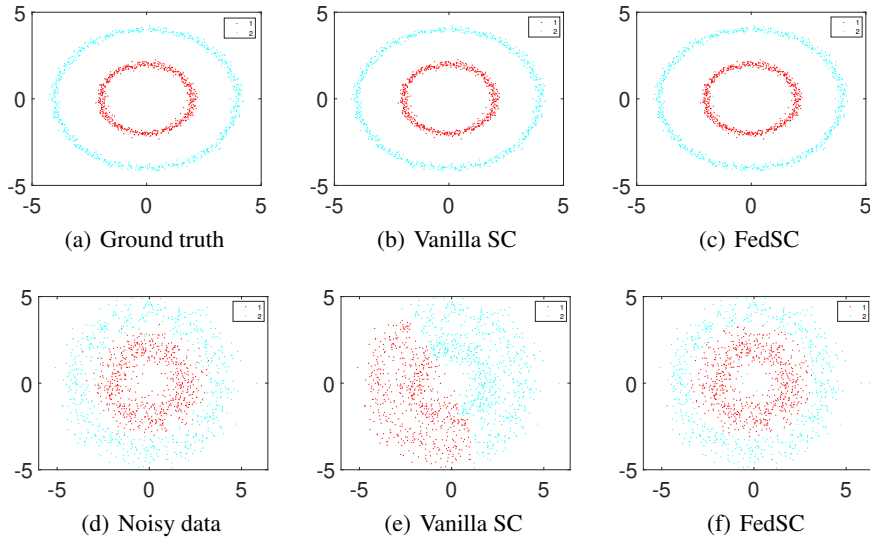


Figure 3: FedSC on concentric circles: (a) ground truth; (b) cluster assignment generated by vanilla SC; (c) cluster assignment generated by FedSC; (d) Noisy ground truth; (e) cluster assignment generated by vanilla SC on noisy data; (f) cluster assignment generated by FedSC on noisy data.

into the raw data, things may change a lot. Figure 3(d) is actually the ground truth Figure 3(a) adding some Gaussian noise. When focusing on this noisy data of concentric circles, we see from Figure 3(e) that the vanilla SC failed to cluster the data points while the proposed FedSC method is still able to correctly cluster the data to some extent as in Figure 3(f). As we know, the similarity graph directly constructed from raw data could be very sensitive to each data point. When we add too much noise, the similarity graph may fail to model the local neighborhood relationships which may be the reason why data points in Figure 3(e) are not separable for vanilla SC. Instead, FedSC is based on matrix factorization in the high-dimensional feature space and has a potential denoising effect. Therefore, it is possible for our method to achieve a better performance. The visualization of COIL20 can be found in the supplementary material.

### 5.3 Comparison with baselines

We compare our method with the clustering method DSC proposed by [Hernández-Pereira *et al.*, 2021]. Because the existing literature on federated spectral clustering is rare, we here select both classic K-means and spectral clustering as the baselines. Two metrics including accuracy and NMI are adopted to evaluate the clustering results on four datasets including iris [Dua and Graff, 2017], COIL20 [Nene *et al.*, 1996], banknote authentication [Dua and Graff, 2017], and USPS [Hull, 1994]. The details are in the supplementary material. Besides the clean data, we also consider adding noise to them to test the performance of methods under the condition of privacy protection. We directly inject Gaussian noise with zero mean and variance  $\sigma^2$  to the raw matrix  $\mathbf{X} \in \mathbb{R}^{m \times n}$  as  $\tilde{\mathbf{X}} = \mathbf{X} + \mathbf{E}$  where  $\mathbf{E}$  is a Gaussian noise matrix, each element  $E_{i,j}$  of which is i.i.d. with  $\mathcal{N}(0, \sigma^2)$ .

Table 1 shows the clustering accuracy. Our FedSC almost always achieves comparable clustering results to vanilla SC. It even outperformed vanilla SC in some cases and K-means in most cases since FedSC has a potential denoising effect by approximating a similarity matrix. More importantly, FedSC significantly outperformed DSC in almost all cases. The reason is that DSC performs spectral clustering on each local dataset, which may lead to very unstable and inaccurate results.

Table 1: Comparison of clustering accuracy ( $\mathbf{X}$  and  $\tilde{\mathbf{X}}$  denote the raw data and corrupted data respectively). The results of NMI are in Section D.4.

		Kmeans	SC	DSC	FedSC
$\mathbf{X}$	Iris	0.8933 $\pm$ 0.0000	0.9000 $\pm$ 0.0000	0.5480 $\pm$ 0.0679	0.9000 $\pm$ 0.0031
	COIL20	0.6113 $\pm$ 0.0534	0.8025 $\pm$ 0.0009	0.1009 $\pm$ 0.0100	0.7828 $\pm$ 0.0231
	Bank	0.6122 $\pm$ 0.0000	0.5918 $\pm$ 0.0000	0.5582 $\pm$ 0.0045	0.7672 $\pm$ 0.1457
	USPS	0.6704 $\pm$ 0.0047	0.6635 $\pm$ 0.0000	0.1686 $\pm$ 0.0014	0.6596 $\pm$ 0.0021
	ORL	0.6325 $\pm$ 0.0270	0.7865 $\pm$ 0.0106	0.1653 $\pm$ 0.0073	0.7235 $\pm$ 0.0170
$\tilde{\mathbf{X}}$ with $0.1\sigma$	Iris	0.8940 $\pm$ 0.0152	0.9120 $\pm$ 0.0332	0.4533 $\pm$ 0.0658	0.8993 $\pm$ 0.0299
	COIL20	0.6283 $\pm$ 0.0484	0.8024 $\pm$ 0.0020	0.0995 $\pm$ 0.0122	0.7790 $\pm$ 0.0213
	Bank	0.6067 $\pm$ 0.0022	0.6067 $\pm$ 0.0944	0.5558 $\pm$ 0.0023	0.7168 $\pm$ 0.1068
	USPS	0.6732 $\pm$ 0.0035	0.6647 $\pm$ 0.0016	0.1690 $\pm$ 0.0007	0.6643 $\pm$ 0.0022
	ORL	0.6195 $\pm$ 0.0329	0.7810 $\pm$ 0.0065	0.1600 $\pm$ 0.0089	0.7323 $\pm$ 0.0250
$\tilde{\mathbf{X}}$ with $0.3\sigma$	Iris	0.8420 $\pm$ 0.0274	0.8327 $\pm$ 0.0267	0.4533 $\pm$ 0.0674	0.8427 $\pm$ 0.0404
	COIL20	0.6422 $\pm$ 0.0366	0.7997 $\pm$ 0.0029	0.0981 $\pm$ 0.0084	0.7793 $\pm$ 0.0240
	Bank	0.6020 $\pm$ 0.0038	0.5859 $\pm$ 0.0105	0.5588 $\pm$ 0.0074	0.6046 $\pm$ 0.0064
	USPS	0.6704 $\pm$ 0.0063	0.6720 $\pm$ 0.0044	0.1673 $\pm$ 0.0019	0.6884 $\pm$ 0.0509
	ORL	0.6098 $\pm$ 0.0167	0.7885 $\pm$ 0.0047	0.1665 $\pm$ 0.0057	0.7417 $\pm$ 0.0280
$\tilde{\mathbf{X}}$ with $0.5\sigma$	Iris	0.7740 $\pm$ 0.0252	0.7313 $\pm$ 0.0494	0.3833 $\pm$ 0.0204	0.7540 $\pm$ 0.0336
	COIL20	0.6389 $\pm$ 0.0296	0.7950 $\pm$ 0.0080	0.1033 $\pm$ 0.0167	0.7403 $\pm$ 0.0294
	Bank	0.6051 $\pm$ 0.0076	0.5923 $\pm$ 0.0094	0.5566 $\pm$ 0.0030	0.6086 $\pm$ 0.0073
	USPS	0.6699 $\pm$ 0.0031	0.7843 $\pm$ 0.0030	0.1683 $\pm$ 0.0017	0.7778 $\pm$ 0.0062
	ORL	0.5983 $\pm$ 0.0295	0.7930 $\pm$ 0.0172	0.1615 $\pm$ 0.0096	0.7107 $\pm$ 0.0345
$\tilde{\mathbf{X}}$ with $0.7\sigma$	Iris	0.6500 $\pm$ 0.0420	0.6087 $\pm$ 0.0468	0.3927 $\pm$ 0.0349	0.6120 $\pm$ 0.0455
	COIL20	0.6220 $\pm$ 0.0627	0.7662 $\pm$ 0.0172	0.0893 $\pm$ 0.0055	0.6803 $\pm$ 0.0198
	Bank	0.6100 $\pm$ 0.0112	0.6046 $\pm$ 0.0144	0.5566 $\pm$ 0.0035	0.6106 $\pm$ 0.0107
	USPS	0.6638 $\pm$ 0.0044	0.7747 $\pm$ 0.0040	0.1675 $\pm$ 0.0004	0.7587 $\pm$ 0.0117
	ORL	0.5723 $\pm$ 0.0360	0.7860 $\pm$ 0.0093	0.1613 $\pm$ 0.0066	0.6910 $\pm$ 0.0232

### 5.4 More numerical result

The tSNE visualization, clustering results in terms of NMI, the performance of FedSC with perturbed factors, etc, are in the supplementary material.

## 6 Conclusion

This paper has proposed a secure kernelized factorization method for federated spectral clustering on distributed data. We provide theoretical guarantees for optimization convergence, correct clustering, and differential privacy. The numerical experiments on synthetic and real image datasets verified the effectiveness of our method. To the best knowledge of the authors, this is the work that successfully addresses the problem of federated spectral clustering. One limitation of this work is that we haven't tested our FedSC on very large datasets, though the moderate-size datasets are sufficient to justify the effectiveness of our FedSC. Note that for large-scale datasets, the bottleneck of clustering is the eigenvalue decomposition of the Laplacian matrix, not our FedSC algorithm.

## Acknowledgments

This work was partially supported by the Youth program 62106211 of the National Natural Science Foundation of China, the General Program JCYJ20210324130208022 of Shenzhen Fundamental Research, the research funding T00120210002 of Shenzhen Research Institute of Big Data, the Guangdong Key Lab of Mathematical Foundations for Artificial Intelligence, and the funding UDF01001770 of The Chinese University of Hong Kong, Shenzhen.

## References

- Jinyu Cai, Jicong Fan, Wenzhong Guo, Shiping Wang, Yunhe Zhang, and Zhao Zhang. Efficient deep embedded subspace clustering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1–10, June 2022.
- Yi-Ruei Chen, Amir Rezapour, and Wen-Guey Tzeng. Privacy-preserving ridge regression on distributed data. *Information Sciences*, 451:34–49, 2018.
- 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.
- Dheeru Dua and Casey Graff. UCI machine learning repository, 2017.
- 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.
- Jicong Fan and Tommy W.S. Chow. Sparse subspace clustering for data with missing entries and high-rank matrix completion. *Neural Networks*, 93:36–44, 2017.
- Jicong Fan and Madeleine Udell. Online high rank matrix completion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.
- Jicong Fan, Zhaoyang Tian, Mingbo Zhao, and Tommy W.S. Chow. Accelerated low-rank representation for subspace clustering and semi-supervised classification on large-scale data. *Neural Networks*, 100:39–48, 2018.
- Jicong Fan, Chengrun Yang, and Madeleine Udell. Robust non-linear matrix factorization for dictionary learning, denoising, and clustering. *IEEE Transactions on Signal Processing*, 69:1755–1770, 2021.
- Jicong Fan, Yiheng Tu, Zhao Zhang, Mingbo Zhao, and Haijun Zhang. A simple approach to automated spectral clustering. In *Advances in Neural Information Processing Systems*, volume 35, pages 9907–9921. Curran Associates, Inc., 2022.
- Jicong Fan. Large-scale subspace clustering via k-factorization. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining, KDD '21*, page 342–352, New York, NY, USA, 2021. Association for Computing Machinery.
- Avishek Ghosh, Jichan Chung, Dong Yin, and Kannan Ramchandran. An efficient framework for clustered federated learning. *Advances in Neural Information Processing Systems*, 33:19586–19597, 2020.
- Chaoyang He, Murali Annavaram, and Salman Avestimehr. Group knowledge transfer: Federated learning of large cnns at the edge. *Advances in Neural Information Processing Systems*, 33:14068–14080, 2020.
- Elena Hernández-Pereira, Oscar Fontenla-Romero, Bertha Guijarro-Berdiñas, and Beatriz Pérez-Sánchez. Federated learning approach for spectral clustering. In *ESANN*, 2021.
- Jonathan J. Hull. A database for handwritten text recognition research. *IEEE Transactions on pattern analysis and machine intelligence*, 16(5):550–554, 1994.
- Peter Kairouz, Sewoong Oh, and Pramod Viswanath. The composition theorem for differential privacy. In *International conference on machine learning*, pages 1376–1385. PMLR, 2015.

- Peter Kairouz, H Brendan McMahan, Brendan Avent, Aurélien Bellet, Mehdi Bennis, Arjun Nitin Bhagoji, Kallista Bonawitz, Zachary Charles, Graham Cormode, Rachel Cummings, et al. Advances and open problems in federated learning. *Foundations and Trends® in Machine Learning*, 14(1–2):1–210, 2021.
- Hyesung Kim, Jihong Park, Mehdi Bennis, and Seong-Lyun Kim. On-device federated learning via blockchain and its latency analysis. *arXiv preprint arXiv:1808.03949*, 2018.
- Beatrice Laurent and Pascal Massart. Adaptive estimation of a quadratic functional by model selection. *Annals of Statistics*, pages 1302–1338, 2000.
- Xiang Li, Kaixuan Huang, Wenhao Yang, Shusen Wang, and Zhihua Zhang. On the convergence of fedavg on non-iid data. *arXiv preprint arXiv:1907.02189*, 2019.
- Li Li, Yuxi Fan, Mike Tse, and Kuo-Yi Lin. A review of applications in federated learning. *Computers & Industrial Engineering*, 149:106854, 2020.
- 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.
- Zitao Li, Bolin Ding, Ce Zhang, Ninghui Li, and Jingren Zhou. Federated matrix factorization with privacy guarantee. *Proceedings of the VLDB Endowment*, 15(4):900–913, 2021.
- Xiangru Lian, Ce Zhang, Huan Zhang, Cho-Jui Hsieh, Wei Zhang, and Ji Liu. Can decentralized algorithms outperform centralized algorithms? a case study for decentralized parallel stochastic gradient descent. *Advances in Neural Information Processing Systems*, 30, 2017.
- Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas. Communication-efficient learning of deep networks from decentralized data. In *Artificial intelligence and statistics*, pages 1273–1282. PMLR, 2017.
- SA Nene, SK Nayar, and H Murase. Columbia university image library (coil-20). *Technical Report CUCS-005-96*, 1996.
- Andrew Ng, Michael Jordan, and Yair Weiss. On spectral clustering: Analysis and an algorithm. In *Advances in Neural Information Processing Systems*, volume 14. MIT Press, 2001.
- Efthymios Tzinis, Jonah Casebeer, Zhepei Wang, and Paris Smaragdis. Separate but together: Unsupervised federated learning for speech enhancement from non-iid data. In *2021 IEEE Workshop on Applications of Signal Processing to Audio and Acoustics (WASPAA)*, pages 46–50. IEEE, 2021.
- Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of machine learning research*, 9(11), 2008.
- Ulrike Von Luxburg. A tutorial on spectral clustering. *Statistics and computing*, 17(4):395–416, 2007.
- 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.
- Hongtao Wang, Ang Li, Bolin Shen, Yuyan Sun, and Hongmei Wang. Federated multi-view spectral clustering. *IEEE Access*, 8:202249–202259, 2020.
- Timothy Yang, Galen Andrew, Hubert Eichner, Haicheng Sun, Wei Li, Nicholas Kong, Daniel Ramage, and Françoise Beaufays. Applied federated learning: Improving google keyboard query suggestions. *arXiv preprint arXiv:1812.02903*, 2018.
- Qiang Yang, Yang Liu, Tianjian Chen, and Yongxin Tong. Federated machine learning: Concept and applications. *ACM Transactions on Intelligent Systems and Technology (TIST)*, 10(2):1–19, 2019.
- Enyue Yang, Yunfeng Huang, Feng Liang, Weike Pan, and Zhong Ming. Fcmf: Federated collective matrix factorization for heterogeneous collaborative filtering. *Knowledge-Based Systems*, 220:106946, 2021.

Fengda Zhang, Kun Kuang, Zhaoyang You, Tao Shen, Jun Xiao, Yin Zhang, Chao Wu, Yuet-ing Zhuang, and Xiaolin Li. Federated unsupervised representation learning. *arXiv preprint arXiv:2010.08982*, 2020.

Weiming Zhuang, Xin Gan, Yonggang Wen, Shuai Zhang, and Shuai Yi. Collaborative unsupervised visual representation learning from decentralized data. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4912–4921, 2021.

## A Framework of the Proposed Algorithm

In the beginning,  $\mathbf{Z}$  and  $C_p$  for  $p = 1, 2, \dots, P$  are initialized randomly in the central server and clients, respectively. Then,  $\mathbf{Z}$  is broadcast to every participating client and helps do the alternate updating of  $\mathbf{Z}_p$  and  $C_p$  based on the update schemes (16) and (14). Afterward, the obtained matrix  $\mathbf{Z}_p$  is sent back to the central server and aggregated for the next round of training. When the tolerance condition is broken, both  $\mathbf{Z}_p^S$  and  $C_p^S$  are sent back to the central server for the subsequent clustering task.

---

### Algorithm 1 Proposed Federated Similarity Reconstruction

---

**Input:** Distributed data  $\{\mathbf{X}_p, : p \in \mathcal{P} := \{1, 2, \dots, P\}\}$ , clients weights  $\{\omega_p : p \in \mathcal{P}\}$ .

