



模式分析与机器智能
工业和信息化部重点实验室
MIT Key Laboratory of
Pattern Analysis & Machine Intelligence

ParNeC | 模式识别与神经计算研究组
Pattern Recognition and Neural Computing

FEDERATED SEMI-SUPERVISED LEARNING WITH INTER-CLIENT CONSISTENCY & DISJOINT LEARNING

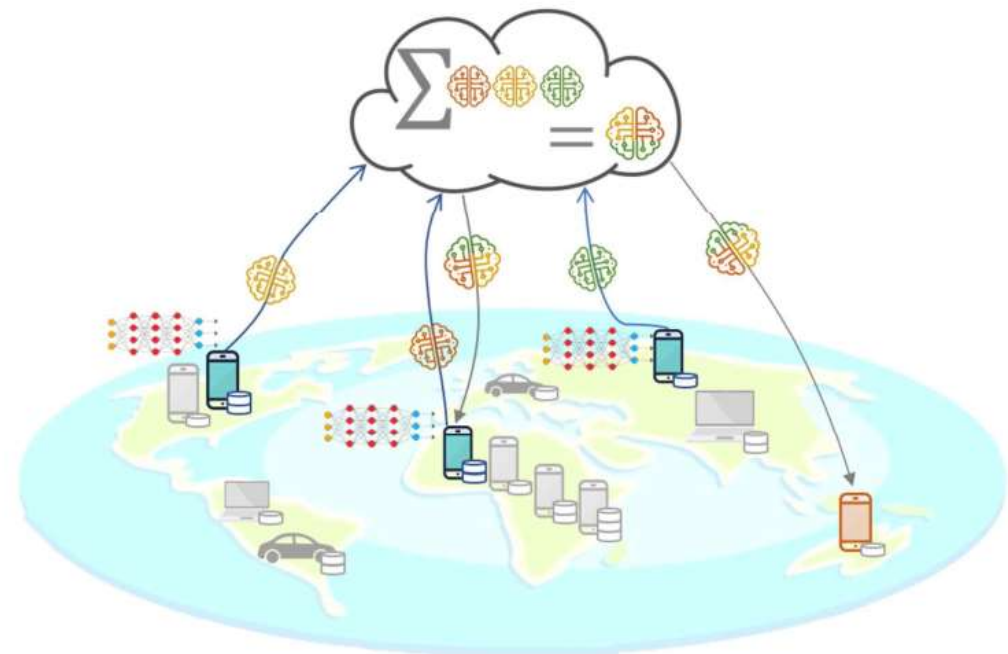
Wonyong Jeong¹, Jaehong Yoon², Eunho Yang^{1,3}, and Sung Ju Hwang^{1,3}
Graduate School of AI¹, KAIST, Seoul, South Korea
School of Computing², KAIST, Daejeon, South Korea
AITRICS³, Seoul, South Korea
{wyjeong, jaehong.yoon, eunhoy, sjhwang82}@kaist.ac.kr

ICLR 2021

Background

Federated Learning (FL)

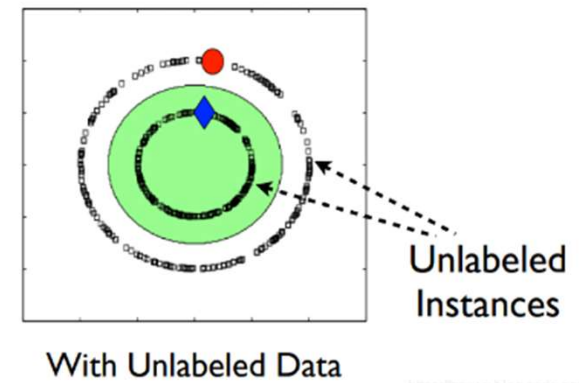
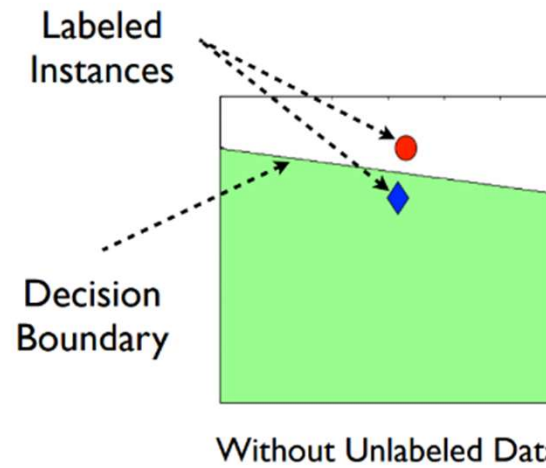
FL can improve model's performance while protecting users' privacy.



