

In the supplemental materials, we provide the theoretical analysis in Appendix A, and the additional experimental details in Appendix B.

A. Theoretical Analysis

A.1. Optimization Consistency of Model Deviation

In the following, we provide the theorem related to optimization consistency of model deviations. When the aggregating weights, i.e., \mathbf{p} , achieve optimal, the model deviation rate is equally contributed to global updating.

Theorem 3 (Optimization consistency of model deviations). *Rethinking the Lagrangian of dual form in Eq. (10),*

$$\mathbf{J} = \min_{\mathbf{p}} \frac{\eta_g^2}{2\phi} \|(\nabla \log \mathbf{u})^\top \mathbf{p}\|^2 + \lambda_E \mathbf{p}^\top \mathbf{1}, \quad (15)$$

it holds $\nabla \log(u_i(\boldsymbol{\theta}_g^t)) = \nabla \log(u_j(\boldsymbol{\theta}_g^t)), \forall i \neq j \in [K]$.

Proof. By deviating Eq. (15) with regarding to \mathbf{p} , we can obtain $(\nabla \log \mathbf{u})^\top \nabla \log(\mathbf{u})^\top \mathbf{p}^* = -\frac{\phi \lambda_E}{\eta_g^2} \mathbf{I}$. Remind that $\mathbf{d}^* = \frac{\eta_g}{\phi} (\nabla \log \mathbf{u})^\top \mathbf{p}^*$, for client $i \neq j \in [K]$ and $\eta_g \rightarrow 0$, we finally have consistent model deviation change rate as below:

$$\begin{aligned} \lim_{\eta_g \rightarrow 0} c_i(\eta_g, \mathbf{d}^*) &= \nabla \log(u_i(\boldsymbol{\theta}_g^*)) \mathbf{d}^* \\ &= \nabla \log(u_j(\boldsymbol{\theta}_g^*)) \mathbf{d}^* = \lim_{\eta_g \rightarrow 0} c_j(\eta_g, \mathbf{d}^*). \end{aligned} \quad (16)$$

□

Therefore, the global model updates with a direction that balances all model deviation change rates, obtaining consistent parameters for server and client models.

A.2. Bound of Client Model Divergence

In this part, we first introduce mild and general assumptions [23], and induct the model updating divergence bound for each client.

Assumption 5. *Let $F_k(\boldsymbol{\theta})$ be the expected model objective for client k , and assume F_1, \dots, F_K are all L -smooth, i.e., for all $\boldsymbol{\theta}_k, F_k(\boldsymbol{\theta}_k) \leq F_k(\boldsymbol{\theta}_k) + (\boldsymbol{\theta}_k - \boldsymbol{\theta}_k)^\top \nabla F_k(\boldsymbol{\theta}_k) + \frac{L}{2} \|\boldsymbol{\theta}_k - \boldsymbol{\theta}_k\|_2^2$.*

Assumption 6. *Let F_1, \dots, F_N are all μ -strongly convex: for all $\boldsymbol{\theta}_k, F_k(\boldsymbol{\theta}_k) \geq F_k(\boldsymbol{\theta}_k) + (\boldsymbol{\theta}_k - \boldsymbol{\theta}_k)^\top \nabla F_k(\boldsymbol{\theta}_k) + \frac{\mu}{2} \|\boldsymbol{\theta}_k - \boldsymbol{\theta}_k\|_2^2$.*

Assumption 7. *Let ξ_k^t be sampled from the k -th client's local data uniformly at random. The variance of stochastic gradients in each client is bounded: $\mathbb{E} \|\nabla F_k(\boldsymbol{\theta}_k^t, \xi_k^t) - \nabla F_k(\boldsymbol{\theta}_k^t)\|^2 \leq \sigma_k^2$.*

Assumption 8. *The expected squared norm of stochastic gradients is uniformly bounded, i.e., $\mathbb{E} \|\nabla F_k(\boldsymbol{\theta}_k^t, \xi_k^t)\|^2 \leq V^2$ for all $k = 1, \dots, N$ and $t = 1, \dots, T - 1$*

Next, we introduce the lemma related to the bound of client model divergence.