- 1: Initialize  $\mathbf{Z}^0$  at server side and  $\{C_p^0\}_{p=1}^P$  at client sides.
  - 2: Randomly choose  $\mathcal{A}^0 \subseteq \mathcal{P}$  with  $|\mathcal{A}^s| = \bar{P}$ .
  - 3: **for** round  $s = 1$  to  $S$  **do**
  - 4:   Server side: compute  $\mathbf{Z}^{s-1} = \frac{1}{P} \sum_{p \in \mathcal{A}^{s-1}} \mathbf{Z}_p^{s-1}$ .
  - 5:   Broadcast  $\mathbf{Z}^{s-1}$  to clients  $c_p, p \in \mathcal{A}^s$ .
  - 6:   Client side:
  - 7:   **for** client  $p = 1$  to  $\bar{P}$  in parallel **do**
  - 8:     set  $\mathbf{Z}_p^{s,0} = \mathbf{Z}^{s-1}$
  - 9:     update local variable  $C_p^s$ :
  - 10:      $C_p^s = (\mathcal{K}(\mathbf{Z}_p^{s,0}, \mathbf{Z}_p^{s,0}) + \lambda \mathbf{I}_d)^{-1} \mathcal{K}(\mathbf{Z}_p^{s,0}, \mathbf{X}_p)$
  - 11:     update local variable  $\mathbf{Z}_p^s$ :
  - 12:     **for**  $t = 1$  to  $Q$  **do**
  - 13:        $\mathbf{Z}_p^{s,t} = \mathbf{Z}_p^{s,t-1} - \eta_s \nabla_{\mathbf{Z}} f_p(\mathbf{Z}_p^{s,t-1})$
  - 14:     **end for**
  - 15:     denote  $\mathbf{Z}_p^s = \mathbf{Z}_p^{s,Q}$
  - 16:     **if** client  $p \in \mathcal{A}^{s-1}$  **then**
  - 17:       upload  $\mathbf{Z}_p^s$  to the server.
  - 18:     **end if**
  - 19:   **end for**
  - 20:   Randomly choose  $\mathcal{A}^s \subseteq \mathcal{P}$  with  $|\mathcal{A}^s| = \bar{P}$ .
  - 21: **end for**
- Output:**  $\mathbf{Z}, C_p, p = 1, 2, \dots, P$
- 

## B More theoretical results about FedSC with perturbed factors

Based on Theorem 3.8 and Theorem 3.6, we can get a bound on the variance of noise for FedSC with perturbed factors:

$$\gamma_{zc}(\sigma_Z, \sigma_C) \leq -1 + 2\sqrt{1 + B_1} \quad (34)$$

where  $B_1 = \frac{\epsilon}{2} - \max_i \frac{1}{4} [\max_k(\mathbf{K}_{ik}^{inter}) - \min_k(\mathbf{K}_{ik}^{intra})]$  and  $\gamma_{zc}(\sigma_Z, \sigma_C)$  is exactly the  $\gamma_{zc}$  in Theorem 3.8 and is clearly a non-decreasing function with respect to  $\sigma_Z$  and  $\sigma_C$ , respectively. Therefore, it is a valid upper bound on both  $\sigma_Z$  and  $\sigma_C$ .

*Proof.* By Theorems 3.8 and 3.6, we have

$$\gamma_{zc}(\gamma_{zc} + 2) \leq B_1 \quad (35)$$

That is, we need to solve a quadratic equation  $\gamma_{zc}^2 + 2\gamma_{zc} - B_1 = 0$  with both  $\gamma_{zc} \geq 0$  and  $B_1 \geq 0$ . Since the discriminant  $\Delta = 4 + 4B_1 \geq 0$ , we have two roots  $-1 \pm 2\sqrt{1 + B_1}$  for this equation. Thus, it has  $0 \leq \gamma_{zc} \leq -1 + 2\sqrt{1 + B_1}$  as desired.  $\square$

This bound indicates that the intensity of noise should not be too strong otherwise it may seriously affect the performance of federated spectral clustering. But, at least under this bound, one can choose to inject as much noise as possible into the raw data to ensure data security and privacy.

Based on Theorem 3.9 and (34), we obtain the following privacy-utility trade-off:

$$\begin{cases} \gamma_{zc}(\sigma_Z, \sigma_C) \leq -1 + 2\sqrt{1 + B_1} \\ \sigma_Z \geq \sqrt{(8S\Delta^2(g_Z) \log(e + (\varepsilon_Z/\delta_Z)))/\varepsilon_Z^2)} \\ \sigma_C \geq 2c\lambda^{-1}\sqrt{d}\tau_X(\tau_X + \Upsilon)/(r^2\varepsilon_C) \end{cases} \quad (36)$$

This ensures both clustering performance and  $(\varepsilon, \delta)$ -differential privacy. In particular, if we increase the intensity of either  $\sigma_Z$  or  $\sigma_C$  to reach the upper bound of  $\gamma_{zc}$ , we can get a strong level of privacy but the worst utility. By the way,  $B_1$  is related to the property of the data. A larger  $B_1$  means a better property for clustering, which further provides a larger upper bound for the noise level, yielding a stronger privacy guarantee.

## C More discussion on the privacy-utility trade-off of FedSC

Although theoretical results are presented in our study to ensure the security of FedSC, we still provide here some insight into methods of using Secure Aggregation or other cryptographic techniques to handle pair-wise client functions (*i.e.*, kernel function in our study). Even though the communication cost will be high, it might be an alternative to reduce the DP noise levels. Specifically, it can be performed as follows.

- Step 1:  $Z_p^S$  is posted to the central server;
- Step 2: Central server aggregates  $Z_p^S$  to get the global  $Z^S$ ;
- Step 3: Central server computes  $\mathcal{K}_{ZZ} = \mathcal{K}(Z^S, Z^S)$  and broadcast it to clients;
- Step 4: For client  $p$  and  $c_i \in C_p$ , if  $c_j \in C_p$ , then client  $p$  directly calculate  $\hat{\mathcal{K}}_{ij} = c_i \mathcal{K}_{ZZ} c_j$ ; if  $c_j \in C_{p'}$ , then client  $p$  firstly encrypts and transfers its  $c_i$  to client  $p'$ , and then client  $p'$  also encrypts its  $c_j$  and use the cipher text to compute  $enc(\hat{\mathcal{K}}_{ij}) = enc(c_i^T) enc(\mathcal{K}_{ZZ}) enc(c_j)$ ; Client  $p'$  transfers the result  $enc(\hat{\mathcal{K}}_{ij})$  back to client  $p$ ; Client  $p$  decrypts  $enc(\hat{\mathcal{K}}_{ij})$  to get  $\hat{\mathcal{K}}_{ij}$ .
- Step 5: Each client  $p$  sends its estimated results  $\hat{\mathcal{K}}_{ij}$  back to the central server without sending its own  $C_p$ ;
- Step 6: Central Server checks whether these posted results from clients are compatible with each other in case of injection attacks and then performs spectral clustering.

This alternative does not send clients'  $C$  to the central server and gives an extra cross-validation process in the central server which may be useful to enhance security.

## D More details and results of the experiments

It should be pointed out that in Table 2 of the main paper, the signal-noise ratio is as high as 12dB, which means the noise is tiny. The parameter  $d$  in our FedSC was automatically determined and has a much larger value than that in the noiseless case. That is why the performance of FedSC in the noisy case is even better than that in the noiseless case of some datasets. In this appendix, we increase the noise level ( $\sigma_e = \beta\sigma$ , where  $\sigma$  denotes the standard deviation of the clean data) and consider one more real dataset.

### D.1 Dataset description

**Synthetic data** The synthetic data is generated from concentric circles. For  $\theta_i \in [0, 2\pi]$ ,  $i = 1, 2, \dots, 1258$ , all the points of this synthetic dataset  $\mathbf{X} \in \mathbb{R}^{2 \times 1258}$  are generated by

$$\mathbf{x}_i = (x_{i1}, x_{i2}) : \begin{cases} x_{i1} = r \cos(\theta_i) + e_{i1} \\ x_{i2} = r \sin(\theta_i) + e_{i2} \end{cases} \quad (37)$$

where  $\theta_0 = 0$ ,  $\theta_{1258} = 2\pi$ , and the remaining  $\theta_i$  are evenly spaced points between  $\theta_0$  and  $\theta_{1258}$ . The additive noise  $e_{i1}$  and  $e_{i2}$  are drawn from  $\mathcal{N}(0, \sigma_e^2)$ . We let  $\sigma_e = 0.1\sigma_x$ , where  $\sigma_x$  denotes the standard deviation of the data without the additive noise  $e$ . In our experiment, we set the hyperparameter  $r$  in (37) to 2 and 4.

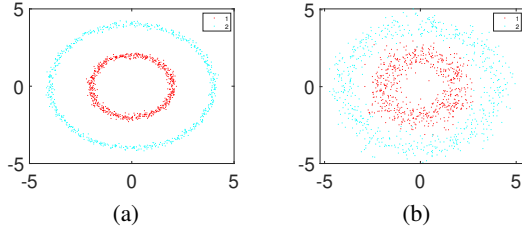


Figure 4: Synthetic dataset of 2 concentric circles: (a) ground truth; (b) noisy data perturbed by Gaussian noise with mean zero and standard deviation  $0.2\text{std}(\mathbf{X})$ .

**Real-world data** Both iris and banknote authentication are from the UCI machine learning library. COIL20 is an image dataset from Columbia Imaging and Vision Laboratory. USPS is a dataset for handwritten text recognition research. The details of the mentioned datasets are shown in Table 2.

Table 2: Summary of four real-world datasets

	# of clusters	# of attributes	# of instances
Iris	3	4	150
COIL20	20	$20 \times 20$	1440
Bank	2	5	1372
USPS	10	$16 \times 16$	9298
ORL	40	$92 \times 112$	400

The visualizations (by t-SNE [Van der Maaten and Hinton, 2008]) of three real-world datasets are shown in Figure 5.

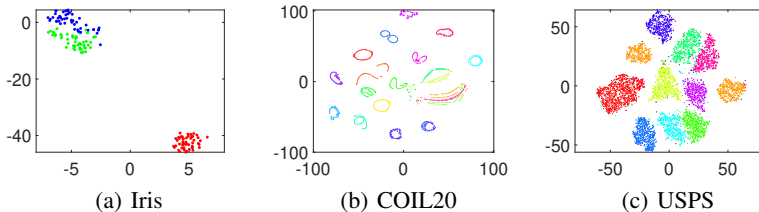


Figure 5: t-SNE visualization of some real-world datasets

## D.2 Evaluation metrics

We use two metrics to evaluate the clustering performance of our method: accuracy and normalized mutual information (NMI). Between them, accuracy is affected by the misclassification rate of cluster assignment. The smaller the misclassification rate, the greater the accuracy. In this study, given data matrix  $\mathbf{X} \in \mathbb{R}^{m \times n}$  with  $n$  sample points of  $m$  features, let  $L_i$  and  $Lr_i$  for  $i = 1, 2, \dots, n$  be the true labels and predictive labels, respectively, the accuracy of predictive cluster assignment  $Lr$  is computed by

$$\text{accuracy} = \frac{1}{n} \text{card}(\{i | L_i = Lr_i, i = 1, \dots, n\}). \quad (38)$$

NMI is a normalized version of mutual information that measures the agreement of two cluster assignments without considering their permutations.

$$\text{NMI} = \frac{2I(L; Lr)}{H(L) + H(Lr)} \quad (39)$$



where  $I(L; Lr)$  is the mutual information between  $L$  and  $Lr$ , and  $H(L)$  and  $H(Lr)$  are the entropy of  $L$  and  $Lr$ , respectively.

### D.3 Parameter settings

In our experiments, we set some hyperparameters including  $\lambda_C$ ,  $d$ ,  $r$ , and  $k$  for implementing the proposed FedSC. Among them,  $\lambda_C$  as the penalty parameter of the regularization term is set to  $1e - 2$ .  $k$  is the hyperparameter of the KNN-based operation on the similarity matrix. We set  $k$  to  $\max(\text{ceil}(\log(n)), 1)$  based on [Von Luxburg, 2007]. The details of methods to set the remaining two parameters are as follows.

**Setting of hyperparameter  $r$**  The hyperparameter  $r$  controls the smoothness of Gaussian kernel function. Due to the distinct characteristics of the datasets, we adopt the following adaptive method to determine the value of  $r$ :

$$r = c * \text{Mean}(\text{Re}(\sqrt{D})) \quad (40)$$

where  $c = 1$  and  $D$  is the distance matrix of data points, of which its element

$$D_{ij} = \|\mathbf{x}_i - \mathbf{x}_j\|_2^2, \quad \text{for } i, j = 1, 2, \dots, n. \quad (41)$$

In reality, the global  $D$  cannot be obtained because FedSC does not allow data transmission. Then we compute  $D_p$  for each client  $c_p, p = 1, \dots, P$  and then determine  $r$  as follows

$$r = \frac{c}{P} \sum_{p=1}^P \text{Mean}(\text{Re}(\sqrt{D_p})), \quad (42)$$

where  $c$  can be tuned.

**Setting of hyperparameter  $d$**  The hyperparameter  $d$  controls the complexity of the approximation  $\phi(\mathbf{X}) \simeq \phi(\mathbf{Z})\mathbf{C}$ . We adaptively determine the value of  $d$  for each dataset:

$$d = \inf_k \left\{ k \in \mathbb{N} \mid \sum_{i=1}^k s_i \geq \text{tol for } k \leq n \right\} \quad (43)$$

where  $s_i$  denotes the  $i$ -th largest eigenvalue of the kernel matrix  $\mathbf{K}_{xx}, i = 1, \dots, n$ . In our experiment, we set the threshold tol to 0.99.

It is worth noting that we use the same  $r$  in the vanilla SC and our FedSC for a fair comparison, though  $D$  cannot be obtained in our FedSC. In addition, the setting of  $d$  relies on  $\mathbf{K}_{xx}$  that is also not obtainable in our FedSC. We use this setting for convenience and in reality we need to determine  $d$  by other methods such as letting  $d = \gamma m$ , where  $\gamma$  is a hyperparameter.

### D.4 Clustering results

**NMI results** The average results of 10 repeated trials are reported in Table 1 (ACC) and Table 3 (NMI). Note that for the USPS dataset, only 5 trials are performed to save time due to its large number of sample points. We see that our FedSC outperformed Kmeans and DSC significantly in almost all cases and has at least comparable performance as SC in most cases.

**Influence of  $d$**  The clustering accuracies with different  $d$  are reported in Table 4.

**Results of FedSC with perturbed factors** The results on COIL20 are reported in Table 5, where  $\sigma_Z = \alpha_z \text{std}(\mathbf{Z}_p)$  and  $\sigma_C = \alpha_c \text{std}(\mathbf{C})$ . Through this experiment, it can be observed that perturbing factors can have a more significant impact on the accuracy of clustering results than perturbing raw data. That is, perturbing factors are more sensitive and can achieve a specified level of differential privacy with weaker noise. Furthermore, it can be seen from Table 5 that the clustering performance is more sensitive to  $\sigma_C$  than  $\sigma_Z$ .

**Malicious attack on  $\mathbf{C}$**  Due to the kernel trick, the optimization problem is nonlinear and nonconvex. Hence, it is very difficult for potential attackers to recover the data  $\mathbf{X}$  from the uploaded factors  $\mathbf{Z}, \mathbf{C}$ , especially when  $\mathbf{X}$  or  $\mathbf{Z}, \mathbf{C}$  are perturbed by noise. Nevertheless, the attacker may perform K-means on the  $\{\mathbf{C}_p\}_{p=1}^P$  to obtain clustering results. However, we find that the clustering accuracy

Table 3: Comparison with existing clustering methods (NMI)

		NMI			
		Kmeans	SC	DSC	FedSC
$X_0$	Iris	0.6356 ± 0.0000	0.6647 ± 0.0000	0.2516 ± 0.0707	0.6708 ± 0.0139
	COIL20	0.7658 ± 0.0143	0.8809 ± 0.0006	0.1259 ± 0.0225	0.8613 ± 0.0144
	Bank	0.1480 ± 0.0000	0.1228 ± 0.0000	0.0122 ± 0.0160	0.4901 ± 0.3050
	USPS	0.6125 ± 0.0026	0.8083 ± 0.0000	0.0064 ± 0.0016	0.7972 ± 0.0121
	ORL	0.8343 ± 0.0136	0.8980 ± 0.0045	0.3970 ± 0.0045	0.8525 ± 0.0097
$X_n$ with $0.1\sigma$	Iris	0.6267 ± 0.0305	0.6878 ± 0.0764	0.1458 ± 0.0794	0.6666 ± 0.0600
	COIL20	0.7721 ± 0.0122	0.8805 ± 0.0012	0.1214 ± 0.0254	0.8529 ± 0.0209
	Bank	0.1383 ± 0.0039	0.1314 ± 0.2010	0.0039 ± 0.0081	0.3853 ± 0.2390
	USPS	0.6147 ± 0.0022	0.8097 ± 0.0016	0.0076 ± 0.0010	0.7920 ± 0.0133
	ORL	0.8280 ± 0.0153	0.8948 ± 0.0019	0.3944 ± 0.0074	0.8554 ± 0.0104
$X_n$ with $0.3\sigma$	Iris	0.5292 ± 0.0528	0.5072 ± 0.0442	0.1403 ± 0.0792	0.5385 ± 0.0594
	COIL20	0.7745 ± 0.0185	0.8767 ± 0.0014	0.1174 ± 0.0224	0.8477 ± 0.0188
	Bank	0.1327 ± 0.0062	0.1031 ± 0.0211	0.0135 ± 0.0253	0.1489 ± 0.0151
	USPS	0.6112 ± 0.0052	0.7927 ± 0.0123	0.0070 ± 0.0016	0.7837 ± 0.0091
	ORL	0.8136 ± 0.0067	0.8974 ± 0.0047	0.3987 ± 0.0069	0.8578 ± 0.0147
$X_n$ with $0.5\sigma$	Iris	0.4316 ± 0.0310	0.3912 ± 0.0426	0.0710 ± 0.0421	0.4046 ± 0.0388
	COIL20	0.7676 ± 0.0148	0.8721 ± 0.0061	0.1238 ± 0.0286	0.8293 ± 0.0168
	Bank	0.1408 ± 0.0133	0.1189 ± 0.0175	0.0059 ± 0.0109	0.1521 ± 0.0113
	USPS	0.6088 ± 0.0036	0.8049 ± 0.0164	0.0059 ± 0.0028	0.7951 ± 0.0127
	ORL	0.8058 ± 0.0150	0.8970 ± 0.0063	0.3977 ± 0.0095	0.8381 ± 0.0133
$X_n$ with $0.7\sigma$	Iris	0.3025 ± 0.0407	0.2755 ± 0.0381	0.0715 ± 0.0459	0.2773 ± 0.0408
	COIL20	0.7589 ± 0.0234	0.8564 ± 0.0141	0.1032 ± 0.0127	0.7772 ± 0.0196
	Bank	0.1506 ± 0.0208	0.1420 ± 0.0262	0.0074 ± 0.0114	0.1577 ± 0.0196
	USPS	0.6007 ± 0.0051	0.8004 ± 0.0110	0.0042 ± 0.0009	0.7636 ± 0.0184
	ORL	0.7877 ± 0.0165	0.8944 ± 0.0043	0.3946 ± 0.0056	0.8236 ± 0.0102

