# Joint Focal Loss and Dominant Gradient Correction for Gradient Conflict in Federated Learning

**Jiajun Wang[b], Yingchi Mao[a,b], Zibo Wang[b], Jun Wu[b] and Jie Wu[c]**

[a] Key Laboratory of Water Big Data Technology of Ministry of Water Resources, Hohai University, Nanjing, China
[b] School of Computer and Information, Hohai University, Nanjing, China
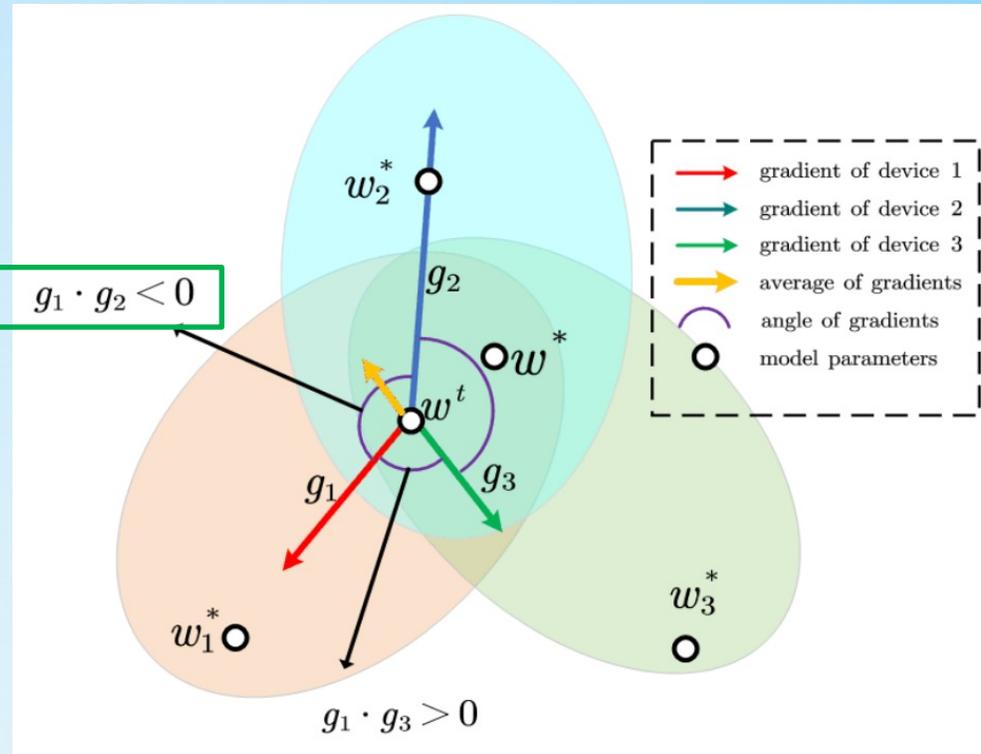[c] Center for Networked Computing, Temple University, Philadelphia USA

# Content

■ Background

　　In federated learning, the data collected by IoT devices are typically heterogeneous .In this case, as the training advances the local model gradually converges to its local target optimum, resulting in a conflict among the upl　　oaded gradients in global aggregation phase.