Lemma 2 (Bound of Client Model Divergence). *With assumption 8, η_t is non-increasing and $\eta_t < 2\eta_{t+E}$ (learning rate of t -th round and E -th epoch) for all $t \geq 0$, there exists $t_0 \leq t$, such that $t - t_0 \leq E - 1$ and $\boldsymbol{\theta}_k^{t_0} = \boldsymbol{\theta}^{t_0}$ for all $k \in [N]$. It follows that*

$$\mathbb{E} \left[\sum_k^K p_k \|\boldsymbol{\theta}^t - \boldsymbol{\theta}_k^t\|^2 \right] \leq 4\eta_t^2 (E - 1)^2 V^2. \quad (17)$$

Proof. Let E be the maximal local epoch. For any round $t > 0$, communication rounds from t to t_0 exist $t - t_0 < E - 1$. and the global model $\boldsymbol{\theta}^{t_0}$ and each local model $\boldsymbol{\theta}_k^{t_0}$ are same at round t_0 .

$$\begin{aligned} &\mathbb{E} \left[\sum_k^K p_k \|\boldsymbol{\theta}^t - \boldsymbol{\theta}_k^t\|^2 \right] \\ &= \mathbb{E} \left[\sum_k^K p_k \|(\boldsymbol{\theta}_k^t - \boldsymbol{\theta}^{t_0}) - (\boldsymbol{\theta}^t - \boldsymbol{\theta}^{t_0})\|^2 \right] \end{aligned} \quad (17a)$$

$$\leq \mathbb{E} \sum_k^K p_k \|\boldsymbol{\theta}_k^t - \boldsymbol{\theta}^{t_0}\|^2 \quad (17b)$$

$$= \mathbb{E} \sum_k^K p_k \left\| \sum_{t=t_0}^{t-1} \eta_t \nabla F_k(\boldsymbol{\theta}_k^t, \xi_k^t) \right\|^2 \quad (17c)$$

$$\leq \mathbb{E} \sum_k^K p_k (t - t_0) \sum_{t=t_0}^{t-1} \eta_{t_0}^2 \|\nabla F_k(\boldsymbol{\theta}_k^t, \xi_k^t)\|^2 \quad (17d)$$

$$\leq 4\eta_t^2 (E - 1)^2 V^2, \quad (17e)$$

where the Eq. (17b) holds since $\mathbb{E}(\boldsymbol{\theta}_k^t - \boldsymbol{\theta}^{t_0}) = \boldsymbol{\theta}^t - \boldsymbol{\theta}^{t_0}$, and $\mathbb{E}\|X - \mathbb{E}(X)\| \leq \mathbb{E}\|X\|$, and Eq. (17d) derives from Jensen inequality. □

A.3. Convergence Error Bound

Definition 1 (Heterogeneity Quantification [23]). Let F^* and F_k^* be the minimum values of F and F_k , respectively. We use the term $\Gamma = F^* - \sum_{k=1}^N p_k F_k^*$ for quantifying the degree of non-IID. If the data are IID, then Γ obviously goes to zero as the number of samples grows. If the data are non-IID, then Γ is nonzero, and its magnitude reflects the heterogeneity of the data distribution.

Theorem 4 (Convergence Error Bound). *Let assumptions 5-8 hold, and L, μ, σ_k, V be defined therein. Let $\kappa = \frac{L}{\mu}, \gamma = \max\{8\kappa, E\}$ and the learning rate $\eta_t = \frac{2}{\mu(\gamma+t)}$. The $F \in \mathcal{D}^2$ with full client participation satisfies*

$$\mathbb{E} \left[F(\bar{\boldsymbol{\theta}}^t) \right] - F^* \leq \frac{\kappa}{\gamma + t} \left(\frac{2B}{\mu} + \frac{\mu(\gamma + 1)}{2} \|\boldsymbol{\theta}^t - \boldsymbol{\theta}^*\|^2 \right),$$

where $B = 4(E - 1)^2V^2 + K + 2\Gamma$.