Background

Semi-supervised learning (SSL)

Semi-supervised learning can learn knowledge from unlabeled data without high labeling cost.



Related work

Federated Learning

FedAvg

(McMahan et al.,2016)

classic FL algorithm requires many communication rounds to train an effective global model.

FedProx

(Li et al.,2020)

adjusts the local training procedure to pull back local models from global model.

Semi-supervised learning

UDA

(Xie et al., 2019)

use two sets of augmentations, weak and strong, and enforce consistency between the weakly and strongly augmented examples.

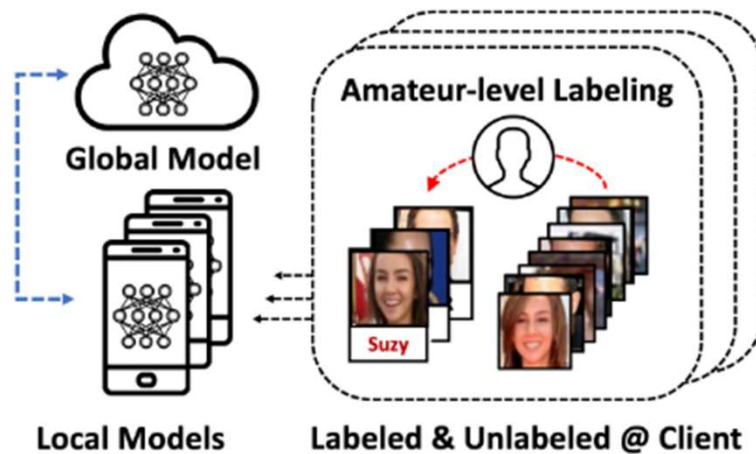
FixMatch

(Sohn et al., 2020)

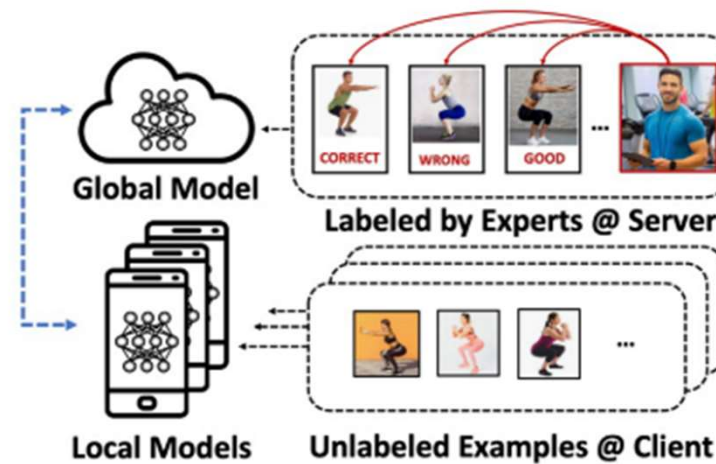
in addition to enforcing consistency between weak-strong augmented pairs, performs pseudo-label refinement on model predictions via thresholding.

Defect in FL

A common limitation is that they only consider supervised learning settings, where the local private data is fully labeled. Yet, the assumption that all of the data examples may include sophisticate annotations is not realistic for real-world applications.



(a) Labels-at-Client Scenario



(b) Labels-at-Server Scenario

Defects in FL + SSL simply

This leads them to practical problems of federated learning with deficiency of labels, namely Federated Semi-Supervised Learning (FSSL).

A naive solution to these scenarios is to simply perform Semi-Supervised Learning (SSL) using any off-the-shelf methods (e.g. FixMatch (Sohn et al., 2020), UDA (Xie et al., 2019)).

This solution has defects:

- does not fully exploit the knowledge of the multiple models trained on heterogeneous data distributions;
- conventional semi-supervised learning approaches are not applicable for scenarios where labeled data is only available at the server;
- even when the labeled data is available at the client, learning from the unlabeled data may lead to forgetting of what the model learned from the labeled data.

