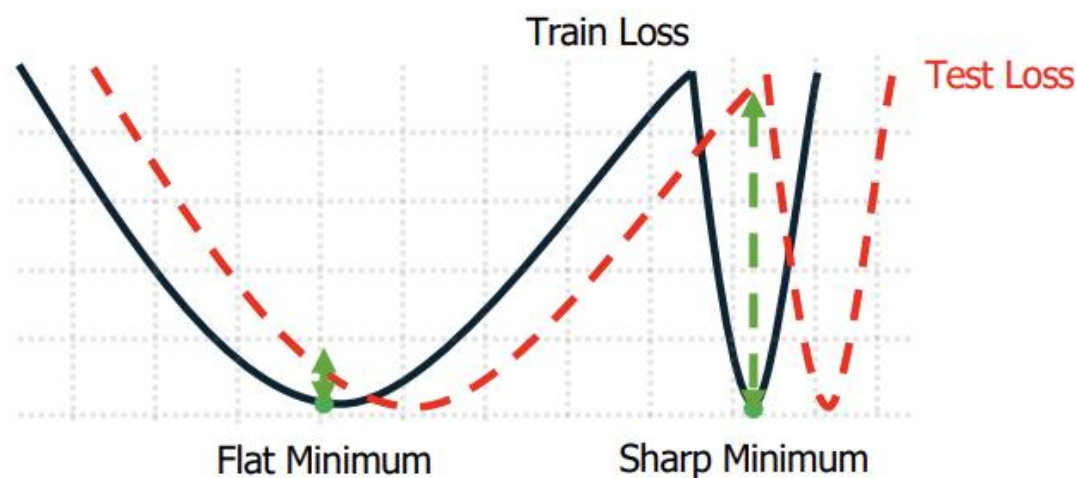


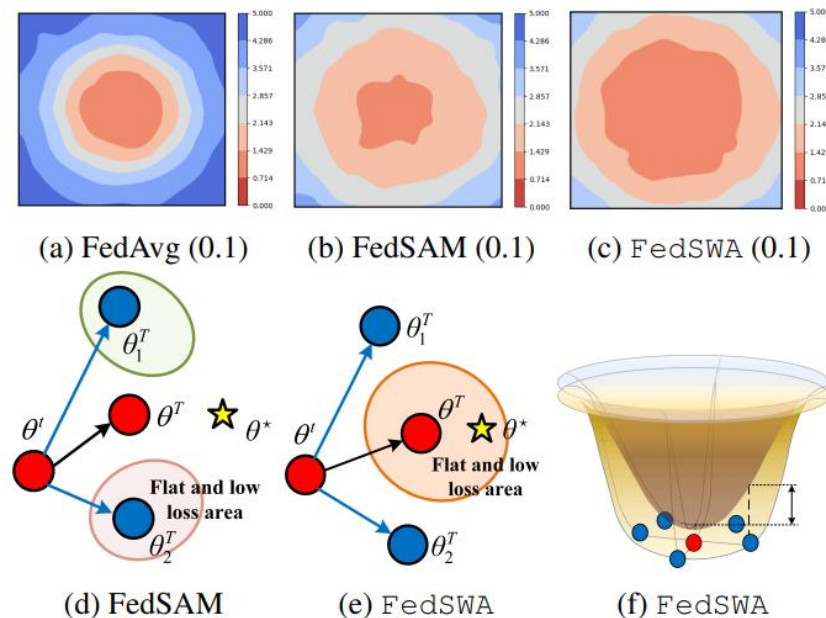
# Improving Generalization in Federated Learning with Highly Heterogeneous Data via Momentum-Based Stochastic Controlled Weight Averaging (ICML'2025)

In order to improve the generalization ability of federated learning in data heterogeneity scenarios, the FedSWA algorithm is proposed based on the SWA optimizer. Compared with the FedSAM proposed by ICML2023, our algorithm is more inclined to find the global flat minimum value.