Proof. By L-smooth assumption 5, we can obtain:

$$\begin{aligned} & \mathbb{E} [F(\boldsymbol{\theta}^t) - F(\boldsymbol{\theta}^*)] \\ & \leq \mathbb{E} \left[(\boldsymbol{\theta}^t - \boldsymbol{\theta}^*)^\top \nabla F(\boldsymbol{\theta}^*) + \frac{L}{2} \|\boldsymbol{\theta}^t - \boldsymbol{\theta}^*\|^2 \right] \quad (18) \\ & = \mathbb{E} \left[\frac{L}{2} \|\boldsymbol{\theta}^t - \boldsymbol{\theta}^*\|^2 \right]. \end{aligned}$$

Since the updating in EUA is $\boldsymbol{\theta}^{t+1} = \boldsymbol{\theta}^t - \eta_t \mathbf{d}^t$ for $\mathbf{d}_t^* = \frac{\eta_t}{\phi} (\nabla \log \mathbf{u}^t)^\top \mathbf{p}_t^* = \sum_k^K p_k \nabla F_k(\boldsymbol{\theta}_k^t)$, we can rewrite it as:

$$\begin{aligned} & \|\boldsymbol{\theta}^{t+1} - \boldsymbol{\theta}^*\|^2 \\ & = \|\boldsymbol{\theta}^t - \eta_t \mathbf{d}_t - \boldsymbol{\theta}^*\|^2 \quad (19) \\ & = \|\boldsymbol{\theta}^t - \boldsymbol{\theta}^*\|^2 - 2 \langle \boldsymbol{\theta}^t - \boldsymbol{\theta}^*, \eta_t \mathbf{d}_t \rangle + \eta_t^2 \|\mathbf{d}_t\|^2. \end{aligned}$$

Next, we induce the bound of the second term.

$$\begin{aligned} & \langle \boldsymbol{\theta}^t - \boldsymbol{\theta}^*, \eta_t \mathbf{d}_t \rangle \\ & = \sum_k^K p_k \eta_t \langle \boldsymbol{\theta}^t - \boldsymbol{\theta}_k^t, \nabla F_k(\boldsymbol{\theta}_k^t) \rangle - \sum_k^K p_k \eta_t \langle \boldsymbol{\theta}_k^t - \boldsymbol{\theta}^*, \nabla F_k(\boldsymbol{\theta}_k^t) \rangle \quad (20) \end{aligned}$$

By Cauchy-Schwarz inequality and AM-GM inequality, we have inequality of the first term of Eq. (20):

$$-2 \langle \boldsymbol{\theta}^t - \boldsymbol{\theta}_k^t, \nabla F_k(\boldsymbol{\theta}_k^t) \rangle \leq \frac{1}{\eta_t} \|\boldsymbol{\theta}^t - \boldsymbol{\theta}_k^t\|^2 + \eta_t \|\nabla F_k(\boldsymbol{\theta}_k^t)\|^2. \quad (21)$$

By the μ -strong convexity of $F_k(\cdot)$, we have

$$- \langle \boldsymbol{\theta}_k^t - \boldsymbol{\theta}^*, \nabla F_k(\boldsymbol{\theta}_k^t) \rangle \leq - (F_k(\boldsymbol{\theta}_k^t) - F_k(\boldsymbol{\theta}^*)) - \frac{\mu}{2} \|\boldsymbol{\theta}_k^t - \boldsymbol{\theta}^*\|^2. \quad (22)$$