INTER-CLIENT CONSISTENCY LOSS

Conventional consistency-regularization methods enforce the predictions from the augmented examples and original (or weakly augmented) instances to output the same class label. $\pi(\mathbf{u})$ performs RandAugment.

$$\|p_{\theta}(\mathbf{y}|\mathbf{u}) - p_{\theta}(\mathbf{y}|\pi(\mathbf{u}))\|_2^2$$

They additionally propose a novel consistency loss called inter-client consistency that regularizes the models learned at multiple clients to output the same prediction, defined as follows:

$$\frac{1}{H} \sum_{j=1}^H \text{KL}[p_{\theta^{h_j}}^*(\mathbf{y}|\mathbf{u}) \| p_{\theta^l}(\mathbf{y}|\mathbf{u})]$$

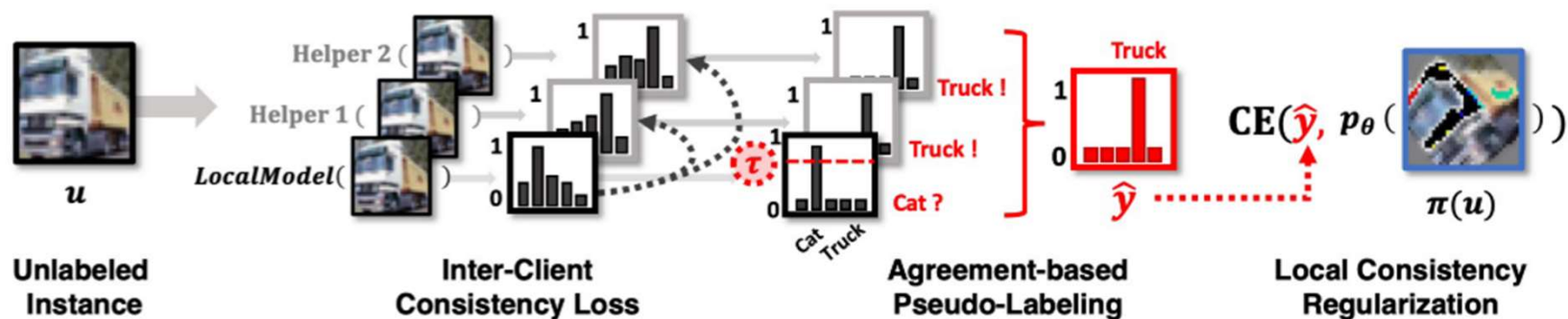
INTER-CLIENT CONSISTENCY LOSS

They also use data-level consistency regularization at each local client similarly to FixMatch. Final consistency regularization term $\Phi(\cdot)$ can be written as follows:

$$\Phi(\cdot) = \text{CrossEntropy}(\hat{\mathbf{y}}, p_{\theta^l}(\mathbf{y}|\pi(\mathbf{u}))) + \frac{1}{H} \sum_{j=1}^H \text{KL}[p_{\theta^{h_j}}^*(\mathbf{y}|\mathbf{u})||p_{\theta^l}(\mathbf{y}|\mathbf{u})]$$

$\hat{\mathbf{y}}$ is the agreement-based pseudo label, defined as follows:

$$\hat{\mathbf{y}} = \text{Max}(\mathbb{1}(p_{\theta^l}^*(\mathbf{y}|\mathbf{u})) + \sum_{j=1}^H \mathbb{1}(p_{\theta^{h_j}}^*(\mathbf{y}|\mathbf{u})))$$



INTER-CLIENT CONSISTENCY LOSS

They select the H helper agents $p_{\theta^{h_j:H}}^*(\mathbf{y}|\mathbf{u})$ for each client as the most relevant models from other clients:

- represent each model by its prediction \mathbf{m} on the same arbitrary input \mathbf{a} located at server (random Gaussian noise). $\mathbf{m}^l = p_{\theta^l}(\mathbf{m}|\mathbf{a})$
- server tries to keep and update all model embeddings $\mathbf{m}^{1:K}$ from clients once each client updates its weights to server.
- server creates K-Dimensional Tree (KD Tree) on \mathbf{m} in the current round r for nearest neighbor search to rapidly select the H helper agents for each client in the next rounds.
- server send helper agents for every 10 rounds, and if a certain client has not yet updated its weights to server in the previous step, then server simply skips sending helpers to the client at the round.

