

Appendix of FGTL

Organization of the Appendix. In this Appendix, we first provide the parameter settings for the algorithm implementation. Then, we show two visualization experiments to complement the results we presented in the main paper.

A Implementation Parameter

Each dataset is randomly divided into training (10%), validation (10%), and testing set (80%). For all experiments, we run 3 times and report the average results. For the model architecture, we set the number of layers for GNN-based encoder to 2, the number of layers for MLP-based classifier to 1, and the hidden dimension to 256. For all methods, we set the local training epoch $E = 10$ and the communication round $R = 1500$. For local training, the learning rate is set to $1e-4$ and the weight decay is set to $5e-4$. In generalization phase of FGTL, the communication round is set to 500 and the path length of random walk is set to 10. In transfer phase of FGTL, confidence gate is set to 0.7. Algorithm 1 describes our proposed method, and the code is available on Github¹.

B Additional Experiment

B.1 Performance and Convergence Visualization

We visually verify the advantages brought by each module in FGTL. Figure 1a and 1b show the performance curves of variants of FGTL on the target client, while Figure 1c shows their convergence curves in generalization phase. We can intuitively observe that FGTL-t not only gains a large performance improvement over FGTL-g-t, but also has a significant advantage in convergence speed. Thus, it indicates that global context embedding promotes the convergence of local contrastive learning. Besides, we also notice that both FGTL-t and FGTL-g-t in Figure 1a and 1b show a slight decrease after reaching the peak and then stabilize. We speculate that this is caused by the global model overfitting the source domain labeled data, which is consistent with our motivation to propose consensus knowledge transfer. We observe that FGTL does not show a decline, but stabilizes and slightly improves in the latter half. Therefore, it indicates that consensus knowledge transfer can effectively solve the overfitting problem and enhance the performance of the global model on the target client.

¹<https://github.com/FedGTL/FGTL>

Algorithm 1 The training framework of FGTL+

Training in generalization phase

Input: Source domains $\{\mathcal{G}_S^k\}_{k=1}^K$; Target domain \mathcal{G}_T

Parameter: Communication round R ; Local epoch E ; Source clients number K

Output: Global GNN Encoder ω_g^R

```

1: Initialize global encoder  $\omega_g^0$ .
2: for  $r = 1, 2, \dots R$  do
3:   for  $k = 1, 2, \dots K$  or  $t$  do
4:     Initialize encoder:  $\omega_k$  or  $\omega_t \leftarrow \omega_g^{r-1}$ .
5:      $\hat{\mathcal{G}}_a, \hat{\mathcal{G}}_b \leftarrow \text{Perturb}(\mathcal{G})$ .
6:      $\tilde{\mathcal{G}} \leftarrow \text{Corrupt}(\mathcal{G})$ .
7:      $H^{(P)}, \tilde{H}^{(P)}, \hat{H}_a^{(P)}, \hat{H}_b^{(P)} \leftarrow \text{Encode}(\mathcal{G}; \tilde{\mathcal{G}}; \hat{\mathcal{G}}_a, \hat{\mathcal{G}}_b)$ .
8:      $\hat{s}_a^{(P)}, \hat{s}_b^{(P)} \leftarrow \text{Readout}(\hat{H}_a^{(P)}; \hat{H}_b^{(P)})$ .
9:     Update encoder  $\omega_k$  or  $\omega_t$  locally for  $E$  times.
10:  end for
11:   $\omega_g^r \leftarrow \text{Aggregate}(\{\omega_k\}_{k=1}^K, \omega_t)$ .
12: end for
13: return  $\omega_g^R$ 

```

Training in transfer phase

Input: Source domains $\{\mathcal{G}_S^k\}_{k=1}^K$; Target domain \mathcal{G}_T

Parameter: Communication round R ; Local epoch E ; Source clients number K ; Group number M

Output: Global Classifier ψ_g^R

```

1:  $\{H_k^{(P)}\}_{k=1}^K, H_t^{(P)} \leftarrow \text{Encode}(\{\mathcal{G}_S^k\}_{k=1}^K; \mathcal{G}_T)$  with  $\omega_g^R$ .
2: Initialize global classifier  $\psi_g^0$ .
3: for  $r = 1, 2, \dots R$  do
4:   for  $k = 1, 2, \dots K$  do
5:     Initialize classifier:  $\psi_k \leftarrow \psi_g^{r-1}$ .
6:     Train  $\psi_k$  with  $\{H_k^{(P)}, Y_k\}$  locally for  $E$  times.
7:   end for
8:    $\{\tilde{\psi}_m\}_{m=1}^M \leftarrow \text{Group}(\{\psi_k\}_{k=1}^K)$ .
9:    $\{P_m\}_{m=1}^M \leftarrow \text{Predict}(H_t^{(P)})$  with  $\{\tilde{\psi}_m\}_{m=1}^M$ .
10:   $\{\tilde{p}_i, n_i\}_{i=1}^{|V^t|} \leftarrow \text{Knowledge Voting with } \{P_m\}_{m=1}^M$ .
11:   $\{H_t^{(P)}, \tilde{P}, n\} \leftarrow \text{LPA}(\{\tilde{p}_i, n_i\}_{i=1}^{|V^t|})$ .
12:  Distill  $\{H_t^{(P)}, \tilde{P}, n\}$  to classifier  $\psi_{M+1}$ .
13:   $\psi_g^r \leftarrow \text{Fuse}(\{\psi_m\}_{m=1}^{M+1})$ .
14: end for
15: return  $\psi_g^R$ 