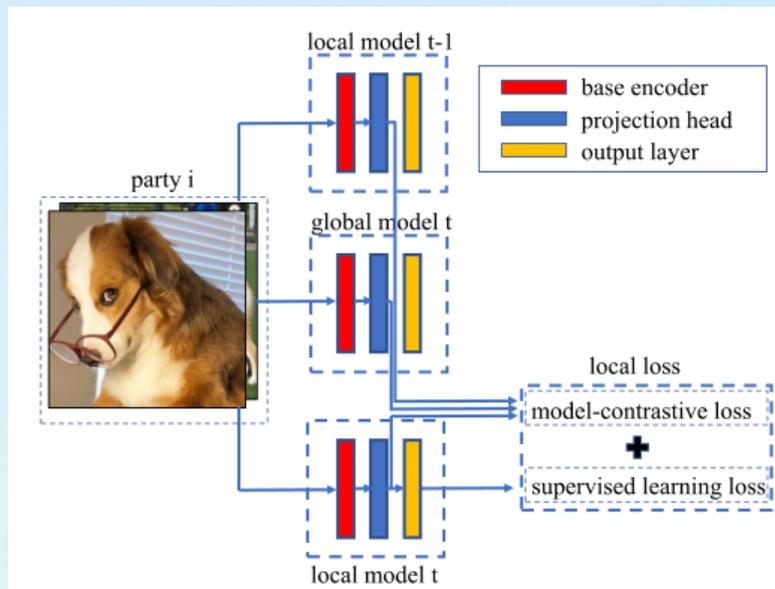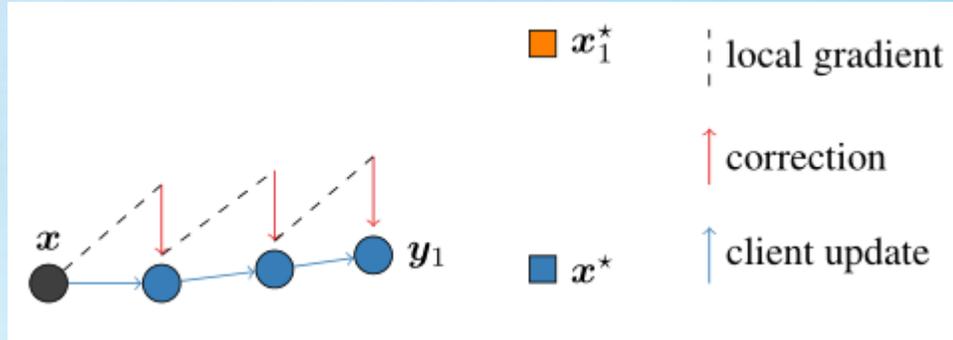Gradient conflict

$g_1 \cdot g_2 < 0$

Impairs the global accuracy

■ Optimization object
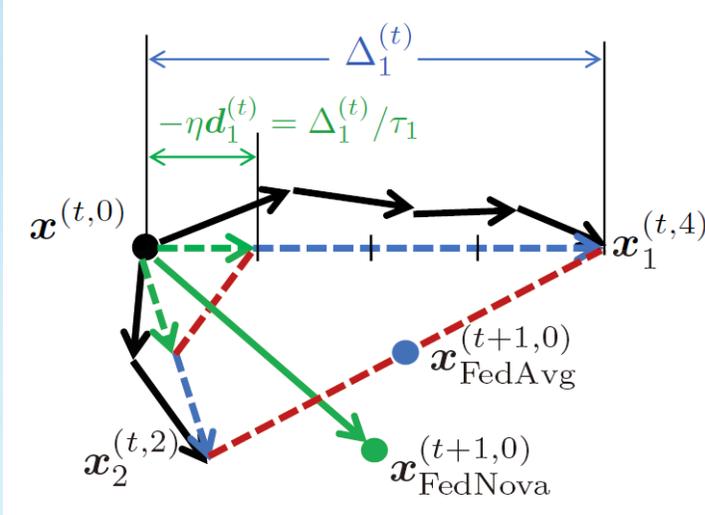
Reduce the conflict among uploaded gradients.

■ Goal

Improve the accuracy of the global model.

# Local training optimization



- SCAFFOLD was proposed as a new stochastic control averaging algorithm

- Li et al corrected local updates of clients by injecting projection heads into the model with a model-level comparison learning method.