PARAMETER DECOMPOSITION FOR DISJOINT LEARNING

Second, they decompose model parameters θ into two variables, σ for supervised learning and ψ for unsupervised learning, such that $\theta = \sigma + \psi$.

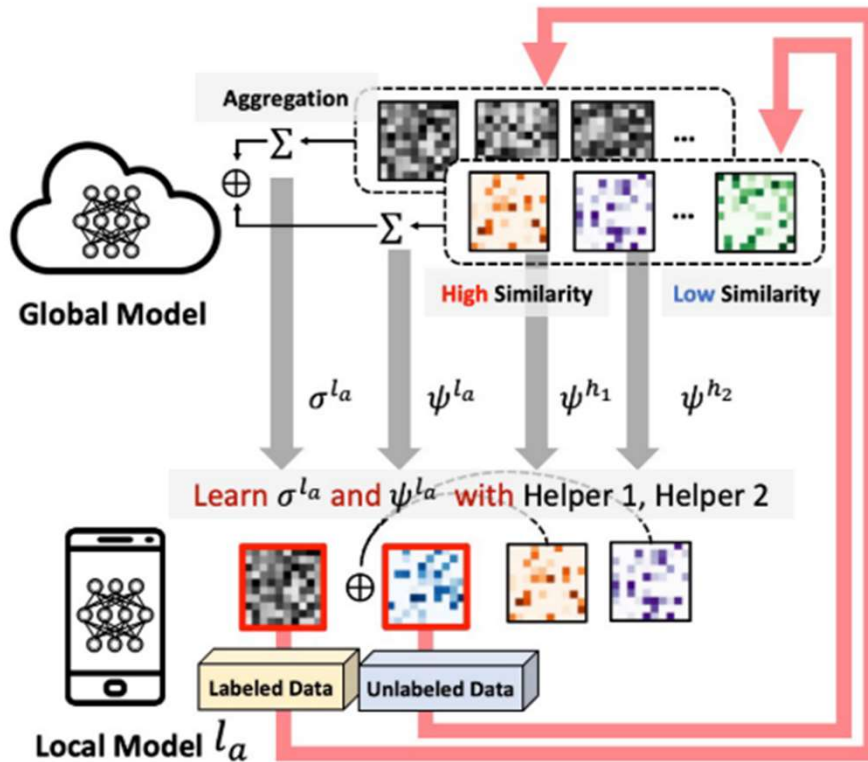
perform standard supervised learning on σ , while keeping ψ fixed during training on labeled data, by minimizing the loss term as follows:

$$\text{minimize } \mathcal{L}_s(\sigma) = \lambda_s \text{CrossEntropy}(\mathbf{y}, p_{\sigma+\psi^*}(\mathbf{y}|\mathbf{x}))$$

perform unsupervised learning conversely on ψ , while keeping σ fixed for the learning phase on unlabeled data, by minimizing the consistency loss terms as follows:

$$\text{minimize } \mathcal{L}_u(\psi) = \lambda_{\text{ICCS}} \Phi_{\sigma^*+\psi}(\cdot) + \lambda_{L_2} \|\sigma^* - \psi\|_2^2 + \lambda_{L_1} \|\psi\|_1$$