Table 4: Accuracy of the proposed algorithm on COIL20 (K = 20)

		d	7K	8K	9K	194 (SVD)	10K	11K	12K
Clean data	Trial 1		0.7806	0.7882	0.7951	0.7944	0.7778	0.7979	0.7694
	Trial 2		0.8035	0.7431	0.7819	0.8007	0.7812	0.8049	0.7792
	Trial 3		0.7562	0.7896	0.7569	0.8097	0.7708	0.7569	0.7764
	Trial 4		0.7444	0.7771	0.7771	0.7889	0.7854	0.7903	0.7958
	Trial 5		0.7965	0.7361	0.8014	0.7833	0.7632	0.7264	0.7694
	Mean		0.7762	0.7668	0.7825	0.7954	0.7757	0.7753	0.7781
		d	10K	20K	30K	749 (SVD)	40K	50K	60K
Noisy data	Trial 1		0.7708	0.7722	0.7764	0.8201	0.7000	0.7306	0.8139
	Trial 2		0.7743	0.7889	0.7646	0.7882	0.7694	0.7451	0.7299
	Trial 3		0.7375	0.7743	0.7167	0.7493	0.7792	0.7965	0.7299
	Trial 4		0.7965	0.7688	0.8014	0.7736	0.7903	0.7382	0.7674
	Trial 5		0.7493	0.7507	0.7576	0.7000	0.7792	0.8076	0.8194
	Mean		0.7657	0.7710	0.7633	0.7662	0.7636	0.7636	0.7721

on  $C$  is lower than those of Kmeans, SC, and FedSC reported in Table 1. For instance, the clustering accuracy on Iris is reported in Table.

**Influence of  $P$**  Table 7 compares the clustering performance between using a single client ( $P = 1$ ) and using multiple ones ( $P = 8$ ). It is clear that the operation of splitting data across multiple clients may lead to an accuracy loss of clustering, which implies that our method is valid.

**More results on MNIST and CIFAR10** We have already included datasets of high-dimensional images in Table 1 like USPS and COIL20 with sizes  $16 \times 16$  and  $20 \times 20$ , respectively. Nevertheless, to further improve the experiment, we added the results of MNIST ( $28 \times 28$ ) and CIFAR10( $32 \times 32$ ) in Table 8.

Table 5: Clustering accuracy (average over 10 trials) of FedSC with perturbed factors on COIL20.

		$\alpha_c$			
		0	0.05	0.1	0.15
		Mean $\pm$ Std	Mean $\pm$ Std	Mean $\pm$ Std	Mean $\pm$ Std
$\alpha_z$	0	0.7831 $\pm$ 0.0268	0.6573 $\pm$ 0.0291	0.6012 $\pm$ 0.0391	0.4835 $\pm$ 0.0526
	0.05	0.7824 $\pm$ 0.0126	0.6573 $\pm$ 0.0243	0.6165 $\pm$ 0.0286	0.4998 $\pm$ 0.0567
	0.1	0.7817 $\pm$ 0.0299	0.6625 $\pm$ 0.0258	0.6453 $\pm$ 0.0186	0.5508 $\pm$ 0.0415
	0.15	0.7881 $\pm$ 0.0223	0.6526 $\pm$ 0.0258	0.6219 $\pm$ 0.0467	0.5901 $\pm$ 0.0326

Table 6: Clustering accuracy among different ways

	Kmeans	SC	DSC	FedSC	Attack on $C$
Iris	0.8920 $\pm$ 0.0028	0.9000 $\pm$ 0.0000	0.5493 $\pm$ 0.1263	0.9027 $\pm$ 0.0064	0.6360 $\pm$ 0.1557

## E Proof for theorem on the convergence of FedSC algorithm

We derive the objective descent with respect to  $C^s$  and  $Z^{s,t}$ , respectively.

**Objective descent with w.r.t.  $C$ :** Based on the update scheme (14) of  $C$ , we have

$$\nabla_C f_p(Z^{s,0}, C_p^s) = (\mathcal{K}(Z^{s,0}, Z^{s,0}) + \lambda I_d)C_p^s - \mathcal{K}(Z^{s,0}, X_p) = 0 \quad (44)$$

where  $Z^{s,0} = Z^{s-1} = \frac{1}{P} \sum_{p \in \mathcal{A}^{s-1}} Z_p^{s-1, Q}$ .

According to the proposed FedSC problem 7, we have

$$\begin{aligned} & f_p(Z^{s,0}, C_p^s) - f_p(Z^{s,0}, C_p^{s-1}) \\ &= \left[ \frac{1}{2} \|\phi(X_p) - \phi(Z^{s,0})C_p^s\|_F^2 + \frac{\lambda}{2} \|C_p^s\|_F^2 \right] \\ & \quad - \left[ \frac{1}{2} \|\phi(X_p) - \phi(Z^{s,0})C_p^{s-1}\|_F^2 + \frac{\lambda}{2} \|C_p^{s-1}\|_F^2 \right] \\ &= -\text{Tr}((C_p^s - C_p^{s-1})^T \mathcal{K}(Z^{s,0}, X_p)) + \frac{\lambda}{2} \text{Tr}(C_p^s (C_p^s)^T - C_p^{s-1} (C_p^{s-1})^T) \\ & \quad + \frac{1}{2} \text{Tr}([C_p^s (C_p^s)^T - C_p^{s-1} (C_p^{s-1})^T] \mathcal{K}(Z^{s,0}, Z^{s,0})) \\ &= -\text{Tr}((C_p^s - C_p^{s-1})^T (\mathcal{K}(Z^{s,0}, Z^{s,0}) + \lambda I_d) C_p^s) \\ & \quad + \frac{1}{2} \text{Tr}([C_p^s (C_p^s)^T - C_p^{s-1} (C_p^{s-1})^T] [\mathcal{K}(Z^{s,0}, Z^{s,0}) + \lambda I_d]) \\ &= -\frac{1}{2} \underbrace{\text{Tr}([C_p^s - C_p^{s-1}]^T [\mathcal{K}(Z^{s,0}, Z^{s,0}) + \lambda I_d] [C_p^s - C_p^{s-1}])}_{\text{T.1}} \\ &\leq -\frac{1}{2} (\gamma_{min}^s + \lambda) \|C_p^s - C_p^{s-1}\|_F^2 \end{aligned} \quad (45)$$

where  $\gamma_{min}^s = \gamma_{min}(\mathcal{K}(Z^{s,0}, Z^{s,0}))$ .

Summing it up from  $p = 1$  to  $P$ , we have

$$F(Z^{s,0}, C^s) - F(Z^{s,0}, C^{s-1}) \leq -\frac{1}{2} (\gamma_{min}^s + \lambda) \sum_{p=1}^P \omega_p \|C_p^s - C_p^{s-1}\|_F^2 \quad (46)$$

**Objective descent with w.r.t.  $Z$ :** According to Assumption 2.1, it implies

$$\begin{aligned} F(Z^{s,t}, C^s) - F(Z^{s,t-1}, C^s) &\leq \underbrace{\langle \nabla_Z F(Z^{s,t-1}, C^s), Z^{s,t} - Z^{s,t-1} \rangle}_{\text{T.2}} \\ &\quad + \frac{L_Z^s}{2} \|Z^{s,t} - Z^{s,t-1}\|_F^2 \end{aligned} \quad (47)$$

Table 7: Clustering accuracy between using a single client and using multiple clients

	Iris	COIL20	Bank	ORL
$P = 1$	$0.9007 \pm 0.0086$	$0.8073 \pm 0.0089$	$0.7536 \pm 0.1271$	$0.7937 \pm 0.0084$
$P = 8$	$0.8993 \pm 0.0165$	$0.8003 \pm 0.0049$	$0.6821 \pm 0.1405$	$0.7112 \pm 0.0258$

Table 8: Clustering accuracy on MNIST and CIFAR10

		Kmeans	SC	DSC	FedSC
$\mathbf{X}$	MNIST	$0.5448 \pm 0.0257$	$0.6265 \pm 0.0439$	$0.1292 \pm 0.0144$	$0.6139 \pm 0.0464$
	CIFAR10	$0.2171 \pm 0.0132$	$0.2182 \pm 0.0133$	$0.1235 \pm 0.0062$	$0.2134 \pm 0.0131$
$\mathbf{X}$ with $0.3\sigma$	MNIST	$0.5402 \pm 0.0225$	$0.5755 \pm 0.0366$	$0.1337 \pm 0.0191$	$0.5606 \pm 0.0457$
	CIFAR10	$0.2209 \pm 0.0154$	$0.2198 \pm 0.0172$	$0.1194 \pm 0.0051$	$0.2187 \pm 0.0125$
$\mathbf{X}$ with $0.7\sigma$	MNIST	$0.5374 \pm 0.0555$	$0.5711 \pm 0.0329$	$0.1340 \pm 0.0135$	$0.5029 \pm 0.0264$
	CIFAR10	$0.2202 \pm 0.0147$	$0.2205 \pm 0.0084$	$0.1209 \pm 0.0045$	$0.2134 \pm 0.0152$

Now, we give a bound on T.2.

**Lemma E.1.** For any  $s$  and  $t$ , it holds that

$$\begin{aligned} \langle \nabla_Z F(\mathbf{Z}^{s,t-1}, \mathbf{C}^s), \mathbf{Z}^{s,t} - \mathbf{Z}^{s,t-1} \rangle &= -L_Z^s \|\mathbf{Z}^{s,t} - \mathbf{Z}^{s,t-1}\|_F^2 \\ &\quad + \langle \nabla_Z F(\mathbf{Z}^{s,t-1}, \mathbf{C}^s), \mathbf{Z}^{s,t} - \mathbf{Z}^{s,t-1} \rangle \end{aligned} \quad (48)$$

*Proof.* Based on the update schemes of  $\mathbf{Z}$ , we have

$$\begin{aligned} \mathbf{Z}_p^{s,t} &= \mathbf{Z}_p^{s,t-1} - \frac{1}{L_Z^s} \nabla_Z f_p(\mathbf{Z}_p^{s,t-1}, \mathbf{C}_p^s) \\ \iff \mathbf{Z}^{s,t} &= \mathbf{Z}^{s,t-1} - \frac{1}{\bar{P}L_Z^s} \sum_{p \in \mathcal{A}^s} \nabla_Z f_p(\mathbf{Z}_p^{s,t-1}, \mathbf{C}_p^s) \\ \iff 0 &= \frac{1}{\bar{P}} \sum_{p \in \mathcal{A}^s} \nabla_Z f_p(\mathbf{Z}_p^{s,t-1}, \mathbf{C}_p^s) + L_Z^s (\mathbf{Z}^{s,t} - \mathbf{Z}^{s,t-1}) \end{aligned} \quad (49)$$

Then, consider the following terms

$$\begin{aligned} &\langle \nabla_Z F(\mathbf{Z}^{s,t-1}, \mathbf{C}^s), \mathbf{Z}^{s,t} - \mathbf{Z}^{s,t-1} \rangle \\ &= \langle \nabla_Z F(\mathbf{Z}^{s,t-1}, \mathbf{C}^s) - \frac{1}{\bar{P}} \sum_{p \in \mathcal{A}^s} \nabla_Z f_p(\mathbf{Z}_p^{s,t-1}, \mathbf{C}_p^s) - L_Z^s (\mathbf{Z}^{s,t} - \mathbf{Z}^{s,t-1}), \mathbf{Z}^{s,t} - \mathbf{Z}^{s,t-1} \rangle \\ &= -L_Z^s \|\mathbf{Z}^{s,t} - \mathbf{Z}^{s,t-1}\|_F^2 \\ &\quad + \langle \nabla_Z F(\mathbf{Z}^{s,t-1}, \mathbf{C}^s) - \frac{1}{\bar{P}} \sum_{p \in \mathcal{A}^s} \nabla_Z f_p(\mathbf{Z}_p^{s,t-1}, \mathbf{C}_p^s), \mathbf{Z}^{s,t} - \mathbf{Z}^{s,t-1} \rangle \end{aligned} \quad (50)$$

□

Based on Lemma 1, we continue to do the following derivation.

$$\begin{aligned} &F(\mathbf{Z}^{s,t}, \mathbf{C}^s) - F(\mathbf{Z}^{s,t-1}, \mathbf{C}^s) \\ &\leq -\frac{L_Z^s}{2} \|\mathbf{Z}^{s,t} - \mathbf{Z}^{s,t-1}\|_F^2 \\ &\quad + \langle \nabla_Z F(\mathbf{Z}^{s,t-1}, \mathbf{C}^s) - \frac{1}{\bar{P}} \sum_{p \in \mathcal{A}^s} \nabla_Z f_p(\mathbf{Z}_p^{s,t-1}, \mathbf{C}_p^s), \mathbf{Z}^{s,t} - \mathbf{Z}^{s,t-1} \rangle \\ &\leq -\frac{L_Z^s}{4} \|\mathbf{Z}^{s,t} - \mathbf{Z}^{s,t-1}\|_F^2 \\ &\quad + \underbrace{\frac{1}{L_Z^s} \left\| \nabla_Z F(\mathbf{Z}^{s,t-1}, \mathbf{C}^s) - \frac{1}{\bar{P}} \sum_{p \in \mathcal{A}^s} \nabla_Z f_p(\mathbf{Z}_p^{s,t-1}, \mathbf{C}_p^s) \right\|_F^2}_{\text{T.3}} \end{aligned} \quad (51)$$

where we used the fact that  $\langle \mathbf{x}, \mathbf{y} \rangle \leq \|\mathbf{x}\|_2 + \|\mathbf{y}\|_2, \forall c > 0$ .

Now, we give a bound on T.3.

**Lemma E.2.** *For any  $s$  and  $t$ , it holds that*

$$\begin{aligned} & \left\| \nabla_Z F(\mathbf{Z}^{s,t-1}, \mathbf{C}^s) - \frac{1}{\bar{P}} \sum_{p \in \mathcal{A}^s} \nabla_Z f_p(\mathbf{Z}_p^{s,t-1}, \mathbf{C}_p^s) \right\|_F^2 \\ & \leq \frac{16\zeta^2}{\bar{P}} + 2\left(1 + \frac{8}{\bar{P}}\right) \sum_{p=1}^P \omega_p (L_{Z_p^s}^s)^2 \|\mathbf{Z}^{s,t-1} - \mathbf{Z}_p^{s,t-1}\|_F^2 \end{aligned} \quad (52)$$

*Proof.* For any  $s$  and  $t$ , it holds that

$$\begin{aligned} & \left\| \nabla_Z F(\mathbf{Z}^{s,t-1}, \mathbf{C}^s) - \frac{1}{\bar{P}} \sum_{p \in \mathcal{A}^s} \nabla_Z f_p(\mathbf{Z}_p^{s,t-1}, \mathbf{C}_p^s) \right\|_F^2 \\ & = \left\| \nabla_Z F(\mathbf{Z}^{s,t-1}, \mathbf{C}^s) - \sum_{p=1}^P \omega_p \nabla_Z f_p(\mathbf{Z}_p^{s,t-1}, \mathbf{C}_p^s) \right. \\ & \quad \left. + \sum_{p=1}^P \omega_p \nabla_Z f_p(\mathbf{Z}_p^{s,t-1}, \mathbf{C}_p^s) - \frac{1}{\bar{P}} \sum_{p \in \mathcal{A}^s} \nabla_Z f_p(\mathbf{Z}_p^{s,t-1}, \mathbf{C}_p^s) \right\|_F^2 \\ & \leq 2 \underbrace{\left\| \nabla_Z F(\mathbf{Z}^{s,t-1}, \mathbf{C}^s) - \sum_{p=1}^P \omega_p \nabla_Z f_p(\mathbf{Z}_p^{s,t-1}, \mathbf{C}_p^s) \right\|_F^2}_{\text{T.4}} \\ & \quad + 2 \underbrace{\left\| \sum_{p=1}^P \omega_p \nabla_Z f_p(\mathbf{Z}_p^{s,t-1}, \mathbf{C}_p^s) - \frac{1}{\bar{P}} \sum_{p \in \mathcal{A}^s} \nabla_Z f_p(\mathbf{Z}_p^{s,t-1}, \mathbf{C}_p^s) \right\|_F^2}_{\text{T.5}} \end{aligned} \quad (53)$$

where we used the fact that  $\|\sum_{i=1}^n a_i\|_2^2 \leq n \sum_{i=1}^n \|a_i\|_2^2$ .