- Although local model divergence is suppressed, model preferences can still lead to gradient conflicts
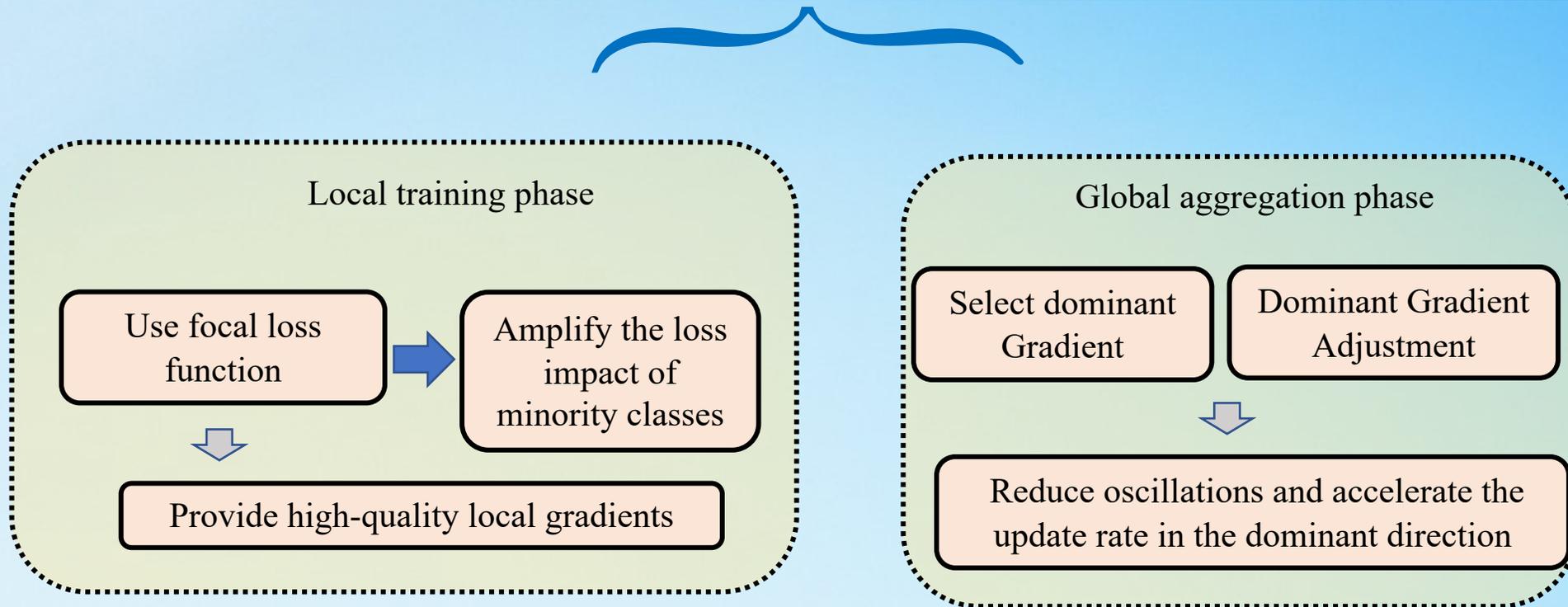
<br>

# Global aggregation optimization



- Wang et al. [10] presented FedNova, which standardized and adjusted local updates in accordance with the amount of local iterations.

- Adaptive schemes such as ADAM, YOGI, and ADAGRAD were introduced on the parameter server to improve the global model accuracy.

- Although the above methods make the global aggregation smoother, the gradient conflict problem is not properly solved.

## Federated Learning Mitigating Gradient Conflict(FedMGC)

## Handling of Class Imbalance
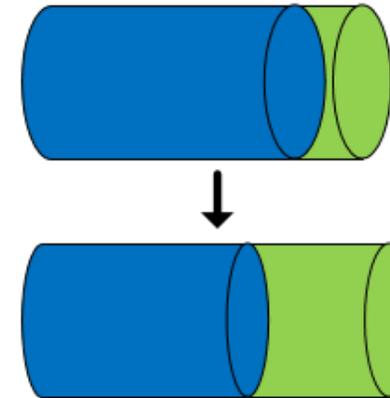
Cross Entropy loss function

↓ Class imbalance

Focal Loss function
$$FL(p_t) = -\beta(1 - p_t)^\gamma \log(p_t)$$

Provide high-quality local gradients



● majority class
● minority class

*Increase the proportion of minority class*

## Dominant Gradient Generation(DGG)

**Algorithm 1** Dominant Gradient Generation

**Input:** Local gradients $g_t = \{g_1^t, ..., g_K^t\}$, loss values $l_t = \{l_1^t, ..., l_K^t\}$, dominant gradient selection ratio $\lambda$.

**Output:** Dominant gradient $dg^t$.

1: **Initialize** $p^t = \{\}$, $z^t = \{\}$
2: **for** $g_i^t \in g_t$ **do**
3:     **for** $g_j^t \in g_t$ **do**
4:       $p_{i,j}^t = \left( \frac{g_i \cdot g_j}{\|g_j\|} + \frac{g_j g_i}{\|g_i\|} \right) / 2$
5:     **end for**
6: **end for**
7: **for** $i = 1, \cdots, K$ **do**
8:     add $p_i^t = \sum_{j=1}^{K} p_{i,j}^t$ to $p^t$
9: **end for**
10: **for** $i = 1, \cdots, K$ **do**
11:     add $z_i^t = \frac{p_i^t}{l_i^t}$ to $z^t$.
12: **end for**
13: sort array $z^t$ in descending order to achieve $z_s = \{z_{s_1}^t, ..., z_{s_K}^t\}$, then choose top $\lceil \lambda K \rceil$ gradients $dg^t = \{g_{s_1}^t, ..., g_{s_{\lceil \lambda K \rceil}}^t\}$ as dominant gradients which map to $z_s$.
14: **return** $dg^t$