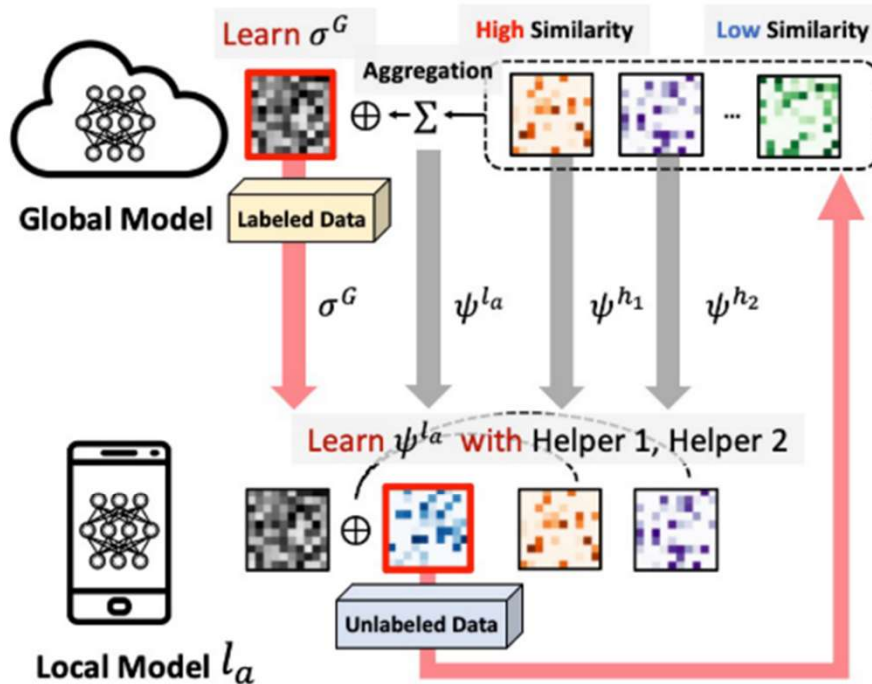
The Overall Algorithm



Algorithm 1 Labels-at-Client Scenario

- 1: **RunServer()**
- 2: initialize σ^0 and ψ^0
- 3: **for** each round $r = 1, 2, \dots, R$ **do**
- 4: $\mathcal{L}^r \leftarrow$ (select random A clients from \mathcal{L})
- 5: **for** each client $l_a \in \mathcal{L}^r$ **in parallel do**
- 6: $\psi_{1:H}^r \leftarrow$ GetNearestNeighbors(ψ^r)
- 7: $\sigma_a^r, \psi_a^r \leftarrow$ RunClient($\sigma^r, \psi^r, \psi_{1:H}^r$)
- 8: EmbedLocalModel(σ_a^r, ψ_a^r)
- 9: **end for**
- 10: $\sigma^{r+1} \leftarrow \frac{1}{A} \sum_{a=1}^A (\sigma_{l_a}^r)$
- 11: $\psi^{r+1} \leftarrow \frac{1}{A} \sum_{a=1}^A (\psi_{l_a}^r)$
- 12: **end for**
- 13: **RunClient**($\sigma, \psi, \psi_{1:H}$)
- 14: $\theta_{l_a} \leftarrow \sigma + \psi, \theta_{h_{1:H}} \leftarrow \sigma + \psi_{1:H}$
- 15: **for** each local epoch e from 1 to E_L **do**
- 16: **for** minibatch $s \in \mathcal{S}_{l_a}$ and $u \in \mathcal{U}_{l_a}$ **do**
- 17: $\theta_{\sigma+\psi^*} \leftarrow \theta_{\sigma+\psi^*} - \eta \nabla \ell_s(\theta_{\sigma+\psi^*}; \theta_{h_{1:H}}, s)$
- 18: $\theta_{\sigma^*+\psi} \leftarrow \theta_{\sigma^*+\psi} - \eta \nabla \ell_u(\theta_{\sigma^*+\psi}; \theta_{h_{1:H}}, u)$
- 19: **end for**
- 20: **end for**