## 研究工作的创新性构思

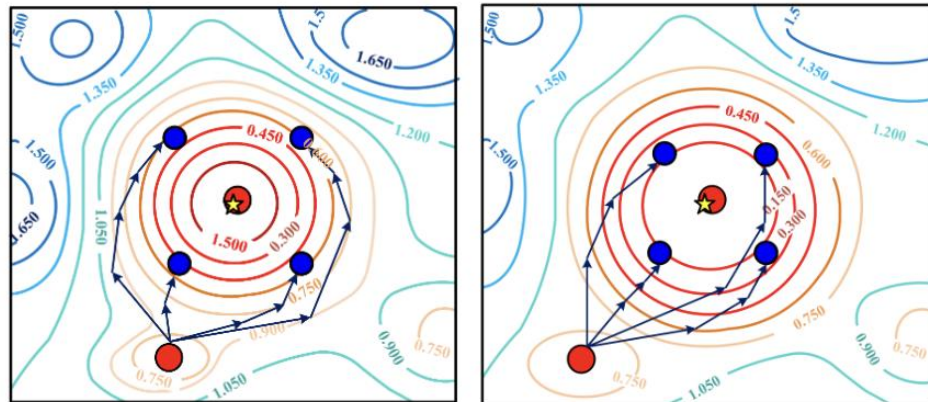
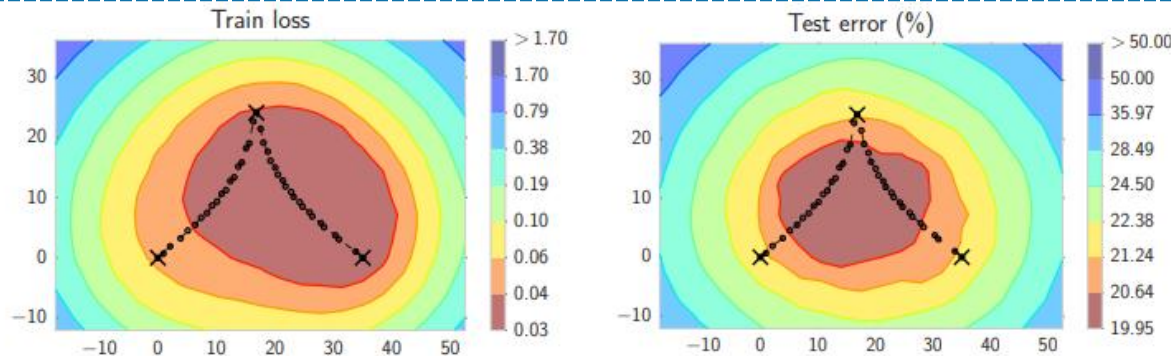


Flat minimum value has better generalization ability



The FedSWA algorithm is superior to FedSAM and flatter

SWA is inspired by a practical observation: at the end of each learning rate cycle, the model often reaches a local minimum near the edge of the loss surface, where the loss is relatively low. By averaging several of these points, we can likely obtain a more general and even lower-loss solution.



(a) FedSWA

(b) FedMoSWA

---

**Algorithm 1** Stochastic Weight Averaging

---

**Require:**

weights  $\hat{w}$ , LR bounds  $\alpha_1, \alpha_2$ ,  
cycle length  $c$  (for constant learning rate  $c = 1$ ), num-  
ber of iterations  $n$

**Ensure:**  $w_{\text{SWA}}$

$w \leftarrow \hat{w}$  {Initialize weights with  $\hat{w}$ }

$w_{\text{SWA}} \leftarrow w$

**for**  $i \leftarrow 1, 2, \dots, n$  **do**

$\alpha \leftarrow \alpha(i)$  {Calculate LR for the iteration}

$w \leftarrow w - \alpha \nabla \mathcal{L}_i(w)$  {Stochastic gradient update}

**if**  $\text{mod}(i, c) = 0$  **then**

$n_{\text{models}} \leftarrow i/c$  {Number of models}

$w_{\text{SWA}} \leftarrow \frac{w_{\text{SWA}} \cdot n_{\text{models}} + w}{n_{\text{models}} + 1}$  {Update average}

**end if**

**end for**

{Compute BatchNorm statistics for  $w_{\text{SWA}}$  weights}

---

**Theoretical Results:** We analyze the generalization ability and convergence rate of the FedSAM algorithm from ICML and our proposed FedSWA algorithm using uniform stability theory. The results show that the proposed FedSWA algorithm outperforms FedSAM.