Calculate the of mutual projection between two gradients, then calculate the gradient projection outlier of client i

By dividing the gradient projection outlier with loss value to get gradient outliers

Sort in descending order and select top ⌈λK⌉ gradients as the dominant gradients

## Dominant Gradient Adjustment

**Algorithm 2** Dominant Gradient Adjustment

**Input:** Local gradients $g_t = \{g_1^t, ..., g_K^t\}$, dominant gradients $dg^t = \left\{ g_{s_1}^t, ..., g_{s_{\lceil \lambda K \rceil}}^t \right\}$.
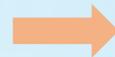
**Output:** Corrected gradients $g^t$.

1: **Initialize** $n = |dg| = \lceil \lambda K \rceil$, $m = |g_t|$, $g_i^{cur} \leftarrow g_i^t$
2: **for** $i < m$ **do**
3:    **for** $j < b$ **do**
4:      **if** $g_i^{cur} \cdot g_{s_j} < 0$ and $i \neq s_j$ **then**
5:        $g_i^{cur} = g_i^{cur} - \frac{g_i^{cur} \cdot g_{s_j}^t}{\|a_{s_j}^t\|^2} g_{s_j}^t$
6:      **end if**
7:    **end for**
8:  **end for**
9: $g^t = \frac{1}{m} \sum_{i=1}^m g_i^{cur}$
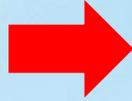10: **return** $g^t$

Detects gradient conflict

Correct the gradient

Aggregate the corrected gradients

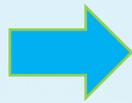| Dataset | Model |
|---|---|
| MNIST | Two convolutional layers and two fully connected layers. |
| CIFAR-10 | Two convolutional layers and three fully connected layers |
| CIFAR-100 | Two convolutional layers and three fully connected layers |

**Dataset & Model**

**Dirichlet distribution: q ~Dir(α).** A smaller value of α indicates stronger data heterogeneity

**Parameter Settings**

| Parameter | Value |
|---|---|
| Batch size | 128 |
| Local epoch | 10 |
| Local learning rate | 0.001 |

**Baselines**

**FedAvg, FedProx, SCAFFOLD, FedNova**

**Indicators**

**Test accuracy and class loss**

## Ablation Experiments

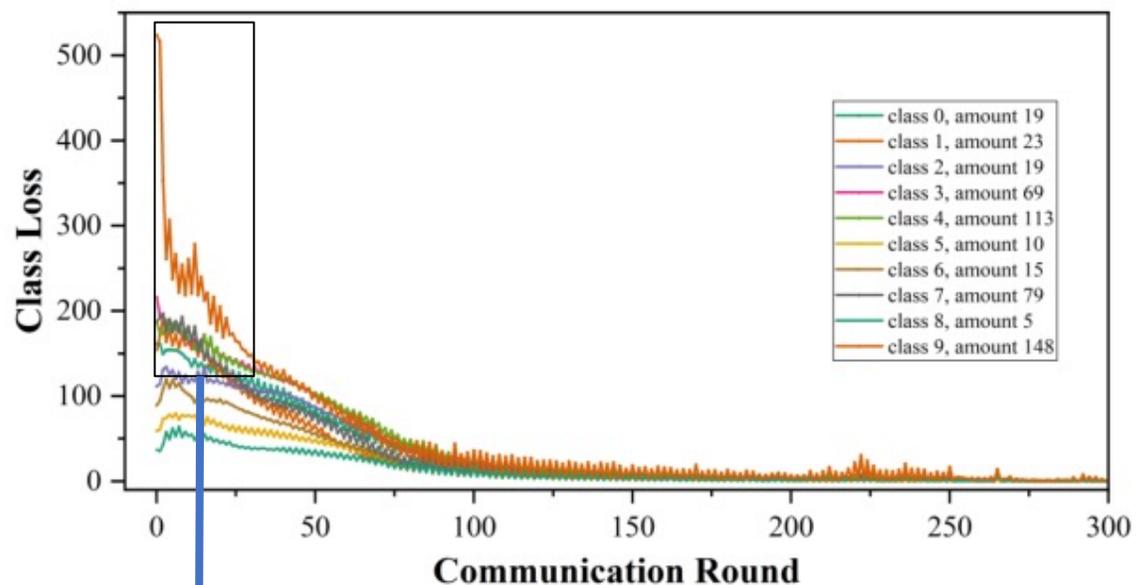**Settings**: $\alpha = 1$, the number of clients is 100, the client participation rate is 10%, CIFAR-10 dataset.