Firstly, we give a bound on T.4:

$$\begin{aligned} & \left\| \nabla_Z F(\mathbf{Z}^{s,t-1}, \mathbf{C}^s) - \sum_{p=1}^P \omega_p \nabla_Z f_p(\mathbf{Z}_p^{s,t-1}, \mathbf{C}_p^s) \right\|_F^2 \\ & = \left\| \sum_{p=1}^P \omega_p [\nabla_Z f(\mathbf{Z}^{s,t-1}, \mathbf{C}_p^s) - \nabla_Z f_p(\mathbf{Z}_p^{s,t-1}, \mathbf{C}_p^s)] \right\|_F^2 \\ & \leq \sum_{p=1}^P \omega_p \|\nabla_Z F(\mathbf{Z}^{s,t-1}, \mathbf{C}_p^s) - \nabla_Z f_p(\mathbf{Z}_p^{s,t-1}, \mathbf{C}_p^s)\|_F^2 \\ & \leq \sum_{p=1}^P \omega_p (L_{Z_p^s}^s)^2 \|\mathbf{Z}^{s,t-1} - \mathbf{Z}_p^{s,t-1}\|_F^2 \end{aligned} \quad (54)$$

where we used the fact that  $\nabla_Z f_p(\cdot, \mathbf{C}_p^s)$  is  $L_{Z_p^s}^s$ -Lipschitz continuous.

Secondly, we give a bound on T.5:

$$\begin{aligned}
& \left\| \sum_{p=1}^P \omega_p \nabla_Z f_p(\mathbf{Z}_p^{s,t-1}, \mathbf{C}_p^s) - \frac{1}{\bar{P}} \sum_{p \in \mathcal{A}^s} \nabla_Z f_p(\mathbf{Z}_p^{s,t-1}, \mathbf{C}_p^s) \right\|_F^2 \\
&= \frac{1}{\bar{P}^2} \left\| \sum_{p' \in \mathcal{A}^s} \left[ \sum_{p=1}^P \omega_p \nabla_Z f_p(\mathbf{Z}_p^{s,t-1}, \mathbf{C}_p^s) - \nabla_Z f_{p'}(\mathbf{Z}_{p'}^{s,t-1}, \mathbf{C}_{p'}^s) \right] \right\|_F^2 \\
&\leq \frac{1}{\bar{P}} \sum_{p' \in \mathcal{A}^s} \left\| \sum_{p=1}^P \omega_p \nabla_Z f_p(\mathbf{Z}_p^{s,t-1}, \mathbf{C}_p^s) - \nabla_Z f_{p'}(\mathbf{Z}_{p'}^{s,t-1}, \mathbf{C}_{p'}^s) \right\|_F^2 \tag{55} \\
&= \frac{1}{\bar{P}} \sum_{p'=1}^P \omega_{p'} \left\| \sum_{p=1}^P \omega_p \nabla_Z f_p(\mathbf{Z}_p^{s,t-1}, \mathbf{C}_p^s) - \nabla_Z f_{p'}(\mathbf{Z}_{p'}^{s,t-1}, \mathbf{C}_{p'}^s) \right\|_F^2 \\
&\leq \frac{1}{\bar{P}} \sum_{p'=1}^P \omega_{p'} \sum_{p=1}^P \omega_p \underbrace{\left\| \nabla_Z f_p(\mathbf{Z}_p^{s,t-1}, \mathbf{C}_p^s) - \nabla_Z f_{p'}(\mathbf{Z}_{p'}^{s,t-1}, \mathbf{C}_{p'}^s) \right\|_F^2}_{\text{T.6}}
\end{aligned}$$

Thirdly, we give a bound on T.6:

$$\begin{aligned}
& \left\| \nabla_Z f_p(\mathbf{Z}_p^{s,t-1}, \mathbf{C}_p^s) - \nabla_Z f_{p'}(\mathbf{Z}_{p'}^{s,t-1}, \mathbf{C}_{p'}^s) \right\|_F^2 \\
&= \left\| \left[ \nabla_Z f_p(\mathbf{Z}_p^{s,t-1}, \mathbf{C}_p^s) - \nabla_Z f_p(\mathbf{Z}^{s,t-1}, \mathbf{C}_p^s) + \nabla_Z f_p(\mathbf{Z}^{s,t-1}, \mathbf{C}_p^s) \right. \right. \\
&\quad \left. \left. - \nabla_Z F(\mathbf{Z}^{s,t-1}, \mathbf{C}^s) + \nabla_Z F(\mathbf{Z}^{s,t-1}, \mathbf{C}^s) \right] \right. \\
&\quad \left. - \left[ \nabla_Z f_{p'}(\mathbf{Z}_{p'}^{s,t-1}, \mathbf{C}_{p'}^s) - \nabla_Z f_{p'}(\mathbf{Z}^{s,t-1}, \mathbf{C}_{p'}^s) + \nabla_Z f_{p'}(\mathbf{Z}^{s,t-1}, \mathbf{C}_{p'}^s) \right. \right. \\
&\quad \left. \left. - \nabla_Z F(\mathbf{Z}^{s,t-1}, \mathbf{C}^s) + \nabla_Z F(\mathbf{Z}^{s,t-1}, \mathbf{C}^s) \right] \right\|_F^2 \\
&= \left\| \left[ \nabla_Z f_p(\mathbf{Z}_p^{s,t-1}, \mathbf{C}_p^s) - \nabla_Z f_p(\mathbf{Z}^{s,t-1}, \mathbf{C}_p^s) \right] \right. \\
&\quad \left. + \left[ \nabla_Z f_p(\mathbf{Z}^{s,t-1}, \mathbf{C}_p^s) - \nabla_Z F(\mathbf{Z}^{s,t-1}, \mathbf{C}^s) \right] \right. \\
&\quad \left. + \left[ \nabla_Z f_{p'}(\mathbf{Z}_{p'}^{s,t-1}, \mathbf{C}_{p'}^s) - \nabla_Z f_{p'}(\mathbf{Z}^{s,t-1}, \mathbf{C}_{p'}^s) \right] \right. \\
&\quad \left. + \left[ \nabla_Z f_{p'}(\mathbf{Z}^{s,t-1}, \mathbf{C}_{p'}^s) - \nabla_Z F(\mathbf{Z}^{s,t-1}, \mathbf{C}^s) \right] \right\|_F^2 \tag{56} \\
&\leq 4 \left\| \nabla_Z f_p(\mathbf{Z}_p^{s,t-1}, \mathbf{C}_p^s) - \nabla_Z f_p(\mathbf{Z}^{s,t-1}, \mathbf{C}_p^s) \right\|_F^2 \\
&\quad + 4 \left\| \nabla_Z f_p(\mathbf{Z}^{s,t-1}, \mathbf{C}_p^s) - \nabla_Z F(\mathbf{Z}^{s,t-1}, \mathbf{C}^s) \right\|_F^2 \\
&\quad + 4 \left\| \nabla_Z f_{p'}(\mathbf{Z}_{p'}^{s,t-1}, \mathbf{C}_{p'}^s) - \nabla_Z f_{p'}(\mathbf{Z}^{s,t-1}, \mathbf{C}_{p'}^s) \right\|_F^2 \\
&\quad + 4 \left\| \nabla_Z f_{p'}(\mathbf{Z}^{s,t-1}, \mathbf{C}_{p'}^s) - \nabla_Z F(\mathbf{Z}^{s,t-1}, \mathbf{C}^s) \right\|_F^2 \\
&\leq 4(L_{Z_p}^s)^2 \left\| \mathbf{Z}^{s,t-1} - \mathbf{Z}_p^{s,t-1} \right\|_F^2 + 4(L_{Z_{p'}}^s)^2 \left\| \mathbf{Z}^{s,t-1} - \mathbf{Z}_{p'}^{s,t-1} \right\|_F^2 + 8\zeta^2
\end{aligned}$$

Here, it is the second time that we use the facts that  $\left\| \sum_{i=1}^n a_i \right\|_2^2 \leq n \sum_{i=1}^n \|a_i\|_2^2$  and that  $\nabla_Z f_p(\cdot, \mathbf{C}_p^s)$  is  $L_{Z_p}^s$ -Lipschitz continuous.

Based on the bound of T.6, we derive the bound on T.5.

$$\begin{aligned}
& \left\| \sum_{p=1}^P \omega_p \nabla_Z f_p(\mathbf{Z}_p^{s,t-1}, \mathbf{C}_p^s) - \frac{1}{\bar{P}} \sum_{p \in \mathcal{A}^s} \nabla_Z f_p(\mathbf{Z}_p^{s,t-1}, \mathbf{C}_p^s) \right\|_F^2 \\
& \leq \frac{1}{\bar{P}} \sum_{p'=1}^P \omega_{p'} \sum_{p=1}^P \omega_p \underbrace{\left[ 4(L_{Z_p}^s)^2 \|\mathbf{Z}^{s,t-1} - \mathbf{Z}_p^{s,t-1}\|_F^2 + 4(L_{Z_{p'}}^s)^2 \|\mathbf{Z}^{s,t-1} - \mathbf{Z}_{p'}^{s,t-1}\|_F^2 + 8\zeta^2 \right]}_{\text{Bound of T.5}} \\
& \leq \frac{8\zeta^2}{\bar{P}} + \frac{8}{\bar{P}} \sum_{p=1}^P \omega_p (L_{Z_p}^s)^2 \|\mathbf{Z}^{s,t-1} - \mathbf{Z}_p^{s,t-1}\|_F^2
\end{aligned} \tag{57}$$

Fourthly, based on the bounds of T.4 and T.5, we continue to derive the final required bound.

$$\begin{aligned}
& \left\| \nabla_Z F(\mathbf{Z}^{s,t-1}, \mathbf{C}^s) - \frac{1}{\bar{P}} \sum_{p \in \mathcal{A}^s} \nabla_Z f_p(\mathbf{Z}_p^{s,t-1}, \mathbf{C}_p^s) \right\|_F^2 \\
& \leq 2 \underbrace{\sum_{p=1}^P \omega_p (L_{Z_p}^s)^2 \|\mathbf{Z}^{s,t-1} - \mathbf{Z}_p^{s,t-1}\|_F^2}_{\text{Bound of T.3}} + 2 \underbrace{\left[ \frac{8\zeta^2}{\bar{P}} + \frac{8}{\bar{P}} \sum_{p=1}^P \omega_p (L_{Z_p}^s)^2 \|\mathbf{Z}^{s,t-1} - \mathbf{Z}_p^{s,t-1}\|_F^2 \right]}_{\text{Bound of T.4}} \\
& = \frac{16\zeta^2}{\bar{P}} + 2\left(1 + \frac{8}{\bar{P}}\right) \sum_{p=1}^P \omega_p (L_{Z_p}^s)^2 \|\mathbf{Z}^{s,t-1} - \mathbf{Z}_p^{s,t-1}\|_F^2
\end{aligned} \tag{58}$$

□

Based on Lemma 2, we continue to do the following derivation.

$$\begin{aligned}
& F(\mathbf{Z}^{s,t}, \mathbf{C}^s) - F(\mathbf{Z}^{s,t-1}, \mathbf{C}^s) \\
& \leq -\frac{L_Z^s}{4} \|\mathbf{Z}^{s,t} - \mathbf{Z}^{s,t-1}\|_F^2 \\
& \quad + \frac{1}{L_Z^s} \underbrace{\left[ \frac{16\zeta^2}{\bar{P}} + 2\left(1 + \frac{8}{\bar{P}}\right) \sum_{p=1}^P \omega_p (L_{Z_p}^s)^2 \|\mathbf{Z}^{s,t-1} - \mathbf{Z}_p^{s,t-1}\|_F^2 \right]}_{\text{Bound of T.3}} \\
& = -\frac{L_Z^s}{4} \|\mathbf{Z}^{s,t} - \mathbf{Z}^{s,t-1}\|_F^2 + \frac{16\zeta^2}{\bar{P}L_Z^s} \\
& \quad + \frac{2}{L_Z^s} \left(1 + \frac{8}{\bar{P}}\right) \sum_{p=1}^P \omega_p (L_{Z_p}^s)^2 \|\mathbf{Z}^{s,t-1} - \mathbf{Z}_p^{s,t-1}\|_F^2
\end{aligned} \tag{59}$$

Then, summing it up from  $t = 1$  to  $Q$  yields

$$\begin{aligned}
& F(\mathbf{Z}^{s,Q}, \mathbf{C}^s) - F(\mathbf{Z}^{s,0}, \mathbf{C}^s) \\
& \leq -\frac{L_Z^s}{4} \sum_{t=1}^Q \|\mathbf{Z}^{s,t} - \mathbf{Z}^{s,t-1}\|_F^2 + \frac{16Q\zeta^2}{\bar{P}L_Z^s} \\
& \quad + \frac{2}{L_Z^s} \left(1 + \frac{8}{\bar{P}}\right) \underbrace{\sum_{t=1}^Q \sum_{p=1}^P \omega_p (L_{Z_p}^s)^2 \|\mathbf{Z}^{s,t-1} - \mathbf{Z}_p^{s,t-1}\|_F^2}_{\text{T.7}}
\end{aligned} \tag{60}$$

Now, we give a bound on T.7.

**Lemma E.3.** For any  $s$ , it holds that

$$\begin{aligned}
& \left\| \mathbf{Z}^{s,t} - \mathbf{Z}_p^{s,t} \right\|_F^2 \\
& \leq \frac{4t}{(L_Z^s)^2 \bar{P}} \sum_{j=0}^{t-1} (L_{Z_p}^s)^2 \left\| \mathbf{Z}^{s,j} - \mathbf{Z}_p^{s,j} \right\|_F^2 + \frac{8t^2 \zeta^2}{(L_Z^s)^2 \bar{P}} \\
& \quad + \frac{4t}{(L_Z^s)^2 \bar{P}} \sum_{j=0}^{t-1} \sum_{p'=1}^P \omega_{p'}(L_{Z_{p'}}^s)^2 \left\| \mathbf{Z}^{s,j} - \mathbf{Z}_{p'}^{s,j} \right\|_F^2
\end{aligned} \tag{61}$$

*Proof.* Based on the update schemes of  $\mathbf{C}$  and  $\mathbf{Z}$ , we have

$$\begin{aligned}
\mathbf{Z}_p^{s,t} &= \mathbf{Z}_p^{s,t-1} - \frac{1}{L_Z^s} \nabla_{\mathbf{Z}} f_p(\mathbf{Z}_p^{s,t-1}, \mathbf{C}_p^s) \\
\iff \mathbf{Z}^{s,t} &= \mathbf{Z}^{s,t-1} - \frac{1}{\bar{P} L_Z^s} \sum_{p \in \mathcal{A}^s} \nabla_{\mathbf{Z}} f_p(\mathbf{Z}_p^{s,t-1}, \mathbf{C}_p^s)
\end{aligned} \tag{62}$$

Consequently, we have

$$\begin{aligned}
\mathbf{Z}_p^{s,t} &= \mathbf{Z}_p^{s,0} - \frac{1}{L_Z^s} \sum_{j=0}^{t-1} \nabla_{\mathbf{Z}} f_p(\mathbf{Z}_p^{s,j}, \mathbf{C}_p^s) \\
\iff \mathbf{Z}^{s,t} &= \mathbf{Z}^{s,0} - \frac{1}{\bar{P} L_Z^s} \sum_{j=0}^{t-1} \sum_{p \in \mathcal{A}^s} \nabla_{\mathbf{Z}} f_p(\mathbf{Z}_p^{s,j}, \mathbf{C}_p^s)
\end{aligned} \tag{63}$$

where  $\mathbf{Z}_p^{s,0} = \mathbf{Z}^{s,0} = \mathbf{Z}^{s-1}$ .



Based on the above identities,

$$\begin{aligned}
& \|\mathbf{Z}^{s,t} - \mathbf{Z}_p^{s,t}\|_F^2 \\
&= \left\| \left[ \mathbf{Z}^{s,0} - \frac{1}{\bar{P}L_Z^s} \sum_{j=0}^{t-1} \sum_{p \in \mathcal{A}^s} \nabla_Z f_p(\mathbf{Z}_p^{s,j}, \mathbf{C}_p^s) \right] - \left[ \mathbf{Z}_p^{s,0} - \frac{1}{L_Z^s} \sum_{j=0}^{t-1} \nabla_Z f_p(\mathbf{Z}_p^{s,j}, \mathbf{C}_p^s) \right] \right\|_F^2 \\
&= \frac{1}{(L_Z^s)^2} \left\| \frac{1}{\bar{P}} \sum_{j=0}^{t-1} \sum_{p \in \mathcal{A}^s} \nabla_Z f_p(\mathbf{Z}_p^{s,j}, \mathbf{C}_p^s) - \sum_{j=0}^{t-1} \nabla_Z f_p(\mathbf{Z}_p^{s,j}, \mathbf{C}_p^s) \right\|_F^2 \\
&\leq \frac{t}{(L_Z^s)^2} \sum_{j=0}^{t-1} \left\| \frac{1}{\bar{P}} \sum_{p \in \mathcal{A}^s} \nabla_Z f_p(\mathbf{Z}_p^{s,j}, \mathbf{C}_p^s) - \nabla_Z f_p(\mathbf{Z}_p^{s,j}, \mathbf{C}_p^s) \right\|_F^2 \\
&= \frac{t}{(L_Z^s)^2 \bar{P}^2} \sum_{j=0}^{t-1} \left\| \sum_{p' \in \mathcal{A}^s} [\nabla_Z f_{p'}(\mathbf{Z}_{p'}^{s,j}, \mathbf{C}_{p'}^s) - \nabla_Z f_p(\mathbf{Z}_p^{s,j}, \mathbf{C}_p^s)] \right\|_F^2 \\
&\leq \frac{t}{(L_Z^s)^2 \bar{P}} \sum_{j=0}^{t-1} \sum_{p' \in \mathcal{A}^s} \left\| \nabla_Z f_{p'}(\mathbf{Z}_{p'}^{s,j}, \mathbf{C}_{p'}^s) - \nabla_Z f_p(\mathbf{Z}_p^{s,j}, \mathbf{C}_p^s) \right\|_F^2 \\
&= \frac{t}{(L_Z^s)^2 \bar{P}} \sum_{j=0}^{t-1} \sum_{p'=1}^P \omega_{p'} \underbrace{\left\| \nabla_Z f_{p'}(\mathbf{Z}_{p'}^{s,j}, \mathbf{C}_{p'}^s) - \nabla_Z f_p(\mathbf{Z}_p^{s,j}, \mathbf{C}_p^s) \right\|_F^2}_{T.6} \\
&\leq \frac{t}{(L_Z^s)^2 \bar{P}} \sum_{j=0}^{t-1} \sum_{p'=1}^P \omega_{p'} \underbrace{\left[ 4(L_{Z_p}^s)^2 \|\mathbf{Z}^{s,j} - \mathbf{Z}_p^{s,j}\|_F^2 + 4(L_{Z_{p'}}^s)^2 \|\mathbf{Z}^{s,j} - \mathbf{Z}_{p'}^{s,j}\|_F^2 + 8\zeta^2 \right]}_{\text{Bound of T.6}} \\
&= \frac{4t}{(L_Z^s)^2 \bar{P}} \sum_{j=0}^{t-1} (L_{Z_p}^s)^2 \|\mathbf{Z}^{s,j} - \mathbf{Z}_p^{s,j}\|_F^2 + \frac{8t^2 \zeta^2}{(L_Z^s)^2 m} \\
&\quad + \frac{4t}{(L_Z^s)^2 \bar{P}} \sum_{j=0}^{t-1} \sum_{p'=1}^P \omega_{p'} (L_{Z_{p'}}^s)^2 \|\mathbf{Z}^{s,j} - \mathbf{Z}_{p'}^{s,j}\|_F^2
\end{aligned} \tag{64}$$

□

Based on Lemma 3, we give a bound on T.7.