Algorithm	Generalization error	Optimization error
FedSAM (Qu et al., 2022)	$\mathcal{O}\left(\frac{L}{mn\beta}e^{1+\frac{1}{T}}(\bar{c}L + \bar{c}\sigma_g + \bar{c}\sigma)\right)$	$\mathcal{O}\left(\frac{\beta F}{\sqrt{TK}s} + \frac{\sqrt{K}\sigma_g^2}{\sqrt{T}s} + \frac{L^2\sigma^2}{T^{3/2}K} + \frac{L^2}{T^2}\right)$
MoFedSAM (Qu et al., 2022)	$\mathcal{O}\left(\frac{L}{mn\beta}e^{1+\frac{1}{T}}(\bar{c}L + \bar{c}\sigma_g + \bar{c}\sigma)\right)$	$\mathcal{O}\left(\frac{\beta LF}{\sqrt{TK}s} + \frac{\beta\sqrt{K}\sigma_g^2}{\sqrt{T}s} + \frac{L^2\sigma^2}{T^{3/2}K} + \frac{\sqrt{K}L^2}{T^{3/2}\sqrt{s}}\right)$
FedSWA (ours)	$\mathcal{O}\left(\frac{L}{mn\beta}e^{1+\frac{1}{T}}(\tilde{c}L + \tilde{c}\sigma_g + \tilde{c}\sigma)\right)$	$\mathcal{O}\left(\frac{\beta(\sigma + \sqrt{K}\sigma_g)\sqrt{F}}{\sqrt{TK}s} + \frac{F^{2/3}(\beta\sigma_g^2)^{1/3}}{T^{2/3}} + \frac{\beta F}{T}\right)$
FedMoSWA (ours)	$\mathcal{O}\left(\frac{L}{mn\beta}e^{1+\frac{1}{T}}(\tilde{c}L + \sigma_g + \tilde{c}\sigma)\right)$	$\mathcal{O}\left(\frac{\sigma\sqrt{F}}{\sqrt{TK}s}(\sqrt{1 + \frac{s}{\alpha^2}}) + \frac{\beta F}{T}\left(\frac{m}{s}\right)^{\frac{2}{3}}\right)$

**Theorem 5.1** (Generalization Error). Assuming all clients participate in each round with Option I:

**Strongly convex:** Under Assumptions 1-5, suppose loss  $\ell(\mathbf{x}, \mathbf{y}; \theta)$  is  $\mu$ -strongly convex. By setting  $\eta_k^t \leq \frac{1}{\beta KT}$ ,  $\tilde{b} = 1 + \left(\frac{\mu}{(\beta + \mu)K}\right)^{K-1} \frac{1}{T}$ , the generalization error satisfies:

$$\text{FedSWA: } \varepsilon_{\text{gen}} \leq \frac{2L}{mn\beta}e^{1 - \frac{\mu}{(\beta + \mu)T}}(\tilde{b}L + \tilde{b}\sigma_g + \tilde{b}\sigma)$$

$$\text{FedMoSWA: } \varepsilon_{\text{gen}} \leq \frac{2L}{mn\beta}e^{1 - \frac{\mu}{(\beta + \mu)T}}(\tilde{b}L + \sigma_g + \tilde{b}\sigma)$$

**Non-convex:** Under Assumptions 2-5, assume  $\ell(\mathbf{x}, \mathbf{y}; \theta)$  is  $\beta$ -smooth. Together with  $\eta_k^t \leq \frac{1}{\beta KT}$ ,  $\tilde{c} = 1 + \left(2 + \frac{1}{KT}\right)^{K-1} \frac{1}{T}$ , the generalization error satisfies:

$$\text{FedSWA: } \varepsilon_{\text{gen}} \leq \frac{2L}{mn\beta}e^{\frac{1}{T}+1}(\tilde{c}L + \tilde{c}\sigma_g + \tilde{c}\sigma).$$

$$\text{FedMoSWA: } \varepsilon_{\text{gen}} \leq \frac{2L}{mn\beta}e^{\frac{1}{T}+1}(\tilde{c}L + \sigma_g + \tilde{c}\sigma).$$