In Theorem 3, we get $(\nabla \log \mathbf{u})^\top \nabla \log(\mathbf{u}) \mathbf{p}^* = -\frac{\phi \lambda_E}{\eta_t} \mathbf{I}$, which indicates that:

$$\begin{aligned} \|\mathbf{d}_t\|^2 & = \frac{\eta_t^2}{\phi} \|(\nabla \log \mathbf{u}^t)^\top \mathbf{p}\|^2 \\ & = \lambda_E \|\mathbf{p}^\top\|^2 \quad (23) \\ & \leq K, \end{aligned}$$

where the last inequation holds due to $\lambda_E < 1$ and $\|\mathbf{p}\| \leq K$. By combining Eq. (21)-(23) and Lemma 2, it follows that

$$\begin{aligned} & \|\boldsymbol{\theta}^{t+1} - \boldsymbol{\theta}^*\|^2 \\ & = \|\boldsymbol{\theta}^t - \eta_t \mathbf{d}_t - \boldsymbol{\theta}^*\|^2 \\ & \leq \|\boldsymbol{\theta}^t - \boldsymbol{\theta}^*\|^2 + \eta_t \sum_k^K p_k \left(\frac{1}{\eta_t} \|\boldsymbol{\theta}^t - \boldsymbol{\theta}_k^t\|^2 + \eta_t \|\nabla F_k(\boldsymbol{\theta}_k^t)\|^2 \right) \\ & \quad + 2\eta_t \sum_k^K p_k \left(- (F_k(\boldsymbol{\theta}_k^t) - F_k(\boldsymbol{\theta}^*)) - \frac{\mu}{2} \|\boldsymbol{\theta}_k^t - \boldsymbol{\theta}^*\|^2 \right) + \eta_t^2 \|\mathbf{d}_t\|^2 \\ & = (1 - \mu\eta_t) \|\boldsymbol{\theta}^t - \boldsymbol{\theta}^*\|^2 + \sum_k^K p_k \|\boldsymbol{\theta}^t - \boldsymbol{\theta}_k^t\|^2 + \eta_t^2 K + 2\eta_t^2 \Gamma \\ & \leq (1 - \mu\eta_t) \|\boldsymbol{\theta}^t - \boldsymbol{\theta}^*\|^2 + 4\eta_t^2 (E - 1)^2 V^2 + \eta_t^2 K + 2\eta_t^2 \Gamma. \quad (24) \end{aligned}$$

Lastly, let $D_t = \mathbb{E} \|\boldsymbol{\theta}^t - \boldsymbol{\theta}^*\|^2$, it follows that

$$D_{t+1} \leq (1 - \eta_t \mu) D_t + \eta_t^2 B, \quad (25)$$

where $B = 4(E - 1)^2V^2 + K + 2\Gamma$.

For a diminishing stepsize, $\eta_t = \frac{\beta}{t+\gamma}$ for some $\beta > \frac{1}{\mu}$ and $\gamma > 0$ such that $\eta_1 \leq \min\{\frac{1}{\mu}, \frac{1}{4L}\} = \frac{1}{4L}$ and $\eta_t \leq 2\eta_{t+E}$. For $v = \max\{\frac{\beta^2 B}{\beta\mu-1}, (\gamma+1)D_1\}$, by definition, it holds $D_t \leq \frac{v}{\gamma+t}$ for $t = 1$. Assume $D_t \leq \frac{v}{\gamma+t}$ holds, then we expand as below:

$$\begin{aligned} D_{t+1} & \leq (1 - \eta_t \mu) D_t + \eta_t^2 B \\ & \leq \left(1 - \frac{\beta\mu}{t+\gamma}\right) \frac{v}{t+\gamma} + \frac{\beta^2 B}{(t+\gamma)^2} \\ & = \frac{t+\gamma-1}{(t+\gamma)^2} v + \left[\frac{\beta^2 B}{(t+\gamma)^2} - \frac{\beta\mu-1}{(t+\gamma)^2} v \right] \\ & \leq \frac{v}{t+\gamma+1}. \quad (26) \end{aligned}$$

Recall Eq. (18), we finally catch:

$$\mathbb{E} [F(\boldsymbol{\theta}^t) - F(\boldsymbol{\theta}^*)] \leq \frac{L}{2} D_t \leq \frac{L}{2} \frac{v}{\gamma+t}. \quad (27)$$