Fig. 1. Test accuracy of FedMGC and related ablation methods.



Fig. 2. Class loss variation in a client with class imbalance.

The loss of majority classes rapidly decline, increasing the proportion of minority classes.

Focal loss and DGC can effectively alleviate the gradient conflict

## Analysis of test accuracy while full client participation



**Settings**: $\alpha = 0.5$, the number of client is 10, client participation rate is 100% , CIFAR-10 dataset.

FedMGC achieves higher test accuracy among all methods, which indicates that the global accuracy can be improved by mitigating the gradient conflict.

Fig. 3. Test accuracy of FedMGC and baselines with full client participation.

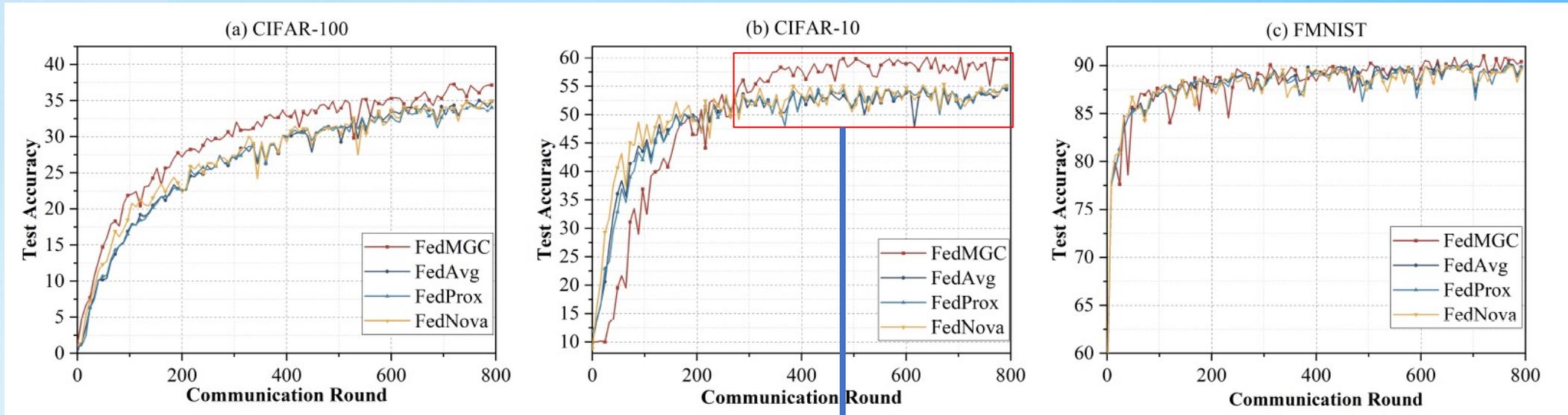## Analysis of test accuracy while client partial participation



Fig. 4. Test Accuracy of FedMGC with baselines for α = 0.5 on CIFAR-100, CIFAR-10 and FMNIST datasets.

**Settings**: α = 0.5, the number of clients is 100, the client participation rate is 10%, on CIFAR-100,CIFAR-10and FMNIST dataset.

6.5%, 5.5%, and 5.2% higher FedAvg, FedProx, and FedNova.

The gradient conflict seriously affects the global accu racy, and FedMGC can effectively mitigate the gradient conflict

## Conclusion

➤ FedMGC increases the loss contribution of minority classes and corrects gradients with conflict.

➤ FedMGC is able to achieve higher test accuracy on various heterogeneity of data.

## Future Work

➤ Convergence analysis of FedMGC.

➤ Further optimize the FL function to reduce the tuning of the hyperparameter.

Thanks for your attention!