The Overall Algorithm



Algorithm 2 Labels-at-Server Scenario

```

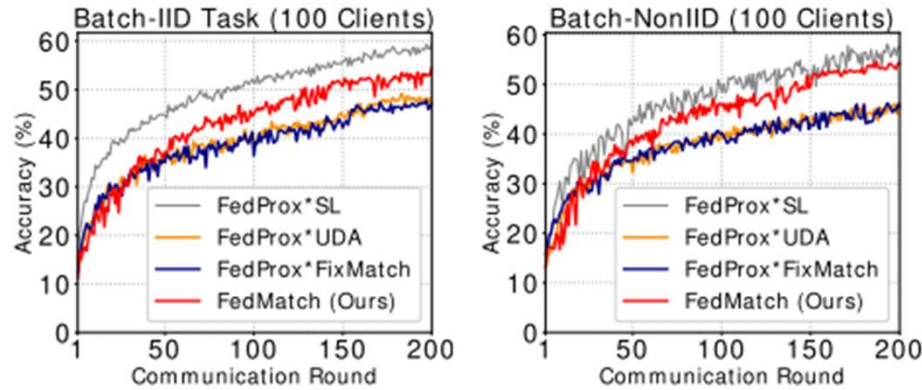
1: RunServer()
2: initialize  $\sigma^0, \psi^0$ 
3: for each round  $r = 1, 2, \dots, R$  do
4:   for each server epoch  $e$  from 1 to  $E_G$  do
5:     for minibatch  $s \in \mathcal{S}_G$  do
6:        $\theta_{\sigma+\psi^*} \leftarrow \theta_{\sigma+\psi^*} - \eta \nabla \ell_s(\theta_{\sigma+\psi^*}; s)$ 
7:     end for
8:   end for
9:    $\mathcal{L}^r \leftarrow$  (select random  $A$  clients from  $\mathcal{L}$ )
10:  for each client  $l_a^r \in \mathcal{L}^r$  in parallel do
11:     $\psi_{1:H}^r \leftarrow$  GetNearestNeighbors( $\psi^r$ )
12:     $\psi_a^r \leftarrow$  RunClient( $\sigma^{r+1}, \psi^r, \psi_{1:H}^r$ )
13:    EmbedLocalModel( $\sigma^{r+1}, \psi_a^r$ )
14:  end for
15:   $\psi^{r+1} \leftarrow \frac{1}{A} \sum_{a=1}^A (\psi_{l_a}^r)$ 
16: end for
17: RunClient( $\sigma, \psi, \psi_{1:H}$ )
18:  $\theta_l \leftarrow \sigma^* + \psi, \theta_{h_{1:H}} \leftarrow \sigma^* + \psi_{1:H}$ 
19: for each local epoch  $e$  from 1 to  $E_L$  do
20:   for minibatch  $u \in \mathcal{U}_{l_a}$  do
21:      $\theta_{\sigma^*+\psi} \leftarrow \theta_{\sigma^*+\psi} - \eta \nabla \ell_u(\theta_{\sigma^*+\psi}; \theta_{h_{1:H}}, u)$ 
22:   end for
23: end for

```

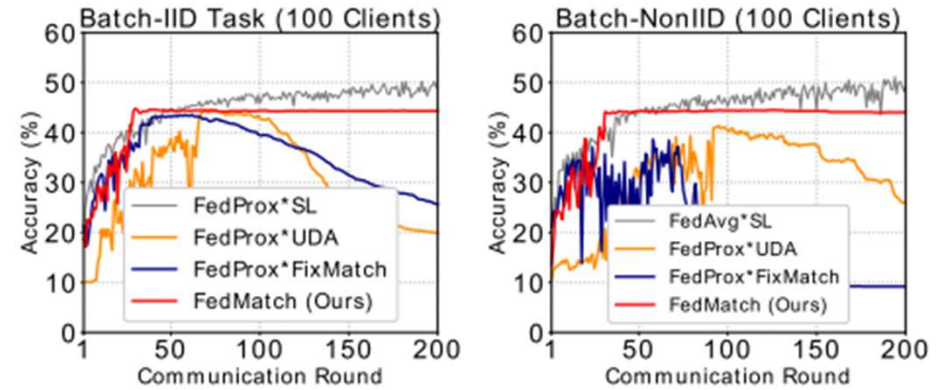
Performance Comparison on Batch-IID & NonIID Tasks

CIFAR-10, Batch-IID Task with 100 Clients ($K=100, F=0.05, H=2$)						
Methods	Labels-at-Client Scenario			Labels-at-Server Scenario		
	Acc.(%)	S2C Cost	C2S Cost	Acc.(%)	S2C Cost	C2S Cost
FedAvg-SL	58.60 ± 0.42	100 %	100 %	52.45 ± 0.23	100 %	100 %
FedProx-SL	59.30 ± 0.31	100 %	100 %	49.11 ± 0.38	100 %	100 %
FedAvg-UDA	46.35 ± 0.29	100 %	100 %	24.81 ± 0.73	100 %	100 %
FedProx-UDA	47.45 ± 0.21	100 %	100 %	19.91 ± 0.31	100 %	100 %
FedAvg-FixMatch	47.01 ± 0.43	100 %	100 %	11.95 ± 0.60	100 %	100 %
FedProx-FixMatch	47.20 ± 0.12	100 %	100 %	25.61 ± 0.32	100 %	100 %
FedMatch (Ours)	52.13 ± 0.34	79 %	46 %	44.95 ± 0.49	45 %	22 %
CIFAR-10, Batch-NonIID Task with 100 Clients ($K=100, F=0.05, H=2$)						
Methods	Labels-at-Client Scenario			Labels-at-Server Scenario		
	Acc.(%)	S2C Cost	C2S Cost	Acc.(%)	S2C Cost	C2S Cost
FedAvg-SL	55.15 ± 0.21	100 %	100 %	51.50 ± 0.51	100 %	100 %
FedProx-SL	57.75 ± 0.15	100 %	100 %	49.31 ± 0.18	100 %	100 %
FedAvg-UDA	44.35 ± 0.39	100 %	100 %	27.61 ± 0.71	100 %	100 %
FedProx-UDA	46.31 ± 0.63	100 %	100 %	26.01 ± 0.78	100 %	100 %
FedAvg-FixMatch	46.20 ± 0.52	100 %	100 %	09.45 ± 0.34	100 %	100 %
FedProx-FixMatch	45.55 ± 0.63	100 %	100 %	09.21 ± 0.24	100 %	100 %
FedMatch (Ours)	52.25 ± 0.81	85 %	49 %	44.17 ± 0.19	42 %	20 %