Following the specific case of [23], we can choose $\beta = \frac{2}{\mu}, \gamma = \max\{8\frac{L}{\mu}, E\} - 1$ and denote $\kappa = \frac{L}{\mu}$, then $\eta_t = \frac{2}{\mu} \frac{1}{\gamma+t}$. One can verify that the choice of η_t satisfies $\eta_t \leq 2\eta_{t+E}$ for $t \geq 1$. Then, we have

$$\begin{aligned} v & = \max \left\{ \frac{\beta^2 B}{\beta\mu-1}, (\gamma+1)\Delta_1 \right\} \\ & \leq \frac{\beta^2 B}{\beta\mu-1} + (\gamma+1)\Delta_1 \quad (28) \\ & \leq \frac{4B}{\mu^2} + (\gamma+1)D_1 \end{aligned}$$

and

$$\mathbb{E} [F(\bar{\boldsymbol{\theta}}^t)] - F^* \leq \frac{L}{2} \frac{v}{\gamma+t} \leq \frac{\kappa}{\gamma+t} \left(\frac{2B}{\mu} + \frac{\mu(\gamma+1)}{2} D_1 \right). \quad (29)$$

As we can see, FedU² similarly converges to a generalization error bound as the FedAvg-like FL model with non-IID data. Discriminatively, benefiting from the optimization of EUA, the communication round multiplies with a smaller B . \square

B. Experimental Supplementary

B.1. Hyper-parameter Sensitivity Analysis

In the following, we study the sensitivity of remaining highly relevant hyper-parameters, i.e., the effect of client numbers and local epochs. Specifically, we compare FedU²-SimCLR and its runner-up method, i.e., FedX-SimCLR, on CIFAR10 $\alpha = 0.1$, by varying the local epochs $E = \{5, 10, 20, 50\}$ in Fig. 7 and the number of clients $K = \{5, 10, 20, 50, 100\}$ in Fig. 8. We train all models until converge to obtain fairly comparable results. As we can see: (1) With the increase of local epochs, each client of FedX-SimCLR obtains a better-performing model, while

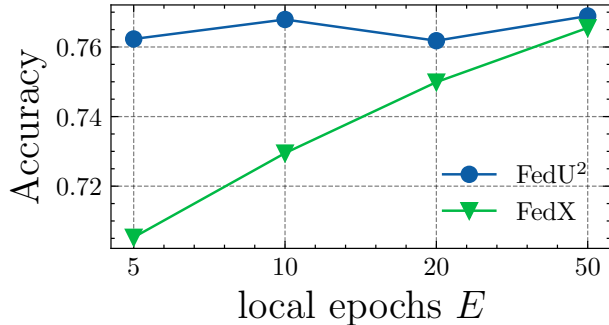


Figure 7. The effect of local epochs E (on CIFAR10 $\alpha = 0.1$).

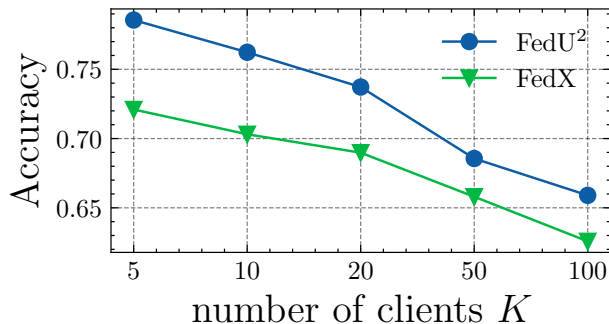


Figure 8. The effect of client number K (on CIFAR10 $\alpha = 0.1$).

each client of FedU²-SimCLR is insensitive. This states that FedU² balances the client model deviation change rate in **EUA**, bringing the benefits of quick convergence. (2) The performance of all methods decreases when the number of clients increases, but FedU²-SimCLR consistently outperforms FedX-SimCLR. It validates that enhancing uniform and unified representations will make **FUSL** methods more generalizable to the cases of various participants amounts.

B.2. Enlarged Figures in Visualization

In our main paper, we depict the top-k singular values of covariance matrix representations in Fig. (3), the corresponding 3-D representation in Fig. (4), and the distribution of data representation in Fig. (5) between global model and randomly sampled local models. The purpose of the above figures is to illustrate the representation enhancement of FedU². In Fig. 9-11, we enlarge these figures to explore the detailed comparisons. In terms of Fig. 11, FedU² keeps the unified representation between global and local models as well as clearer decision boundary for each class.