**Lemma E.4.** *For any  $s$ , it holds that*

$$\sum_{t=1}^Q \sum_{p=1}^P \omega_p (L_{Z_p}^s)^2 \|\mathbf{Z}^{s,t-1} - \mathbf{Z}_p^{s,t-1}\|_F^2 \leq \frac{8Q(Q-1)(2Q-1)\zeta^2}{\bar{P} - 4(Q-1)^2(1 + \bar{L}_Z^2/L_Z^2)} \tag{65}$$

*Proof.* Based on Lemma 3, we have

$$\begin{aligned}
& \sum_{t=1}^Q \sum_{p=1}^P \omega_p(L_{Z_p}^s)^2 \|\mathbf{Z}^{s,t-1} - \mathbf{Z}_p^{s,t-1}\|_F^2 \\
& \leq \sum_{t=1}^Q \sum_{p=1}^P \omega_p(L_{Z_p}^s)^2 \left[ \frac{4(t-1)}{(L_Z^s)^2 \bar{P}} \sum_{j=0}^{t-2} (L_{Z_p}^s)^2 \left\| \mathbf{Z}^{s,j} - \mathbf{Z}_p^{s,j} \right\|_F^2 + \frac{8(t-1)^2 \zeta^2}{(L_Z^s)^2 \bar{P}} \right. \\
& \quad \left. + \frac{4(t-1)}{(L_Z^s)^2 \bar{P}} \sum_{j=0}^{t-2} \sum_{p'=1}^P \omega_{p'}(L_{Z_{p'}}^s)^2 \left\| \mathbf{Z}^{s,j} - \mathbf{Z}_{p'}^{s,j} \right\|_F^2 \right] \\
& = \sum_{t=1}^Q \frac{4(t-1)}{\bar{P}} \sum_{p=1}^P \omega_p(L_{Z_p}^s)^2 \left( \frac{L_{Z_p}^s}{L_Z^s} \right)^2 \sum_{j=0}^{t-2} \left\| \mathbf{Z}^{s,j} - \mathbf{Z}_p^{s,j} \right\|_F^2 + \sum_{t=1}^Q \frac{8(t-1)^2 \zeta^2}{\bar{P}} \\
& \quad + \sum_{t=1}^Q \frac{4(t-1)}{\bar{P}} \sum_{j=0}^{t-2} \sum_{p'=1}^P \omega_{p'}(L_{Z_{p'}}^s)^2 \left\| \mathbf{Z}^{s,j} - \mathbf{Z}_{p'}^{s,j} \right\|_F^2 \tag{66} \\
& \leq \sum_{t=1}^Q \frac{4(t-1)}{\bar{P}} \sum_{p=1}^P \omega_p(L_{Z_p}^s)^2 \left( \frac{\bar{L}_Z}{L_Z} \right)^2 \sum_{j=0}^{t-2} \left\| \mathbf{Z}^{s,j} - \mathbf{Z}_p^{s,j} \right\|_F^2 + \sum_{t=1}^Q \frac{8(t-1)^2 \zeta^2}{\bar{P}} \\
& \quad + \sum_{t=1}^Q \frac{4(t-1)}{\bar{P}} \sum_{j=0}^{t-2} \sum_{p'=1}^P \omega_{p'}(L_{Z_{p'}}^s)^2 \left\| \mathbf{Z}^{s,j} - \mathbf{Z}_{p'}^{s,j} \right\|_F^2 \\
& = \sum_{t=1}^Q \frac{4(t-1)}{\bar{P}} \left( 1 + \frac{\bar{L}_Z^2}{L_Z^2} \right) \sum_{j=0}^{t-2} \sum_{p=1}^P \omega_p(L_{Z_p}^s)^2 \left\| \mathbf{Z}^{s,j} - \mathbf{Z}_p^{s,j} \right\|_F^2 + \frac{8Q(Q-1)(2Q-1)\zeta^2}{\bar{P}} \\
& \leq \frac{4(Q-1)^2}{\bar{P}} \left( 1 + \frac{\bar{L}_Z^2}{L_Z^2} \right) \sum_{t=1}^Q \sum_{p=1}^P \omega_p(L_{Z_p}^s)^2 \left\| \mathbf{Z}^{s,t-1} - \mathbf{Z}_p^{s,t-1} \right\|_F^2 + \frac{8Q(Q-1)(2Q-1)\zeta^2}{\bar{P}}
\end{aligned}$$

Here, we used the inequality as

$$\sum_{t=1}^Q \frac{4(t-1)}{\bar{P}} \sum_{j=0}^{t-2} \left\| \mathbf{Z}^{s,j} - \mathbf{Z}_p^{s,j} \right\|_F^2 \leq \frac{4(Q-1)^2}{\bar{P}} \sum_{t=1}^Q \left\| \mathbf{Z}^{s,t-1} - \mathbf{Z}_p^{s,t-1} \right\|_F^2 \tag{67}$$

Rearranging the above inequality, we have

$$\sum_{t=1}^Q \sum_{p=1}^P \omega_p(L_{Z_p}^s)^2 \left\| \mathbf{Z}^{s,t-1} - \mathbf{Z}_p^{s,t-1} \right\|_F^2 \leq \frac{8Q(Q-1)(2Q-1)\zeta^2}{\bar{P} - 4(Q-1)^2(1 + \bar{L}_Z^2/L_Z^2)} \tag{68}$$

□

Based on Lemma 4, we continue to do the following derivation:

$$\begin{aligned}
& F(\mathbf{Z}^{s,Q}, \mathbf{C}^s) - F(\mathbf{Z}^{s,0}, \mathbf{C}^s) \\
& \leq -\frac{L_Z^s}{4} \sum_{t=1}^Q \|\mathbf{Z}^{s,t} - \mathbf{Z}^{s,t-1}\|_F^2 + \frac{16Q\zeta^2}{\bar{P}L_Z^s} \\
& \quad + \underbrace{\frac{2}{L_Z^s} \left(1 + \frac{8}{\bar{P}}\right) \sum_{t=1}^Q \sum_{p=1}^P \omega_p (L_{Z_p}^s)^2 \|\mathbf{Z}^{s,t-1} - \mathbf{Z}_p^{s,t-1}\|_F^2}_{T.7} \\
& \leq -\frac{L_Z^s}{4} \sum_{t=1}^Q \|\mathbf{Z}^{s,t} - \mathbf{Z}^{s,t-1}\|_F^2 + \frac{16Q\zeta^2}{\bar{P}L_Z^s} \\
& \quad + \underbrace{\frac{2}{L_Z^s} \left(1 + \frac{8}{\bar{P}}\right) \frac{8Q(Q-1)(2Q-1)\zeta^2}{\bar{P} - 4(Q-1)^2(1 + \bar{L}_Z^2/L_Z^2)}}_{\text{Bound of T.7}} \\
& \leq -\frac{L_Z^s}{4} \sum_{t=1}^Q \|\mathbf{Z}^{s,t} - \mathbf{Z}^{s,t-1}\|_F^2 + \frac{16Q\zeta^2\psi}{\bar{P}L_Z^s}
\end{aligned} \tag{69}$$

where  $\psi = 1 + \frac{(\bar{P}+8)(Q-1)(2Q-1)}{\bar{P}-4(Q-1)^2(1+\bar{L}_Z^2/L_Z^2)}$ .

Then, combining (46) and (69) yields

$$\begin{aligned}
& \frac{1}{2}(\gamma_{min}^s + \lambda) \sum_{p=1}^P \omega_p \|\mathbf{C}_p^s - \mathbf{C}_p^{s-1}\|_F^2 + \frac{L_Z^s}{4} \sum_{t=1}^Q \|\mathbf{Z}^{s,t} - \mathbf{Z}^{s,t-1}\|_F^2 \\
& \leq [F(\mathbf{Z}^{s,0}, \mathbf{C}^{s-1}) - F(\mathbf{Z}^{s,Q}, \mathbf{C}^s)] + \frac{16Q\zeta^2\psi}{\bar{P}L_Z^s}
\end{aligned} \tag{70}$$

**Derivation of the main result:** Based on (70), we derive the convergence in terms of the iterative terms,  $T_C(\mathbf{Z}^{s,0}, \mathbf{C}^s)$  and  $T_Z(\mathbf{Z}^{s,t}, \mathbf{C}^s)$  for  $s = 1, 2, \dots, S$ .

$$T_C(\mathbf{Z}^{s,0}, \mathbf{C}^s) \leq \frac{2}{\gamma_{min}^s + \lambda} [F(\mathbf{Z}^{s,0}, \mathbf{C}^{s-1}) - F(\mathbf{Z}^{s,Q}, \mathbf{C}^s)] + \frac{32Q\zeta^2\psi}{(\gamma_{min}^s + \lambda)\bar{P}L_Z^s} \tag{71}$$

Then, summing it up from  $s = 1$  to  $S$  yields

$$\sum_{s=1}^S T_C(\mathbf{Z}^{s,0}, \mathbf{C}^s) \leq \frac{2}{\underline{\gamma}_{min} + \lambda} [F(\mathbf{Z}^{1,0}, \mathbf{C}^0) - \underline{f}] + \frac{32SQ\zeta^2\psi}{(\underline{\gamma}_{min} + \lambda)\bar{P}L_Z} \tag{72}$$

Similarly, we have

$$\sum_{t=1}^Q T_Z(\mathbf{Z}^{s,t}, \mathbf{C}^s) \leq \frac{4}{L_Z^s} [F(\mathbf{Z}^{s,0}, \mathbf{C}^{s-1}) - F(\mathbf{Z}^{s,Q}, \mathbf{C}^s)] + \frac{64Q\zeta^2\psi}{\bar{P}(L_Z^s)^2} \tag{73}$$

Then, summing it up from  $s = 1$  to  $S$  yields

$$\sum_{s=1}^S \sum_{t=1}^Q T_Z(\mathbf{Z}^{s,t}, \mathbf{C}^s) \leq \frac{4}{\underline{L}_Z} [F(\mathbf{Z}^{1,0}, \mathbf{C}^0) - \underline{f}] + \frac{64SQ\zeta^2\psi}{\bar{P}\underline{L}_Z^2} \tag{74}$$

Lastly, combining (72) and (74) and dividing two sides of it by  $T = S(1 + Q)$  yields

$$\frac{1}{T} \left[ \sum_{s=1}^S T_C(\mathbf{Z}^{s,0}, \mathbf{C}^s) + \sum_{s=1}^S \sum_{t=1}^Q T_Z(\mathbf{Z}^{s,t}, \mathbf{C}^s) \right] \leq \frac{D}{T} [F(\mathbf{Z}^{1,0}, \mathbf{C}^0) - \underline{f}] + \frac{16\zeta^2\psi D}{\bar{P}\underline{L}_Z} \tag{75}$$

where  $D = \frac{2}{\underline{\gamma}_{min} + \lambda} + \frac{4}{\underline{L}_Z}$ .

## F Proof for theorem on error bound on noisy similarity matrix

*Proof.* It follows from the triangle inequality of matrix norm that

$$\begin{aligned} \left\| \hat{\mathbf{K}}_{\tilde{\mathbf{x}}\tilde{\mathbf{x}}} - \mathbf{K}_{xx} \right\|_{\infty} &= \left\| \hat{\mathbf{K}}_{\tilde{\mathbf{x}}\tilde{\mathbf{x}}} - \mathbf{K}_{\tilde{\mathbf{x}}\tilde{\mathbf{x}}} + \mathbf{K}_{\tilde{\mathbf{x}}\tilde{\mathbf{x}}} - \mathbf{K}_{xx} \right\|_{\infty} \\ &\leq \left\| \hat{\mathbf{K}}_{\tilde{\mathbf{x}}\tilde{\mathbf{x}}} - \mathbf{K}_{\tilde{\mathbf{x}}\tilde{\mathbf{x}}} \right\|_{\infty} + \left\| \mathbf{K}_{\tilde{\mathbf{x}}\tilde{\mathbf{x}}} - \mathbf{K}_{xx} \right\|_{\infty}. \end{aligned} \quad (76)$$

Since  $\tilde{\mathbf{X}} = \mathbf{X} + \mathbf{E}$  and  $\phi(\tilde{\mathbf{X}}) = \phi(\mathbf{Z})\mathbf{C}$ , we have

$$\begin{cases} \hat{\mathbf{K}}_{\tilde{\mathbf{x}}\tilde{\mathbf{x}}} = (\phi(\mathbf{Z})\mathbf{C})^T (\phi(\mathbf{Z})\mathbf{C}) = \mathbf{C}^T \mathcal{K}(\mathbf{Z}, \mathbf{Z}) \mathbf{C} \\ \mathbf{K}_{\tilde{\mathbf{x}}\tilde{\mathbf{x}}} = \phi(\tilde{\mathbf{X}})^T \phi(\tilde{\mathbf{X}}) = \mathcal{K}(\tilde{\mathbf{X}}, \tilde{\mathbf{X}}) \\ \mathbf{K}_{xx} = \phi(\mathbf{X})^T \phi(\mathbf{X}) = \mathcal{K}(\mathbf{X}, \mathbf{X}) \end{cases} \quad (77)$$

where  $\mathbf{E}$  is a Gaussian noise matrix, of which  $\mathbf{E}_{i,j} \sim \mathcal{N}(0, \sigma^2)$ . Hence, we obtain

$$\left\| \hat{\mathbf{K}}_{\tilde{\mathbf{x}}\tilde{\mathbf{x}}} - \mathbf{K}_{xx} \right\|_{\infty} \leq \left\| \mathbf{C}^T \mathcal{K}(\mathbf{Z}, \mathbf{Z}) \mathbf{C} - \mathcal{K}(\tilde{\mathbf{X}}, \tilde{\mathbf{X}}) \right\|_{\infty} + \left\| \mathcal{K}(\tilde{\mathbf{X}}, \tilde{\mathbf{X}}) - \mathcal{K}(\mathbf{X}, \mathbf{X}) \right\|_{\infty}. \quad (78)$$

Denote by  $\mathbf{x}_i$  (or  $\mathbf{e}_i$ ) the  $i$ -th column of  $\mathbf{X} \in \mathbb{R}^{m \times n}$  (or  $\mathbf{E} \in \mathbb{R}^{m \times n}$ ),  $i = 1, \dots, n$ , we first derive the upper bound of the second term of the RHS of (78) as follows:

$$\begin{aligned} \left\| \mathbf{K}_{\tilde{\mathbf{x}}\tilde{\mathbf{x}}} - \mathbf{K}_{xx} \right\|_{\infty} &= \left\| \mathcal{K}(\tilde{\mathbf{X}}, \tilde{\mathbf{X}}) - \mathcal{K}(\mathbf{X}, \mathbf{X}) \right\|_{\infty} \\ &= \max_{i,j} |\mathcal{K}(\tilde{\mathbf{x}}_i, \tilde{\mathbf{x}}_j) - \mathcal{K}(\mathbf{x}_i, \mathbf{x}_j)| \\ &= \max_{i,j} \left| \exp\left(-\frac{\|(\mathbf{x}_i + \mathbf{e}_i) - (\mathbf{x}_j + \mathbf{e}_j)\|_2^2}{2r^2}\right) - \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|_2^2}{2r^2}\right) \right| \\ &\leq \max_{i,j} \frac{1}{2r^2} \left| -\|(\mathbf{x}_i + \mathbf{e}_i) - (\mathbf{x}_j + \mathbf{e}_j)\|_2^2 + \|\mathbf{x}_i - \mathbf{x}_j\|_2^2 \right| \\ &= \max_{i,j} \frac{1}{2r^2} \left| \|(\mathbf{x}_i - \mathbf{x}_j) + (\mathbf{e}_i - \mathbf{e}_j)\|_2^2 - \|\mathbf{x}_i - \mathbf{x}_j\|_2^2 \right| \\ &= \max_{i,j} \frac{1}{2r^2} \left| \|\mathbf{e}_i - \mathbf{e}_j\|_2^2 + 2\langle \mathbf{x}_i - \mathbf{x}_j, \mathbf{e}_i - \mathbf{e}_j \rangle \right| \\ &\leq \max_{i,j} \frac{1}{2r^2} \left( \|\mathbf{e}_i - \mathbf{e}_j\|_2^2 + 2|\langle \mathbf{x}_i - \mathbf{x}_j, \mathbf{e}_i - \mathbf{e}_j \rangle| \right) \\ &\leq \max_{i,j} \frac{1}{2r^2} \left( \|\mathbf{e}_i - \mathbf{e}_j\|_2^2 + 2\|\mathbf{x}_i - \mathbf{x}_j\|_2 \|\mathbf{e}_i - \mathbf{e}_j\|_2 \right), \end{aligned} \quad (79)$$

where for the first inequality we have used the fact that the exponential function is locally Lipschitz continuous, i.e.,

$$|e^x - e^y| < |x - y| \text{ for } x, y < 0.$$

Now, let us figure out the upper bound of  $\|\mathbf{e}_i - \mathbf{e}_j\|_2^2$ . Note that

$$\|\mathbf{e}_i - \mathbf{e}_j\|_2^2 = \sum_{l=1}^m (e_{li} - e_{lj})^2 = 2\sigma^2 \sum_{l=1}^m \left( \frac{e_{li} - e_{lj}}{\sqrt{2}\sigma} \right)^2 \quad (80)$$

where  $e_{li}$  represents the  $l$ -th element of the column vector  $\mathbf{e}_i$ ,  $k = 1, \dots, m$ . It is clear that for  $l = 1, \dots, m$ ,

$$\begin{aligned} \mathbb{E}[e_{li} - e_{lj}] &= 0, \\ \text{var}[e_{li} - e_{lj}] &= 2\sigma^2. \end{aligned} \quad (81)$$

Hence,  $\frac{e_{li} - e_{lj}}{\sqrt{2}\sigma}$  is a standard Gaussian random variable drawn from  $\mathcal{N}(0, 1)$ . Based on this, we can define a random variable as

$$Q = \sum_{l=1}^m \left( \frac{e_{li} - e_{lj}}{\sqrt{2}\sigma} \right)^2, \quad (82)$$

which is distributed according to the Chi-squared distribution with  $m$  degrees of freedom.