**Theorem 5.2** (Optimization Error of FedMoSWA). For  $\beta$ -smooth functions  $\{F_i\}$ , which satisfy Assumptions 6-9, and are the same as in the SCAFFOLD (Karimireddy et al., 2020) algorithm (see the Appendix for details), the output of FedMoSWA has expected error smaller than  $\epsilon$ .

**Strongly convex:**  $\eta_k^t \leq \min\left(\frac{1}{\beta K\alpha}, \frac{s}{\mu m K\alpha}\right)$ ,  $T \geq \max\left(\frac{\beta}{\mu}, \frac{m}{s}\right)$  then

$$\mathcal{O}\left(\frac{\sigma^2}{\mu TKs}\left(1 + \frac{s}{\alpha^2}\right) + \frac{m\mu}{s}D^2 \exp\left(-\left\{\frac{s}{m} + \frac{\mu}{\beta}\right\}T\right)\right)$$

**Non-convex:**  $\eta_k^t \leq \frac{1}{K\alpha\beta}\left(\frac{s}{m}\right)^{\frac{2}{3}}$ , and  $T \geq 1$ , then

$$\mathcal{O}\left(\frac{\sigma\sqrt{F}}{\sqrt{TK}s}\left(\sqrt{1 + \frac{m}{\alpha^2}}\right) + \frac{\beta F}{T}\left(\frac{m}{s}\right)^{\frac{2}{3}}\right).$$

Here  $D^2 := \|\theta^0 - \theta^*\|^2 + \frac{1}{2m\beta^2} \sum_{i=1}^m \|c_i^0 - \nabla F_i(\theta^*)\|^2$ ,  $F := F(\theta^0) - F(\theta^*)$ .

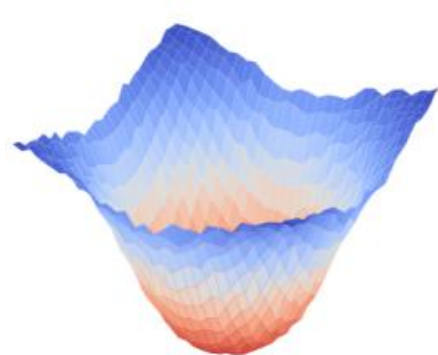
**Algorithm 1** FedSWA, FedMoSWA algorithm.

- 1: **Input:**  $\lambda, \rho$ , initial server model  $\theta_0$ , number of clients  $N$ , number of communication rounds  $T$ , number of local iterations  $K$ , local learning rate  $\eta_l$ .
- 2: **for**  $t = 0, \dots, T$  **do**
- 3:   Communicate  $(\theta_{t-1})$  to selected clients  $i \in [s]$ .
- 4:   Communicate  $(\theta_{t-1}, \mathbf{m})$  to selected clients  $i \in [s]$ .
- 5:   **for**  $i = 1, \dots, s$  clients in parallel **do**
- 6:     **for**  $k = 0, \dots, K$  local update **do**
- 7:       Compute mini-batch gradient  $g_i(\theta_{i,k}^t)$ .
- 8:        $\eta_k^t = \eta_l \left(1 - \frac{k}{K}\right) + \frac{k}{K} \rho \eta_l$ .
- 9:        $\theta_{i,k+1}^t \leftarrow \theta_{i,k}^t - \eta_k^t (g_i(\theta_{i,k}^t))$ .
- 10:       $\theta_{i,k+1}^t \leftarrow \theta_{i,k}^t - \eta_k^t (g_i(\theta_{i,k}^t) - c_i + \mathbf{m})$ .
- 11:    **end for**
- 12:    Communicate  $(\theta_{i,K}^t)$  to server.
- 13:     $c_i^+ \leftarrow$  (i)  $g_i(\mathbf{x})$  or (ii)  $c_i - \mathbf{m} + \frac{1}{\sum_k \eta_k^t} (\theta_{t-1} - \theta_{i,k}^t)$ .
- 14:    Communicate  $(\theta_{i,K}^t, c_i^+ - \mathbf{m})$  to server,  $c_i \leftarrow c_i^+$ .
- 15:    **end for**
- 16:     $\mathbf{m} \leftarrow \mathbf{m} + \gamma \frac{1}{s} \sum_{i \in [s]} \Delta c_i$ .
- 17:     $\mathbf{v}_t = \frac{1}{s} \sum_{i=1}^s \theta_{i,K}^t$ ,  $\theta_t = \theta_{t-1} + \alpha (\mathbf{v}_t - \theta_{t-1})$ .
- 18: **end for**