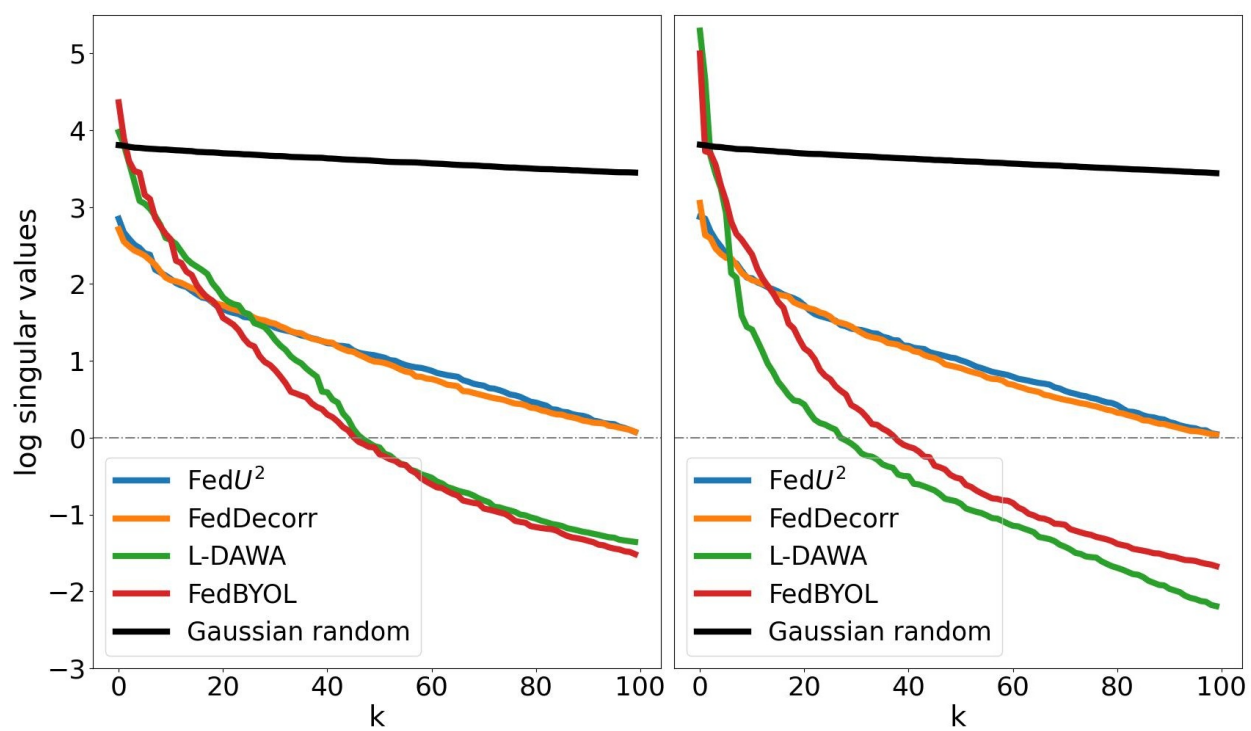


Figure 9. Top k log singular values of the covariance matrix of global model (left) and local model (right) representations.