From Laurent and Massart [2000], we know that for any positive  $t$ , the Chi-squared variable  $Q$  with  $m$  degrees of freedom satisfies

$$\mathbb{P}(Q > m + 2\sqrt{mt} + 2t) \leq 1 - e^{-t}. \quad (83)$$

Hence, a bound on  $\|e_i - e_j\|_2^2$  with probability  $1 - e^{-t}$  is

$$\|e_i - e_j\|_2^2 = 2\sigma^2 Q \leq 2\sigma^2(m + 2\sqrt{mt} + 2t). \quad (84)$$

Using union bound for (84), we have

$$\max_{i,j} \|e_i - e_j\|_2^2 = 2\sigma^2 Q \leq 2\sigma^2(m + 2\sqrt{mt} + 2t). \quad (85)$$

holds with probability at least  $1 - n(n-1)e^{-t}$ . Assume  $\|\mathbf{x}_i\|_2 \leq \theta$ , then  $\|\mathbf{x}_i - \mathbf{x}_j\|_2 \leq 2\theta$ . For convenience, let  $\xi = \sqrt{m + 2\sqrt{mt} + 2t}$ . It follows from (79) and (85) that, with probability at least  $1 - n(n-1)e^{-t}$ ,

$$\begin{aligned} \|\mathbf{K}_{\tilde{x}\tilde{x}} - \mathbf{K}_{xx}\|_\infty &\leq \frac{1}{2r^2} \left[ 2\sigma^2 \xi^2 + 2\|\mathbf{x}_i - \mathbf{x}_j\|_2 \sqrt{2\sigma\xi} \right] \\ &\leq \frac{1}{r^2} \left[ \sigma^2 \xi^2 + 2\sqrt{2\sigma\xi}\theta \right] \\ &= \frac{1}{r^2} \left[ (\sigma\xi + \sqrt{2}\theta)^2 - 2\theta^2 \right]. \end{aligned} \quad (86)$$

Now, we figure out the upper bound of the first term of the RHS of (78). We have

$$\begin{aligned} \|\hat{\mathbf{K}}_{\tilde{x}\tilde{x}} - \mathbf{K}_{\tilde{x}\tilde{x}}\|_\infty &= \left\| \mathbf{C}^T \mathcal{K}(\mathbf{Z}, \mathbf{Z}) \mathbf{C} - \mathcal{K}(\tilde{\mathbf{X}}, \tilde{\mathbf{X}}) \right\|_\infty \\ &= \left\| (\phi(\mathbf{Z})\mathbf{C})^T (\phi(\mathbf{Z})\mathbf{C}) - \phi(\tilde{\mathbf{X}})^T \phi(\tilde{\mathbf{X}}) \right\|_\infty \\ &= \left\| (\phi(\mathbf{Z})\mathbf{C})^T (\phi(\mathbf{Z})\mathbf{C}) - (\phi(\mathbf{Z})\mathbf{C})^T \phi(\tilde{\mathbf{X}}) + (\phi(\mathbf{Z})\mathbf{C})^T \phi(\tilde{\mathbf{X}}) - \phi(\tilde{\mathbf{X}})^T \phi(\tilde{\mathbf{X}}) \right\|_\infty \\ &= \left\| (\phi(\mathbf{Z})\mathbf{C})^T (\phi(\mathbf{Z})\mathbf{C} - \phi(\tilde{\mathbf{X}})) + (\phi(\mathbf{Z})\mathbf{C} - \phi(\tilde{\mathbf{X}}))^T \phi(\tilde{\mathbf{X}}) \right\|_\infty \\ &\leq \left\| (\phi(\mathbf{Z})\mathbf{C})^T (\phi(\mathbf{Z})\mathbf{C} - \phi(\tilde{\mathbf{X}})) \right\|_\infty + \left\| (\phi(\mathbf{Z})\mathbf{C} - \phi(\tilde{\mathbf{X}}))^T \phi(\tilde{\mathbf{X}}) \right\|_\infty \\ &\leq \|\phi(\mathbf{Z})\mathbf{C}\|_{2,\infty} \cdot \left\| \phi(\mathbf{Z})\mathbf{C} - \phi(\tilde{\mathbf{X}}) \right\|_{2,\infty} + \left\| \phi(\mathbf{Z})\mathbf{C} - \phi(\tilde{\mathbf{X}}) \right\|_{2,\infty} \cdot \left\| \phi(\tilde{\mathbf{X}}) \right\|_{2,\infty} \\ &= (\|\phi(\mathbf{Z})\mathbf{C}\|_{2,\infty} + \left\| \phi(\tilde{\mathbf{X}}) \right\|_{2,\infty}) \left\| \phi(\mathbf{Z})\mathbf{C} - \phi(\tilde{\mathbf{X}}) \right\|_{2,\infty} \end{aligned} \quad (87)$$

where  $\|\cdot\|_{2,\infty}$  is a norm such that  $\|\mathbf{X}\|_{2,\infty} = \max_i \|\mathbf{X}_{:,i}\|_2$  for a real matrix  $\mathbf{X}$ . Here we can just assume that  $\left\| \phi(\mathbf{Z})\mathbf{C} - \phi(\tilde{\mathbf{X}}) \right\|_{2,\infty} \leq \gamma$ ,  $\gamma$  is some constant; this relies on the optimization.

Moreover, assume  $\|\mathbf{C}\|_2 \leq \tau_C$ , we have

$$\begin{aligned} \left\| \phi(\tilde{\mathbf{X}}) \right\|_{2,\infty} &= 1 \\ \left\| \phi(\mathbf{Z})\mathbf{C} \right\|_{2,\infty} &\leq \|\phi(\mathbf{Z})\|_F \max_j \|\mathbf{C}_{:,j}\| \leq \sqrt{d}\tau_C \end{aligned} \quad (88)$$

Hence, we can continue to do the derivation of the preceding inequality.

$$\begin{aligned} \|\hat{\mathbf{K}}_{\tilde{x}\tilde{x}} - \mathbf{K}_{\tilde{x}\tilde{x}}\|_\infty &\leq \left( \|\phi(\mathbf{Z})\mathbf{C}\|_{2,\infty} + \left\| \phi(\tilde{\mathbf{X}}) \right\|_{2,\infty} \right) \left\| \phi(\mathbf{Z})\mathbf{C} - \phi(\tilde{\mathbf{X}}) \right\|_{2,\infty} \\ &\leq (\sqrt{d}\tau_C + 1)\gamma \end{aligned} \quad (89)$$

As a result, the overall bound is

$$\begin{aligned} \left\| \hat{\mathbf{K}}_{\tilde{x}\tilde{x}} - \mathbf{K}_{xx} \right\|_{\infty} &\leq \left\| \hat{\mathbf{K}}_{\tilde{x}\tilde{x}} - \mathbf{K}_{\tilde{x}\tilde{x}} \right\|_{\infty} + \left\| \mathbf{K}_{\tilde{x}\tilde{x}} - \mathbf{K}_{xx} \right\|_{\infty} \\ &\leq \frac{1}{r^2} \left[ (\sigma\xi + \sqrt{2}\theta)^2 - 2\theta^2 \right] + (\sqrt{d}\tau_C + 1)\gamma \end{aligned} \quad (90)$$

where  $\xi = \sqrt{(m + 2\sqrt{mt} + 2t)}$ .

□

## G Proof for Theorem 3.6

*Proof.* For a data matrix  $\mathbf{X} \in \mathbb{R}^{m \times n}$ , it is perturbed by a Gaussian noise matrix  $\mathbf{E} \in \mathbb{R}^{m \times n}$  with  $\mathcal{N}(0, \sigma^2)$  to form the noisy data matrix  $\tilde{\mathbf{X}} = \mathbf{X} + \mathbf{E}$ . Let  $\mathbf{K}_{xx} = \mathcal{K}(\mathbf{X}, \mathbf{X})$  be the ground truth and  $\hat{\mathbf{K}}_{\tilde{x}\tilde{x}} = \mathbf{C}^T \mathcal{K}(\tilde{\mathbf{Z}}, \tilde{\mathbf{Z}}) \mathbf{C}$  be the approximated similarity matrix by FedKMF algorithm 1. Then, consider any two data points in  $\mathbf{X}$ ,  $\mathbf{x}_i$  and  $\mathbf{x}_j$ , with the identical label, we know that

$$\mathbf{K}_{iu} \leq \mathbf{K}_{ij} \quad (91)$$

where  $\mathbf{K}_{iu} = \max_k ((\mathbf{K}_{xx})_{ik}^{inter})$  and  $\mathbf{K}_{ij} = (\mathbf{K}_{xx})_{ij}$ .

After running FedKMF (Algorithm 1) on the noisy matrix  $\tilde{\mathbf{X}}$ , if the approximated similarity matrix satisfies  $\left\| \hat{\mathbf{K}}_{\tilde{x}\tilde{x}} - \mathbf{K}_{xx} \right\| < B(\sigma)$ , then consider two points in  $\tilde{\mathbf{X}}$ ,  $\tilde{\mathbf{x}}_i$  and  $\tilde{\mathbf{x}}_j$ , that are actually  $\mathbf{x}_i$  and  $\mathbf{x}_j$  perturbed by some noise, we have

$$\begin{aligned} \hat{\mathbf{K}}_{iu} - \hat{\mathbf{K}}_{ij} &= \hat{\mathbf{K}}_{iu} - \mathbf{K}_{iu} + \mathbf{K}_{iu} - \hat{\mathbf{K}}_{ij} + \mathbf{K}_{ij} - \mathbf{K}_{ij} \\ &= (\hat{\mathbf{K}}_{iu} - \mathbf{K}_{iu}) + (\mathbf{K}_{ij} - \hat{\mathbf{K}}_{ij}) + (\mathbf{K}_{iu} - \mathbf{K}_{ij}) \\ &\leq |\hat{\mathbf{K}}_{iu} - \mathbf{K}_{iu}| + |\mathbf{K}_{ij} - \hat{\mathbf{K}}_{ij}| + (\mathbf{K}_{iu} - \mathbf{K}_{ij}) \\ &\leq B(\sigma) + B(\sigma) + (\mathbf{K}_{iu} - \mathbf{K}_{ij}) \\ &= 2B(\sigma) + (\mathbf{K}_{iu} - \mathbf{K}_{ij}) \end{aligned} \quad (92)$$

where  $\hat{\mathbf{K}}_{iu} = \max_k ((\hat{\mathbf{K}}_{\tilde{x}\tilde{x}})_{ik}^{inter})$  and  $\hat{\mathbf{K}}_{ij} = (\hat{\mathbf{K}}_{\tilde{x}\tilde{x}})_{ij}$ .

Based on Definition 3.5,  $\tilde{\mathbf{x}}_i$  and  $\tilde{\mathbf{x}}_j$  can be correctly clustered only if the following inequality holds with some tolerance  $\epsilon > 0$ .

$$\hat{\mathbf{K}}_{iu} - \hat{\mathbf{K}}_{ij} \leq \epsilon \quad (93)$$

Thus, combining inequalities (92) and (93), the bound function  $B(\sigma)$  satisfies

$$B(\sigma) \leq \frac{1}{2} [\epsilon - (\mathbf{K}_{iu} - \mathbf{K}_{ij})] \quad (94)$$

Similarly, if we consider any two data points in  $\mathbf{X}$ ,  $\mathbf{x}_i$  and  $\mathbf{x}_j$ , with different labels, we know that

$$\mathbf{K}_{ij} \leq \mathbf{K}_{iv} \quad (95)$$

where  $\mathbf{K}_{iv} = \min_k ((\mathbf{K}_{xx})_{ik}^{intra})$ .

After running FedKMF algorithm 1 on the noisy matrix  $\tilde{\mathbf{X}}$ , we have

$$\begin{aligned} \hat{\mathbf{K}}_{ij} - \hat{\mathbf{K}}_{iv} &= \hat{\mathbf{K}}_{ij} - \mathbf{K}_{ij} + \mathbf{K}_{ij} - \hat{\mathbf{K}}_{iv} + \mathbf{K}_{iv} - \mathbf{K}_{iv} \\ &= (\hat{\mathbf{K}}_{ij} - \mathbf{K}_{ij}) + (\mathbf{K}_{iv} - \hat{\mathbf{K}}_{iv}) + (\mathbf{K}_{ij} - \mathbf{K}_{iv}) \\ &\leq |\hat{\mathbf{K}}_{ij} - \mathbf{K}_{ij}| + |\mathbf{K}_{iv} - \hat{\mathbf{K}}_{iv}| + (\mathbf{K}_{ij} - \mathbf{K}_{iv}) \\ &\leq B(\sigma) + B(\sigma) + (\mathbf{K}_{ij} - \mathbf{K}_{iv}) \\ &= 2B(\sigma) + (\mathbf{K}_{ij} - \mathbf{K}_{iv}) \end{aligned} \quad (96)$$

where  $\hat{\mathbf{K}}_{iv} = \max_k ((\hat{\mathbf{K}}_{\tilde{x}\tilde{x}})_{ik}^{intra})$ .

Based on Definition 3.5,  $\tilde{\mathbf{x}}_i$  and  $\tilde{\mathbf{x}}_j$  can be correctly clustered only if the following inequality holds with some tolerance  $\epsilon > 0$ .

$$\hat{\mathbf{K}}_{ij} - \hat{\mathbf{K}}_{iv} \leq \epsilon \quad (97)$$

Thus, combining inequalities (96) and (97), the bound function  $B(\sigma)$  satisfies

$$B(\sigma) \leq \frac{1}{2}[\epsilon - (\mathbf{K}_{ij} - \mathbf{K}_{iv})] \quad (98)$$

Then, with two upper bounds on  $B(\sigma)$ , (94) and (98), we have

$$B(\sigma) \leq \frac{1}{2} \min_i \{\epsilon - (\mathbf{K}_{iu} - \mathbf{K}_{ij}), \epsilon - (\mathbf{K}_{ij} - \mathbf{K}_{iv})\}. \quad (99)$$

where  $\epsilon$  is the parameter of tolerance.

Alternatively, a slightly looser version is like

$$\begin{aligned} B(\sigma) &\leq \min_i \frac{1}{4} [2\epsilon - (\mathbf{K}_{iu} - \mathbf{K}_{iv})] \\ \text{or } B(\sigma) &\leq \frac{\epsilon}{2} - \max_i \frac{1}{4} (\mathbf{K}_{iu} - \mathbf{K}_{iv}). \end{aligned} \quad (100)$$

where  $\mathbf{K}_{iu} = \max_k ((\mathbf{K}_{xx})_{ik}^{inter})$  and  $\mathbf{K}_{iv} = \min_k ((\mathbf{K}_{xx})_{ik}^{intra})$ .

□

## H Proof for Proposition 3.7

*Proof.* The  $\ell_2$ -sensitivity [Dwork *et al.*, 2014] of a function  $f : \mathbb{N}^{|\mathcal{X}|} \rightarrow \mathbb{R}^k$  is:

$$\Delta_2(f) = \max_{x \sim y} \|f(x) - f(y)\|_2,$$

where  $x \sim y$  denotes that  $x$  and  $y$  are neighboring datasets. In our case, the function is  $f(x) = x$ . Then

$$\|f(x) - f(y)\|_2 = \|x - y\|_2 \leq 2\tau_X.$$

It means  $\Delta_2(f) \leq 2\tau_X$ . Now using Theorem 3.22 in [Dwork *et al.*, 2014] and Lemma H.1, we get the desired result.