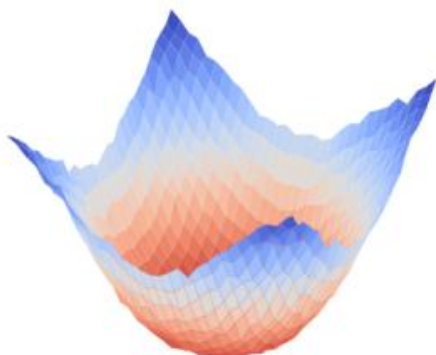
Improved accuracy by 6.3% on the CIFAR100 dataset. Improved accuracy by 1.8% on the ImageNet dataset.

Table 3: Comparison of each algorithm on the CIFAR100 and Tiny ImageNet datasets with different data heterogeneity.

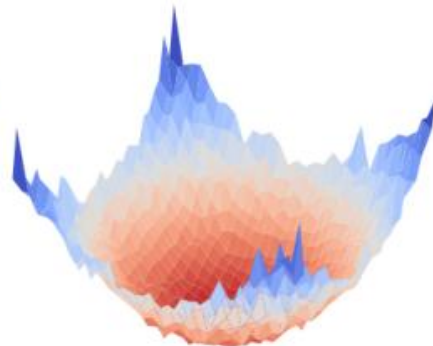
Method	CIFAR100 (ResNet-18)						Tiny ImageNet (ViT-Base)					
	Dirichlet-0.1		Dirichlet-0.3		Dirichlet-0.6		Dirichlet-0.1		Dirichlet-0.3		Dirichlet-0.6	
	Acc.(%)	Rounds	Acc.(%)	Rounds	Acc.(%)	Rounds	Acc.(%)	Rounds	Acc.(%)	Rounds	Acc.(%)	Rounds
	1000R	55%	1000R	55%	1000R	55%	400R	70%	400R	70%	400R	70%
FedAvg	45.8 $\pm$ 0.3	1000+	52.5 $\pm$ 0.3	1000+	54.2 $\pm$ 0.2	1000+	70.9 $\pm$ 0.1	258	71.8 $\pm$ 0.1	223	72.8 $\pm$ 0.1	208
FedDyn	45.8 $\pm$ 0.2	1000+	45.9 $\pm$ 0.3	1000+	46.5 $\pm$ 0.2	1000+	67.5 $\pm$ 0.3	400+	68.2 $\pm$ 0.3	400+	69.3 $\pm$ 0.3	400+
SCAFFOLD	44.3 $\pm$ 0.3	1000+	50.3 $\pm$ 0.3	1000+	52.3 $\pm$ 0.2	1000+	71.6 $\pm$ 0.1	202	72.5 $\pm$ 0.1	192	73.1 $\pm$ 0.2	169
FedSAM	40.1 $\pm$ 0.4	1000+	49.0 $\pm$ 0.3	1000+	51.9 $\pm$ 0.5	1000+	71.4 $\pm$ 0.2	212	72.2 $\pm$ 0.2	194	72.9 $\pm$ 0.4	180
MoFedSAM	51.5 $\pm$ 0.2	1000+	57.5 $\pm$ 0.2	770	60.1 $\pm$ 0.1	603	71.6 $\pm$ 0.4	229	72.4 $\pm$ 0.3	214	72.5 $\pm$ 0.4	209
FedLESAM	48.7 $\pm$ 0.2	1000+	53.3 $\pm$ 0.4	1000+	52.1 $\pm$ 0.1	1000+	71.9 $\pm$ 0.3	210	72.1 $\pm$ 0.2	188	72.5 $\pm$ 0.3	182
FedASAM	47.7 $\pm$ 0.3	1000+	46.6 $\pm$ 0.2	1000+	49.8 $\pm$ 0.1	1000+	69.2 $\pm$ 0.3	400+	71.3 $\pm$ 0.2	234	72.1 $\pm$ 0.3	196
FedACG	52.2 $\pm$ 0.4	1000+	57.7 $\pm$ 0.2	717	61.7 $\pm$ 0.4	518	66.2 $\pm$ 0.2	400+	68.5 $\pm$ 0.1	400+	70.2 $\pm$ 0.3	386
FedSWA	50.3 $\pm$ 0.3	1000+	55.5 $\pm$ 0.4	889	59.8 $\pm$ 0.3	574	71.9 $\pm$ 0.3	199	72.6 $\pm$ 0.2	179	73.2 $\pm$ 0.2	168
FedMoSWA	<b>61.9<math>\pm</math>0.5</b>	<b>577</b>	<b>66.2<math>\pm</math>0.4</b>	<b>468</b>	<b>67.9<math>\pm</math>0.4</b>	<b>330</b>	<b>73.8<math>\pm</math>0.3</b>	<b>161</b>	<b>74.4<math>\pm</math>0.3</b>	<b>152</b>	<b>74.7<math>\pm</math>0.1</b>	<b>144</b>