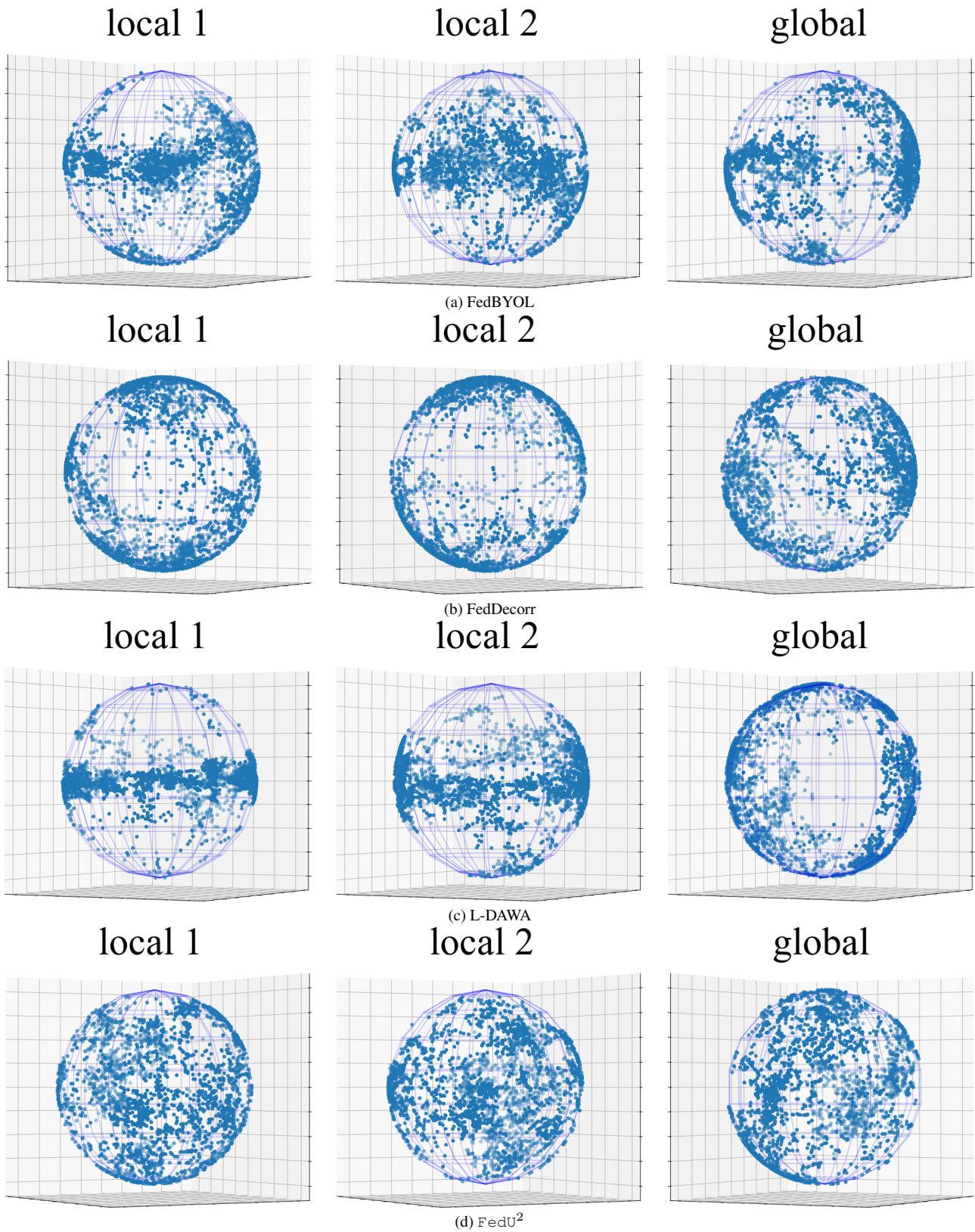
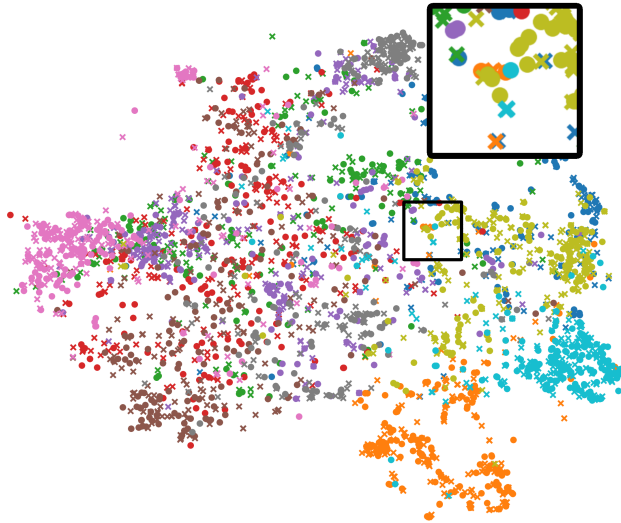
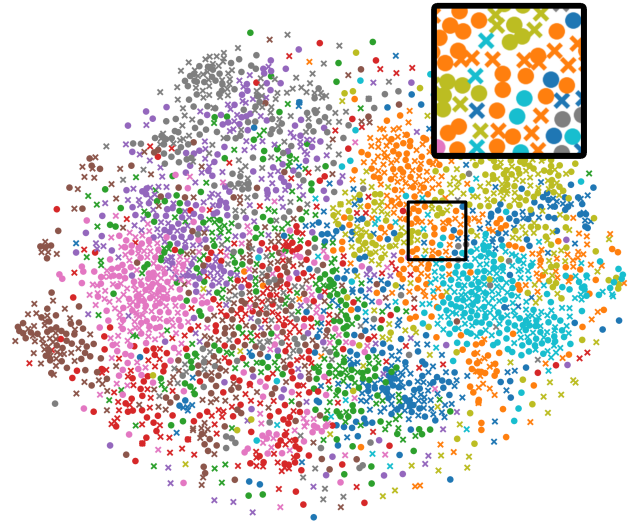