**Lemma H.1** (Post-Processing [Dwork *et al.*, 2014]). *Let  $\mathcal{M} : \mathbb{N}^{|\mathcal{X}|} \rightarrow R$  be a randomized algorithm that is  $(\epsilon, \delta)$ -differentially private. Let  $h : R \rightarrow R'$  be an arbitrary randomized mapping. Then  $h \circ \mathcal{M} : \mathbb{N}^{|\mathcal{X}|} \rightarrow R'$  is  $(\epsilon, \delta)$ -differentially private.*

□

## I Proof for Theorem 3.8

*Proof.* Based on the assumptions  $\|\mathbf{C}\|_{2,\infty} \leq \tau_C$ ,  $\|\phi(\mathbf{Z})\mathbf{C} - \phi(\mathbf{X})\|_{2,\infty} \leq \gamma$ ,  $\tilde{\mathbf{C}} = \mathbf{C} + \mathbf{E}_C$  for the entry  $(\mathbf{E}_C)_{ij} \sim \mathcal{N}(0, \sigma_C^2)$ , and  $\tilde{\mathbf{Z}} = \mathbf{Z} + \mathbf{E}_Z$  for the entry  $(\mathbf{E}_Z)_{ij} \sim \mathcal{N}(0, \sigma_Z^2)$ , we obtain

$$\begin{aligned} \left\| \hat{\mathbf{K}}_{xx} - \mathbf{K}_{xx} \right\|_{\infty} &= \left\| \tilde{\mathbf{C}}^T \mathcal{K}(\tilde{\mathbf{Z}}, \tilde{\mathbf{Z}}) \tilde{\mathbf{C}} - \mathcal{K}(\mathbf{X}, \mathbf{X}) \right\|_{\infty} \\ &= \left\| (\phi(\tilde{\mathbf{Z}}) \tilde{\mathbf{C}})^T (\phi(\tilde{\mathbf{Z}}) \tilde{\mathbf{C}}) - \phi^T(\mathbf{X}) \phi(\mathbf{X}) \right\|_{\infty} \\ &= \left\| (\phi(\tilde{\mathbf{Z}}) \tilde{\mathbf{C}})^T (\phi(\tilde{\mathbf{Z}}) \tilde{\mathbf{C}}) - (\phi(\tilde{\mathbf{Z}}) \tilde{\mathbf{C}})^T \phi(\mathbf{X}) + (\phi(\tilde{\mathbf{Z}}) \tilde{\mathbf{C}})^T \phi(\mathbf{X}) - \phi^T(\mathbf{X}) \phi(\mathbf{X}) \right\|_{\infty} \\ &= \left\| (\phi(\tilde{\mathbf{Z}}) \tilde{\mathbf{C}})^T (\phi(\tilde{\mathbf{Z}}) \tilde{\mathbf{C}}) - \phi(\mathbf{X}) + (\phi(\tilde{\mathbf{Z}}) \tilde{\mathbf{C}} - \phi(\mathbf{X}))^T \phi(\mathbf{X}) \right\|_{\infty} \\ &\leq \left\| \phi(\tilde{\mathbf{Z}}) \tilde{\mathbf{C}} \right\|_{2,\infty} \left\| \phi(\tilde{\mathbf{Z}}) \tilde{\mathbf{C}} - \phi(\mathbf{X}) \right\|_{2,\infty} + \left\| \phi(\tilde{\mathbf{Z}}) \tilde{\mathbf{C}} - \phi(\mathbf{X}) \right\|_{2,\infty} \|\phi(\mathbf{X})\|_{2,\infty} \\ &\leq \left( \left\| \phi(\tilde{\mathbf{Z}}) \tilde{\mathbf{C}} \right\|_{2,\infty} + \|\phi(\mathbf{X})\|_{2,\infty} \right) \underbrace{\left\| \phi(\tilde{\mathbf{Z}}) \tilde{\mathbf{C}} - \phi(\mathbf{X}) \right\|_{2,\infty}}_{\text{T.1}} \end{aligned} \quad (101)$$

where  $\|\phi(\mathbf{X})\|_{2,\infty} = 1$ . An upper bound on  $\|\phi(\tilde{\mathbf{Z}})\tilde{\mathbf{C}}\|_{2,\infty}$  can be obtained only if we derive the upper bound on  $\|\phi(\tilde{\mathbf{Z}})\tilde{\mathbf{C}} - \phi(\mathbf{X})\|_{2,\infty}$ . That is, if  $\|\phi(\tilde{\mathbf{Z}})\tilde{\mathbf{C}} - \phi(\mathbf{X})\|_{2,\infty} \leq \gamma_{zc}$ , it implies that

$$\|\phi(\tilde{\mathbf{Z}})\tilde{\mathbf{C}}\|_{2,\infty} \leq \|\phi(\tilde{\mathbf{Z}})\tilde{\mathbf{C}} - \phi(\mathbf{X})\|_{2,\infty} + \|\phi(\mathbf{X})\|_{2,\infty} = \gamma_{zc} + 1 \quad (102)$$

which consequently gives

$$\|\tilde{\mathbf{K}}_{xx} - \mathbf{K}_{xx}\|_{\infty} \leq \gamma_{zc}(\gamma_{zc} + 2) \quad (103)$$

Hence, we derive the upper bound on  $\|\phi(\tilde{\mathbf{Z}})\tilde{\mathbf{C}} - \phi(\mathbf{X})\|_{2,\infty}$  for the remaining proof.

$$\begin{aligned} \|\phi(\tilde{\mathbf{Z}})\tilde{\mathbf{C}} - \phi(\mathbf{X})\|_{2,\infty} &= \|\phi(\tilde{\mathbf{Z}})\tilde{\mathbf{C}} - \phi(\mathbf{Z})\mathbf{C} + \phi(\mathbf{Z})\mathbf{C} - \phi(\mathbf{X})\|_{2,\infty} \\ &\leq \underbrace{\|\phi(\tilde{\mathbf{Z}})\tilde{\mathbf{C}} - \phi(\mathbf{Z})\mathbf{C}\|_{2,\infty}}_{T2} + \|\phi(\mathbf{Z})\mathbf{C} - \phi(\mathbf{X})\|_{2,\infty} \end{aligned} \quad (104)$$

where the second term  $\|\phi(\mathbf{Z})\mathbf{C} - \phi(\mathbf{X})\|_{2,\infty} \leq \gamma$  is the assumption. Next, we derive the upper bound on  $T.2$ .

$$\begin{aligned} \|\phi(\tilde{\mathbf{Z}})\tilde{\mathbf{C}} - \phi(\mathbf{Z})\mathbf{C}\|_{2,\infty} &= \|\phi(\tilde{\mathbf{Z}})\tilde{\mathbf{C}} - \phi(\tilde{\mathbf{Z}})\mathbf{C} + \phi(\tilde{\mathbf{Z}})\mathbf{C} - \phi(\mathbf{Z})\mathbf{C}\|_{2,\infty} \\ &= \|\phi(\tilde{\mathbf{Z}})(\tilde{\mathbf{C}} - \mathbf{C}) + (\phi(\tilde{\mathbf{Z}}) - \phi(\mathbf{Z}))\mathbf{C}\|_{2,\infty} \\ &\leq \|\phi(\tilde{\mathbf{Z}})(\tilde{\mathbf{C}} - \mathbf{C})\|_{2,\infty} + \|(\phi(\tilde{\mathbf{Z}}) - \phi(\mathbf{Z}))\mathbf{C}\|_{2,\infty} \\ &\leq \sqrt{d}\|\mathbf{E}_C\|_{2,\infty} + \|\phi(\tilde{\mathbf{Z}}) - \phi(\mathbf{Z})\|_F \|\mathbf{C}\|_{2,\infty} \\ &\leq \sigma_C \xi_d \sqrt{d} + \tau_C \underbrace{\|\phi(\tilde{\mathbf{Z}}) - \phi(\mathbf{Z})\|_F}_{T3} \end{aligned} \quad (105)$$

where we used  $\frac{1}{\sigma_C^2} \|\mathbf{E}_C\|_{2,\infty}^2 \leq \xi_d^2 = d + 2\sqrt{dt} + 2t$  with probability at least  $1 - ne^{-t}$  [Laurent and Massart, 2000] since it is the fact that  $\frac{1}{\sigma_C^2} \|\mathbf{E}_C\|_{2,\infty}^2 \sim \chi_d^2$  where the entry  $(\mathbf{E}_C)_{ij} \sim \mathcal{N}(0, \sigma_C^2)$ . For  $T.3$ , we have

$$\begin{aligned} \|\phi(\tilde{\mathbf{Z}}) - \phi(\mathbf{Z})\|_F^2 &= \text{Tr}((\phi(\tilde{\mathbf{Z}}) - \phi(\mathbf{Z}))^T (\phi(\tilde{\mathbf{Z}}) - \phi(\mathbf{Z}))) \\ &= \text{Tr}(\phi^T(\tilde{\mathbf{Z}})\phi(\tilde{\mathbf{Z}}) - \phi^T(\tilde{\mathbf{Z}})\phi(\mathbf{Z}) - \phi^T(\mathbf{Z})\phi(\tilde{\mathbf{Z}}) + \phi^T(\mathbf{Z})\phi(\mathbf{Z})) \\ &= 2d - 2 \underbrace{\text{Tr}(\phi^T(\tilde{\mathbf{Z}})\phi(\mathbf{Z}))}_{T4} \end{aligned} \quad (106)$$

For  $T.4$ , we can obtain

$$\begin{aligned} \text{Tr}(\phi^T(\tilde{\mathbf{Z}})\phi(\mathbf{Z})) &= \sum_{j=1}^d \phi^T(\tilde{\mathbf{z}}_j)\phi(\mathbf{z}_j) = \sum_{j=1}^d \exp\left(-\frac{\|\mathbf{z}_j + (\mathbf{E}_Z)_{:,j} - \mathbf{z}_j\|_2^2}{2r^2}\right) \\ &= \sum_{j=1}^d \exp\left(-\frac{\|(\mathbf{E}_Z)_{:,j}\|_2^2}{2r^2}\right) \\ &\geq d \exp\left(-\frac{\sigma_Z^2 \xi_d^2}{2r^2}\right) \end{aligned} \quad (107)$$

where we use  $\frac{1}{\sigma_Z^2} \|(\mathbf{E}_Z)_{:,j}\|_2^2 \leq \xi_d^2 = d + 2\sqrt{dt} + 2t$  with probability at least  $1 - de^{-t}$  [Laurent and Massart, 2000] since it is the fact that  $\frac{1}{\sigma_Z^2} \|(\mathbf{E}_Z)_{:,j}\|_2^2 \sim \chi_d^2$  where the entry  $(\mathbf{E}_Z)_{ij} \sim \mathcal{N}(0, \sigma_Z^2)$ .



Hence, we can go back to give an upper bound on  $\left\| \phi(\tilde{\mathbf{Z}})\tilde{\mathbf{C}} - \phi(\mathbf{X}) \right\|_{2,\infty}$ :

$$\left\| \phi(\tilde{\mathbf{Z}})\tilde{\mathbf{C}} - \phi(\mathbf{X}) \right\|_{2,\infty} \leq \gamma + \sigma_C \xi_d \sqrt{d} + \tau_C \sqrt{2d \left( 1 - \exp \left( -\frac{\sigma_Z^2 \xi_d^2}{2r^2} \right) \right)} \quad (108)$$

Finally, we have

$$\left\| \tilde{\mathbf{K}}_{xx} - \mathbf{K}_{xx} \right\|_{\infty} \leq \gamma_{zc}(\gamma_{zc} + 2) \quad (109)$$

where

$$\gamma_{zc} = \gamma + \sqrt{d} \left( \sigma_C \xi_d + \tau_C \sqrt{2 \left( 1 - \exp \left( -\frac{\sigma_Z^2 \xi_d^2}{2r^2} \right) \right)} \right). \quad (110)$$

□

## J Proof for Theorem 3.9

**Lemma J.1.** Assume  $\Upsilon = \max_{\mathbf{z}_j, \mathbf{x}'} \|\mathbf{z}_j - \mathbf{x}'\|_2$  and  $\|\mathbf{x} - \mathbf{x}'\|_2 \leq 2\tau_X$ , let  $\{\mathbf{C}_p^S\}_{p=1}^P$  be perturbed by noise drawn from  $\mathcal{N}(0, \sigma^2)$  with the parameter  $\sigma \geq \frac{2c\lambda^{-1}\sqrt{d}\tau_X(\tau_X + \Upsilon)}{r^2\varepsilon}$  for  $c^2 > 2\ln(1.25/\delta)$ . Then, the Gaussian Mechanism that adds noise to  $\{\mathbf{C}_p^S\}_{p=1}^P$  is  $(\varepsilon, \delta)$ -differentially private.

**Lemma J.2.** Assume  $\max_{(p,j)} \{\|\mathbf{x}_{p_j}\|_2, \|\mathbf{x}'_{p_j}\|_2\} \leq \tau_X$ ,  $\max_{(i,j)} \|\mathbf{z}_i - \mathbf{x}_j\|_{\infty} = \Upsilon$ ,  $\|\mathbf{Z}_p^s\|_{sp} \leq \tau_Z \forall s$ , and  $\|\mathbf{C}^S\|_{2,\infty} \leq \tau_C$ , let  $\{\mathbf{Z}_p^s\}_{p=1}^P$  for  $s = 1, \dots, S$  be perturbed by noise drawn from  $\mathcal{N}(0, \sigma^2)$  with variance  $(8S\Delta^2(g_Z) \log(e + (\varepsilon/\delta))/\varepsilon^2)$  where  $\Delta(g_Z) = \frac{2\sqrt{d}\tau_C\tau_X\eta_k}{r^2} \left\{ 1 + (\tau_X + \tau_Z) \frac{(\tau_X + \Upsilon)}{r^2} \right\}$ . Then, the Gaussian Mechanism that adds noise to  $\{\mathbf{Z}_p^s\}_{p=1}^P$  for  $s = 1, \dots, S$  is  $(\varepsilon, \delta)$ -differentially private.

By Lemma J.2, the mechanism that adds Gaussian noise to  $\mathbf{Z}_p^s$  for  $s = 1, \dots, S$  with variance  $(8S\Delta^2(g_Z) \log(e + (\varepsilon_Z/\delta_Z))/\varepsilon_Z^2)$  satisfies  $(\varepsilon_Z, \delta_Z)$ -differential privacy under  $S$ -fold adaptive composition for any  $\varepsilon_Z > 0$  and  $\delta_Z \in (0, 1]$ . By Lemma J.1, the Gaussian Mechanism that injects noise to  $\mathbf{C}_p^S$  with parameter  $\sigma \geq 2c\lambda^{-1}\sqrt{d}\tau_X(\tau_X + \Upsilon)/(r^2\varepsilon_C)$  is  $(\varepsilon_C, \delta_C)$ -differentially private. Therefore, by Theorem 3.16 of [Dwork *et al.*, 2014], the proposed algorithm that adds Gaussian noise to  $\mathbf{Z}_p^s$  for  $s = 1, \dots, S$  and  $\mathbf{C}_p^S$  is  $(\varepsilon_C + \varepsilon_Z, \delta_C + \delta_Z)$ -differentially private. This finished the proof.

## K Proof for Lemma J.1

*Proof.* In our FedSC, for each column of  $\mathbf{C}$ , we have

$$\mathbf{c} = g_{\mathbf{C}}(\mathbf{x}) = \mathbf{G}\mathcal{K}(\mathbf{Z}, \mathbf{x}) = \mathbf{G} \begin{bmatrix} \exp \left( -\frac{\|\mathbf{z}_1 - \mathbf{x}\|_2^2}{2r^2} \right) \\ \vdots \\ \exp \left( -\frac{\|\mathbf{z}_d - \mathbf{x}\|_2^2}{2r^2} \right) \end{bmatrix}, \quad (111)$$

where  $\mathbf{G} = (\mathcal{K}(\mathbf{Z}, \mathbf{Z}) + \lambda\mathbf{I}_d)^{-1}$ . We have

$$\|g_{\mathbf{C}}(\mathbf{x}) - g_{\mathbf{C}}(\mathbf{x}')\|_2 \leq \|\mathbf{G}\|_{sp} \|\mathcal{K}(\mathbf{Z}, \mathbf{x}) - \mathcal{K}(\mathbf{Z}, \mathbf{x}')\|_2, \quad (112)$$

where  $\|\cdot\|_{sp}$  denotes the spectral norm of matrix. Since  $\exp(z)$  is locally Lipschitz continuous when  $z < 0$ , we have

$$\begin{aligned} & \left( \exp \left( -\frac{\|\mathbf{z}_j - \mathbf{x}\|_2^2}{2r^2} \right) - \exp \left( -\frac{\|\mathbf{z}_j - \mathbf{x}'\|_2^2}{2r^2} \right) \right)^2 \\ & \leq \left| \frac{1}{2r^2} \left( \|\mathbf{z}_j - \mathbf{x}\|_2^2 - \|\mathbf{z}_j - \mathbf{x}'\|_2^2 \right) \right|^2 \\ & = \left| \frac{1}{2r^2} \left( \|\mathbf{x} - \mathbf{x}'\|_2^2 + 2\langle \mathbf{z}_j - \mathbf{x}', \mathbf{x}' - \mathbf{x}' \rangle \right) \right|^2 \\ & \leq \left| \frac{1}{2r^2} \left( \|\mathbf{x} - \mathbf{x}'\|_2^2 + 2\|\mathbf{z}_j - \mathbf{x}'\|_2 \|\mathbf{x} - \mathbf{x}'\|_2 \right) \right|^2. \end{aligned} \quad (113)$$

Let  $\Upsilon = \max_{\mathbf{z}_j, \mathbf{x}'} \|\mathbf{z}_j - \mathbf{x}'\|_2$  and  $\|\mathbf{x} - \mathbf{x}'\|_2 \leq 2\tau_X$ . Then the  $\ell_2$ -sensitivity of  $g$  is

$$\begin{aligned}
\Delta_2(g_C) &\leq \sup_{\mathbf{x}, \mathbf{x}'} \|\mathbf{G}\|_\sigma \sqrt{\sum_{j=1}^d \left| \frac{1}{2r^2} (\|\mathbf{x} - \mathbf{x}'\|_2^2 + 2\|\mathbf{z}_j - \mathbf{x}'\|_2 \|\mathbf{x}' - \mathbf{x}'\|_2) \right|^2} \\
&\leq \|\mathbf{G}\|_\sigma \sqrt{\sum_{j=1}^d \left| \frac{1}{2r^2} (4\tau_X^2 + 4\Upsilon\tau_X) \right|^2} \\
&= \|\mathbf{G}\|_\sigma \frac{2\sqrt{d}\tau_X(\tau_X + \Upsilon)}{r^2} \\
&\leq \frac{2\lambda^{-1}\sqrt{d}\tau_X(\tau_X + \Upsilon)}{r^2}.
\end{aligned} \tag{114}$$

Then according to Theorem 3.22 in [Dwork *et al.*, 2014], for  $c^2 > 2\ln(1.25/\delta)$  the Gaussian Mechanism with parameter  $\sigma \geq \frac{2c\lambda^{-1}\sqrt{d}\tau_X(\tau_X + \Upsilon)}{r^2\varepsilon}$  is  $(\varepsilon, \delta)$ -differentially private. This finished the proof.  $\square$

## L Proof for Lemma J.2

*Proof. Proof Sketch* The  $(\varepsilon, \delta)$ -differential privacy of the proposed algorithm can be achieved by injecting noise into  $\mathbf{Z}$  for each local update and into  $\mathbf{C}$  at the final round. To prove this, we first compute the sensitivity of  $\mathbf{Z}$  and  $\mathbf{C}$  for determining the differential privacy of them. Then, we use the adaptive composition [Kairouz *et al.*, 2015] to get the superposition of them which will give the final theoretical result.