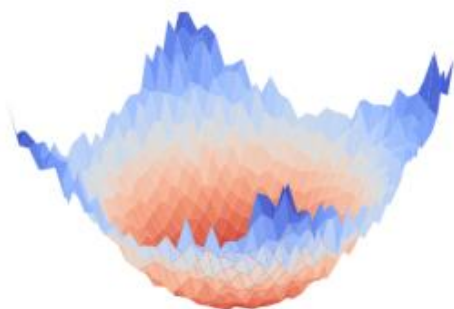
(a) SCAFFOLD



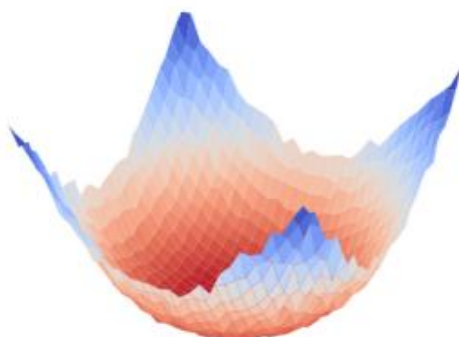
(b) FedLESAM



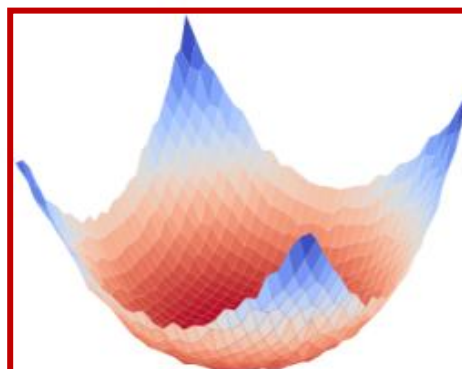
(c) FedSAM



(d) MoFedSAM



(e) FedSWA



(f) FedMoSWA

Through the visualization of model loss landscapes, it can also be observed that our algorithm has flatter loss landscapes and lower loss values.

□ **Junkang Liu**, Fanhua Shang, Yuanyuan Liu, Hongying Liu. **Improving Generalization in Federated Learning with Heterogeneous Data via Momentum-Based Stochastic Controlled Weight Averaging.** *ICML 2025*.