Figure 10. The representations collapse issue on the sphere using BYOL model (on CIFAR10 $\alpha = 0.1$ Cross-silo). The more blank representation space, the more severe collapse issue is.

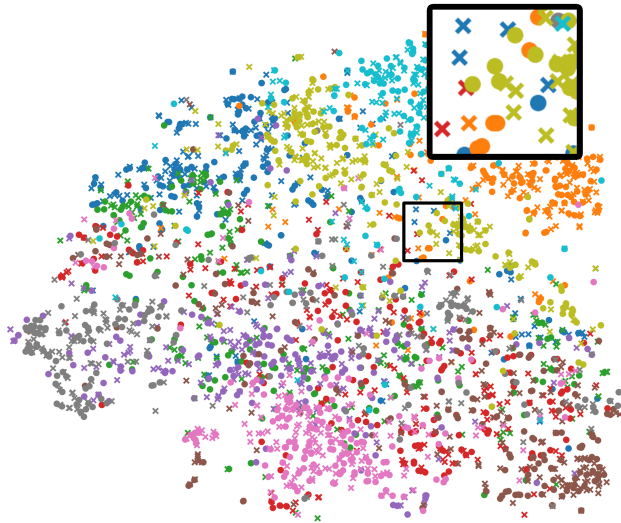
class ID ● 0 ● 1 ● 2 ● 3 ● 4 ● 5 ● 6 ● 7 ● 8 ● 9
 ● global model × local model



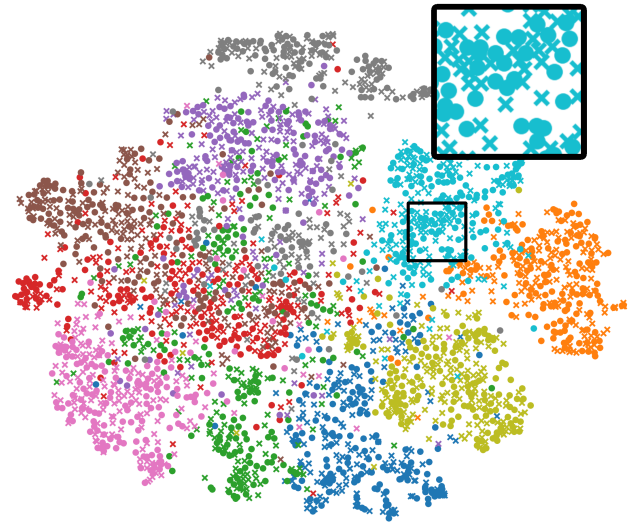
(a) FedBYOL



(b) FedDecorr



(c) L-DAWA



(d) FedU²

Figure 11. The distributions of data representations using global and local BYOL model (on CIFAR10 $\alpha = 0.1$ Cross-silo).