```

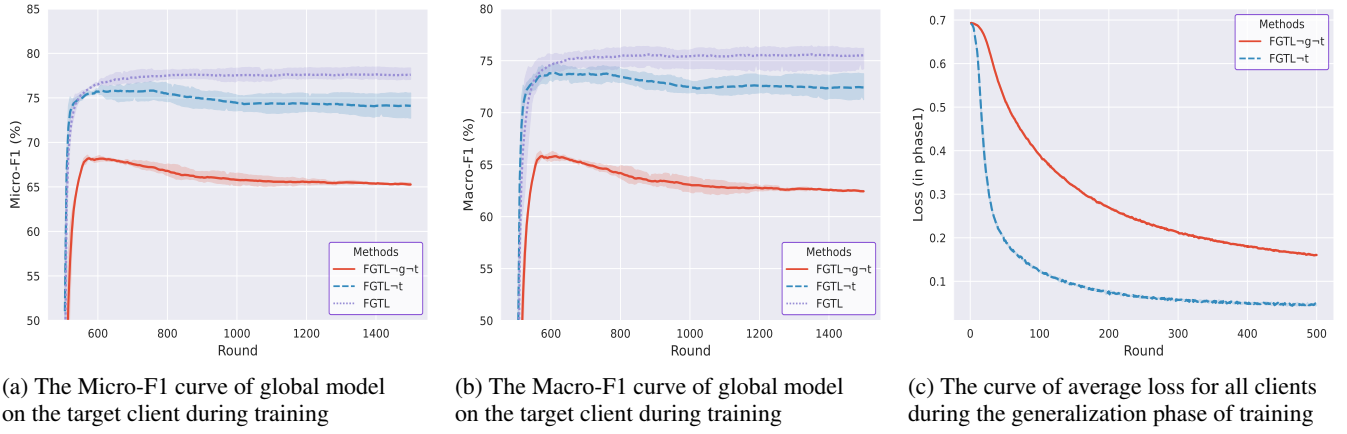


Figure 1: Comparison of performance and convergence curves among variants of FGTL. Target: DBLPv7, Source: ACMv9, Citationv1.

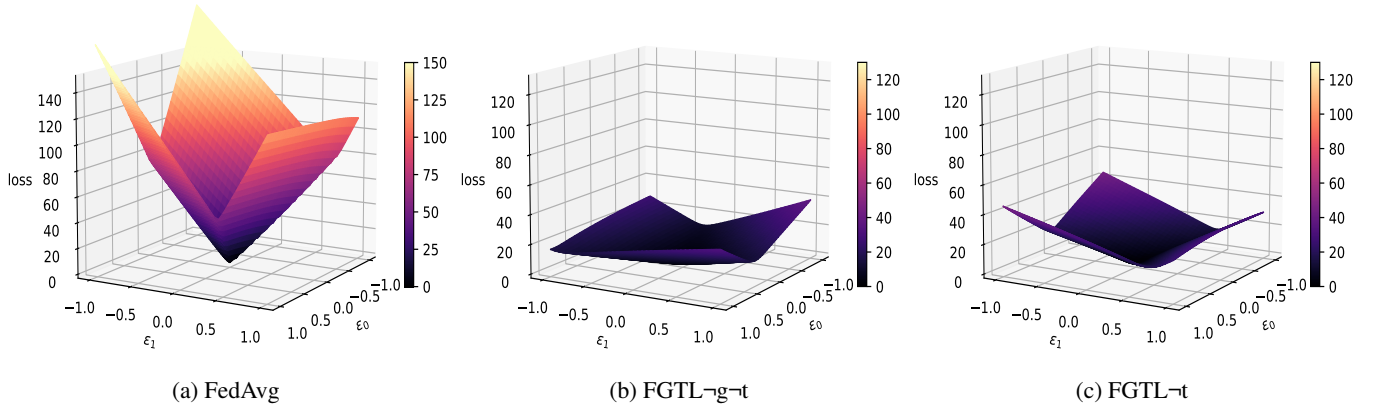


Figure 2: Visualization of the parametric loss landscape on target client with Hessian eigenvectors ϵ_0 and ϵ_1 for global model of each method. Target: Citationv1, Source: ACMv9, DBLPv7.

B.2 Generalization Visualization

To intuitively illustrate that local contrastive learning can improve the generalization of GNN-based encoder, we visually compare the loss landscapes of FedAvg, FGTL-g-t, and FGTL-t on the target client in Figure 2. Recent works in network generalization [Jiang *et al.*, 2019; Keskar *et al.*, 2016] show that networks with lower the top Hessian eigenvalue and Hessian trace are generally less sensitive to small perturbations in the network weights, i.e., good generalization. It can smooth the loss space during training, thus facilitating convergence. We can observe in Figure 2 that FGTL-g-t and FGTL-t are able to smooth the loss landscape, compared to FedAvg. Thus, it suggests that local contrastive learning promotes the generalization of the global model, thus alleviating the client heterogeneity problem.

References

[Jiang *et al.*, 2019] Yiding Jiang, Behnam Neyshabur, Hossein Mobahi, Dilip Krishnan, and Samy Bengio. Fantastic generalization measures and where to find them. *arXiv preprint arXiv:1912.02178*, 2019.

[Keskar *et al.*, 2016] Nitish Shirish Keskar, Dheevatsa Mudigere, Jorge Nocedal, Mikhail Smelyanskiy, and Ping Tak Peter Tang. On large-batch training for deep learning: Generalization gap and sharp minima. *arXiv preprint arXiv:1609.04836*, 2016.