Now, the formal proof is as follows.

In our FedSC, consider one-step update of  $\mathbf{Z}$  at client of  $p$

$$\mathbf{Z}_p^{s,t} = \mathbf{Z}_p^{s,t-1} - \eta_t \frac{\partial f}{\partial \mathbf{Z}}(\mathbf{Z}_p^{s,t-1}) \tag{115}$$

where the derivative is given by

$$\frac{\partial f}{\partial \mathbf{Z}}(\mathbf{Z}_p^{s,t-1}) = \frac{1}{r^2}(\mathbf{X}_p \mathbf{W}_Z - \mathbf{Z}_p \bar{\mathbf{W}}_Z) + \frac{2}{r^2}(\mathbf{Z}_p \mathbf{Q}_Z - \mathbf{Z}_p \bar{\mathbf{Q}}_Z) \tag{116}$$

and  $\mathbf{X}_p = [\mathbf{x}_{p_1}, \dots, \mathbf{x}_{p_{j-1}}, \mathbf{x}_{p_j}, \mathbf{x}_{p_{j+1}}, \dots, \mathbf{x}_{p_{N_p}}]$ .

For the simplicity of the proof, we omit the pair of iteration parameters  $(s, t)$  and instead denote two adjacent local updates by  $k$  and  $k-1$  for nonnegative  $k \geq 1$ . Thus, we have the equivalent version of a one-step update of  $\mathbf{Z}_p$  at client  $p$

$$\mathbf{Z}_p^k = \mathbf{Z}_p^{k-1} - \eta_k \left\{ \frac{1}{r^2}(\mathbf{X}_p \mathbf{W}_Z - \mathbf{Z}_p^{k-1} \bar{\mathbf{W}}_Z) + \frac{2}{r^2}(\mathbf{Z}_p^{k-1} \mathbf{Q}_Z - \mathbf{Z}_p^{k-1} \bar{\mathbf{Q}}_Z) \right\} \tag{117}$$

where  $\mathbf{X}_p = [\mathbf{x}_{p_1}, \dots, \mathbf{x}_{p_{j-1}}, \mathbf{x}_{p_j}, \mathbf{x}_{p_{j+1}}, \dots, \mathbf{x}_{p_{N_p}}]$ .

To compute the sensitivity of  $\mathbf{Z}_p$ , we give the counterpart of the above update as

$$(\mathbf{Z}_p^k)' = \mathbf{Z}_p^{k-1} - \eta_k \left\{ \frac{1}{r^2}(\mathbf{X}_p' \mathbf{W}'_Z - \mathbf{Z}_p^{k-1} \bar{\mathbf{W}}'_Z) + \frac{2}{r^2}(\mathbf{Z}_p^{k-1} \mathbf{Q}_Z - \mathbf{Z}_p^{k-1} \bar{\mathbf{Q}}_Z) \right\} \tag{118}$$

where  $\mathbf{X}_p' = [\mathbf{x}_{p_1}, \dots, \mathbf{x}_{p_{j-1}}, \mathbf{x}'_{p_j}, \mathbf{x}_{p_{j+1}}, \dots, \mathbf{x}_{p_{N_p}}]$ .

Next, let's start to derive the upper bound on  $\|\mathbf{Z}_p^k - (\mathbf{Z}_p^k)'\|_F$  term by term.

$$\begin{aligned}
\|\mathbf{Z}_p^k - (\mathbf{Z}_p^k)'\|_F &= \left\| -\frac{\eta_k}{r^2} \{ (\mathbf{X}_p \mathbf{W}_Z - \mathbf{Z}_p^{k-1} \bar{\mathbf{W}}_Z) - (\mathbf{X}_p' \mathbf{W}'_Z - \mathbf{Z}_p^{k-1} \bar{\mathbf{W}}'_Z) \} \right\|_F \\
&= \frac{\eta_k}{r^2} \|(\mathbf{X}_p \mathbf{W}_Z - \mathbf{X}_p' \mathbf{W}'_Z) - \mathbf{Z}_p^{k-1} (\bar{\mathbf{W}}_Z - \bar{\mathbf{W}}'_Z)\|_F \\
&\leq \frac{\eta_k}{r^2} \left\{ \underbrace{\|\mathbf{X}_p \mathbf{W}_Z - \mathbf{X}_p' \mathbf{W}'_Z\|_F}_{\text{T.1}} + \underbrace{\|\mathbf{Z}_p^{k-1} (\bar{\mathbf{W}}_Z - \bar{\mathbf{W}}'_Z)\|_F}_{\text{T.2}} \right\}
\end{aligned} \tag{119}$$

For  $T.1$ , we have

$$\begin{aligned}\|\mathbf{X}_p \mathbf{W}_Z - \mathbf{X}'_p \mathbf{W}'_Z\|_F &= \|\mathbf{X}_p \mathbf{W}_Z - \mathbf{X}'_p \mathbf{W}_Z + \mathbf{X}'_p \mathbf{W}_Z - \mathbf{X}'_p \mathbf{W}'_Z\|_F \\ &\leq \underbrace{\|(\mathbf{X}_p - \mathbf{X}'_p) \mathbf{W}_Z\|_F}_{T.3} + \underbrace{\|\mathbf{X}'_p (\mathbf{W}_Z - \mathbf{W}'_Z)\|_F}_{T.4}\end{aligned}\quad (120)$$

For  $T.3$ , we have

$$\begin{aligned}\|(\mathbf{X}_p - \mathbf{X}'_p) \mathbf{W}_Z\|_F &= \|(\mathbf{x}_{p_j} - \mathbf{x}'_{p_j})(-\mathbf{c}_j^T \odot \mathcal{K}(\mathbf{x}_{p_j}, \mathbf{Z}_p^{k-1}))\|_F \\ &\leq \|\mathbf{x}_{p_j} - \mathbf{x}'_{p_j}\|_2 \|\mathbf{c}_j^T \odot \mathcal{K}(\mathbf{x}_{p_j}, \mathbf{Z}_p^{k-1})\|_2 \\ &= \|\mathbf{x}_{p_j} - \mathbf{x}'_{p_j}\|_2 \sqrt{(\mathbf{c}_j \odot \mathbf{c}_j)^T (\mathcal{K}(\mathbf{Z}_p^{k-1}, \mathbf{x}_{p_j}) \odot \mathcal{K}(\mathbf{Z}_p^{k-1}, \mathbf{x}_{p_j}))} \\ &\leq \|\mathbf{x}_{p_j} - \mathbf{x}'_{p_j}\|_2 \sqrt{\|\mathbf{c}_j \odot \mathbf{c}_j\|_2 \|\mathcal{K}(\mathbf{Z}_p^{k-1}, \mathbf{x}_{p_j}) \odot \mathcal{K}(\mathbf{Z}_p^{k-1}, \mathbf{x}_{p_j})\|_2} \quad (121) \\ &= \|\mathbf{x}_{p_j} - \mathbf{x}'_{p_j}\|_2 \|\mathbf{c}_{p_j}\|_4 \|\mathcal{K}(\mathbf{Z}_p^{k-1}, \mathbf{x}_{p_j})\|_4 \\ &\leq \|\mathbf{x}_{p_j} - \mathbf{x}'_{p_j}\|_2 \|\mathbf{c}_{p_j}\|_2 \|\mathcal{K}(\mathbf{Z}_p^{k-1}, \mathbf{x}_{p_j})\|_2 \\ &\leq \|\mathbf{x}_{p_j} - \mathbf{x}'_{p_j}\|_2 \|C\|_{2,\infty} \|\mathcal{K}(\mathbf{Z}_p^{k-1}, \mathbf{x}_{p_j})\|_2\end{aligned}$$

Here, we use  $\|\mathbf{a} \odot \mathbf{b}\|_2 = \sqrt{(\mathbf{a} \odot \mathbf{a})^T (\mathbf{b} \odot \mathbf{b})}$  for  $\mathbf{a}, \mathbf{b} \in \mathbb{R}^d$  for the second equality; Cauchy-Schwarz inequality for the second inequality;  $\sqrt{\|\mathbf{a} \odot \mathbf{a}\|_2} = \|\mathbf{a}\|_4$  for  $\mathbf{a} \in \mathbb{R}^d$  for the third inequality.

Let  $\Delta_{\mathbf{Z},x} = \mathcal{K}(\mathbf{Z}, \mathbf{x}) - \mathcal{K}(\mathbf{Z}, \mathbf{x}')$ , then we have for  $T.4$ ,

$$\begin{aligned}\|\mathbf{X}'_p (\mathbf{W}_Z - \mathbf{W}'_Z)\|_F &= \|\mathbf{x}'_{p_j} \left\{ -\mathbf{c}_j^T \odot (\mathcal{K}(\mathbf{x}_{p_j}, \mathbf{Z}_p^{k-1}) - \mathcal{K}(\mathbf{x}'_{p_j}, \mathbf{Z}_p^{k-1})) \right\}\|_F \\ &\leq \|\mathbf{x}'_{p_j}\|_2 \|\mathbf{c}_j^T \odot (\mathcal{K}(\mathbf{x}_{p_j}, \mathbf{Z}_p^{k-1}) - \mathcal{K}(\mathbf{x}'_{p_j}, \mathbf{Z}_p^{k-1}))\|_2 \\ &= \|\mathbf{x}'_{p_j}\|_2 \sqrt{(\mathbf{c}_j \odot \mathbf{c}_j)^T (\Delta_{\mathbf{Z}_p^{k-1}, \mathbf{x}_{p_j}} \odot \Delta_{\mathbf{Z}_p^{k-1}, \mathbf{x}_{p_j}})} \\ &\leq \|\mathbf{x}'_{p_j}\|_2 \|C\|_{2,\infty} \|\Delta_{\mathbf{Z}_p^{k-1}, \mathbf{x}_{p_j}}\|_2\end{aligned}\quad (122)$$

Therefore, assume  $\max\{\|\mathbf{x}_{p_j}\|_2, \|\mathbf{x}'_{p_j}\|_2\} \leq \tau_X$ , which means  $\|\mathbf{x} - \mathbf{x}'\|_2 \leq 2\tau_X$ , we can get an upper bound on  $T.1$ .

$$\begin{aligned}\|\mathbf{X}_p \mathbf{W}_Z - \mathbf{X}'_p \mathbf{W}'_Z\|_F &\leq \underbrace{\|(\mathbf{X}_p - \mathbf{X}'_p) \mathbf{W}_Z\|_F}_{T.5} + \underbrace{\|\mathbf{X}'_p (\mathbf{W}_Z - \mathbf{W}'_Z)\|_F}_{T.6} \\ &\leq \|\mathbf{x}_{p_j} - \mathbf{x}'_{p_j}\|_2 \|C\|_{2,\infty} \|\mathcal{K}(\mathbf{Z}_p^{k-1}, \mathbf{x}_{p_j})\|_2 + \|\mathbf{x}'_{p_j}\|_2 \|C\|_{2,\infty} \|\Delta_{\mathbf{Z}_p^{k-1}, \mathbf{x}_{p_j}}\|_2 \\ &= \|C\|_{2,\infty} \left( \|\mathbf{x}_{p_j} - \mathbf{x}'_{p_j}\|_2 \|\mathcal{K}(\mathbf{Z}_p^{k-1}, \mathbf{x}_{p_j})\|_2 + \|\mathbf{x}'_{p_j}\|_2 \|\Delta_{\mathbf{Z}_p^{k-1}, \mathbf{x}_{p_j}}\|_2 \right) \\ &\leq 2\sqrt{d}\tau_C\tau_X \left( 1 + \frac{\tau_X(\tau_X + \Upsilon)}{r^2} \right)\end{aligned}\quad (123)$$

Here, we also used the fact that  $\|\mathcal{K}(\mathbf{Z}_p^{k-1}, \mathbf{x}_{p_j})\|_2 \leq \sqrt{d}$  and the derived bound  $\frac{2\sqrt{d}\tau_X(\tau_X + \Upsilon)}{r^2}$  on  $\|\Delta_{\mathbf{Z}_p^{k-1}, \mathbf{x}_{p_j}}\|_2$  given by K, and  $\|C\|_{2,\infty} \leq \tau_C$  given in Proof F.

Assume  $\|\mathbf{Z}_p^k\|_{sp} \leq \tau_Z \forall k$ , we have for  $T.2$

$$\begin{aligned}
\|\mathbf{Z}_p^{k-1}(\bar{\mathbf{W}}_Z - \bar{\mathbf{W}}'_Z)\|_F &= \|\mathbf{Z}_p^{k-1}(\text{diag}(\mathbf{1}_n^T \mathbf{W}_Z) - \text{diag}(\mathbf{1}_n^T \mathbf{W}'_Z))\|_F \\
&= \|\mathbf{Z}_p^{k-1} \text{diag}(\mathbf{1}_n^T (\mathbf{W}_Z - \mathbf{W}'_Z))\|_F \\
&\leq \|\mathbf{Z}_p^{k-1}\|_{sp} \|\text{diag}(\mathbf{1}_n^T (\mathbf{W}_Z - \mathbf{W}'_Z))\|_F \\
&= \|\mathbf{Z}_p^{k-1}\|_{sp} \|\mathbf{1}_n^T (\mathbf{W}_Z - \mathbf{W}'_Z)\|_2 \\
&\leq \|\mathbf{Z}_p^{k-1}\|_{sp} \|\mathbf{c}_j^T \odot (\mathcal{K}(\mathbf{x}_{p_j}, \mathbf{Z}_p^{k-1}) - \mathcal{K}(\mathbf{x}'_{p_j}, \mathbf{Z}_p^{k-1}))\|_2 \\
&\leq \|\mathbf{Z}_p^{k-1}\|_{sp} \|\mathbf{C}\|_{2,\infty} \|\Delta_{\mathbf{Z}_p^{k-1}, \mathbf{x}_{p_j}}\|_2 \\
&\leq \frac{2\sqrt{d}\tau_Z\tau_C\tau_X(\tau_X + \Upsilon)}{r^2}
\end{aligned} \tag{124}$$

Thus, we get the upper bounds on  $T.1$  and  $T.2$ , respectively, and finally give an upper bound on  $\|\mathbf{Z}_p^k - (\mathbf{Z}_p^k)'\|_F$ .

$$\begin{aligned}
\|\mathbf{Z}_p^k - (\mathbf{Z}_p^k)'\|_F &\leq \frac{\eta_k}{r^2} \left\{ \underbrace{\|\mathbf{X}_p \mathbf{W}_Z - \mathbf{X}'_p \mathbf{W}'_Z\|_F}_{T.1} + \underbrace{\|\mathbf{Z}_p^{k-1}(\bar{\mathbf{W}}_Z - \bar{\mathbf{W}}'_Z)\|_F}_{T.2} \right\} \\
&\leq \frac{\eta_k}{r^2} \left\{ 2\sqrt{d}\tau_C\tau_X \left(1 + \frac{\tau_X(\tau_X + \Upsilon)}{r^2}\right) + \frac{2\sqrt{d}\tau_Z\tau_C\tau_X(\tau_X + \Upsilon)}{r^2} \right\} \\
&= \frac{2\sqrt{d}\tau_C\tau_X\eta_k}{r^2} \left\{ 1 + (\tau_X + \tau_Z) \frac{(\tau_X + \Upsilon)}{r^2} \right\}
\end{aligned} \tag{125}$$

Therefore, if we define  $\mathbf{Z}_p = g_Z(\mathbf{X}_p)$ , the  $\ell_2$ -sensitivity of  $g_Z$  is

$$\begin{aligned}
\Delta_2(g_Z) &= \sup_{\mathbf{X}_p, \mathbf{X}'_p} \|g_Z(\mathbf{X}_p) - g_Z(\mathbf{X}'_p)\|_2 \\
&= \sup_{\mathbf{X}_p, \mathbf{X}'_p} \|\mathbf{Z}_p^k - (\mathbf{Z}_p^k)'\|_F \\
&\leq \frac{2\sqrt{d}\tau_C\tau_X\eta_k}{r^2} \left\{ 1 + (\tau_X + \tau_Z) \frac{(\tau_X + \Upsilon)}{r^2} \right\}
\end{aligned} \tag{126}$$

By Theorem 4.3 of [Kairouz *et al.*, 2015], the mechanism that adds Gaussian noise to  $\mathbf{Z}_p^s$  for  $s = 1, \dots, S$  with variance  $(8S\Delta^2(g_Z) \log(e + (\varepsilon_Z/\delta_Z)))/\varepsilon_Z^2$  satisfies  $(\varepsilon_Z, \delta_Z)$ -differential privacy under  $S$ -fold adaptive composition for any  $\varepsilon_Z > 0$  and  $\delta_Z \in (0, 1]$ . This finished the proof.

□