Performance Comparison on Batch-IID & NonIID Tasks



(a) Labels-at-Client Scenario



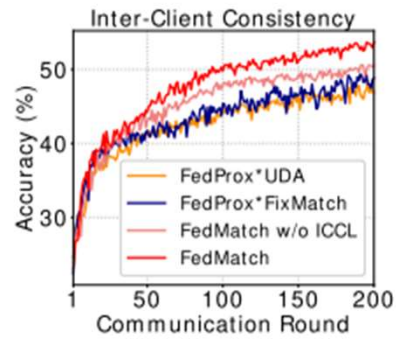
(b) Labels-at-Server Scenario

Averaged Local Performance on Streaming-NonIID Task

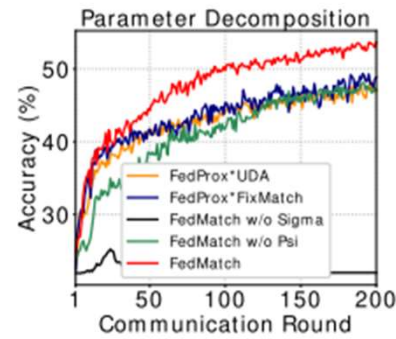
Fashion-MNIST, Streaming-NonIID Task with 10 Clients ($K=10, F=1.0, H=2$)						
Methods	<i>Labels-at-Client Scenario</i>			<i>Labels-at-Server Scenario</i>		
	Acc.(%)	S2C Cost	C2S Cost	Acc.(%)	S2C Cost	C2S Cost
Local-SL	87.19 ± 0.36	N/A	N/A	N/A	N/A	N/A
Local-UDA	70.70 ± 0.28	N/A	N/A	N/A	N/A	N/A
Local-FixMatch	62.62 ± 0.32	N/A	N/A	N/A	N/A	N/A
FedProx-SL	82.06 ± 0.26	100 %	100 %	77.43 ± 0.42	100 %	100 %
FedProx-UDA	73.71 ± 0.17	100 %	100 %	83.34 ± 0.21	100 %	100 %
FedProx-FixMatch	62.40 ± 0.43	100 %	100 %	73.71 ± 0.32	100 %	100 %
FedMatch (Ours)	77.95 ± 0.14	37 %	48 %	84.15 ± 0.31	14 %	63 %

Experiment

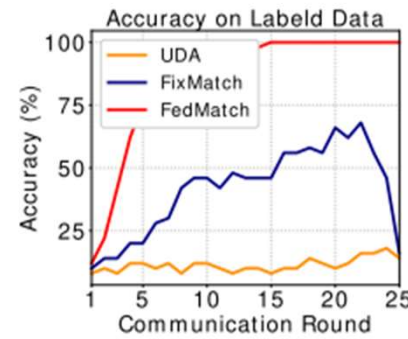
Ablation Study



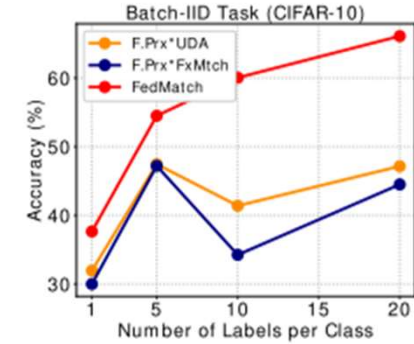
(a) Inter-Client Consistency



(b) Sigma & Psi



(c) Inter-Task Interference



(d) Number of Labels

Thanks