

Momentum-Driven Adaptivity: Towards Tuning-Free Asynchronous Federated Learning

Wenjing Yan
The Chinese University of Hong Kong

(Joint work with Xiangyu Zhong, Xiaolu Wang, Yingjun Angela Zhang)

ICML 2025

Outline

1 Background
and
Motivation



2. Algorithm
Design



3. Simulation
Results

Problem Formulation of Federated Learning (FL)

• Federated Learning:

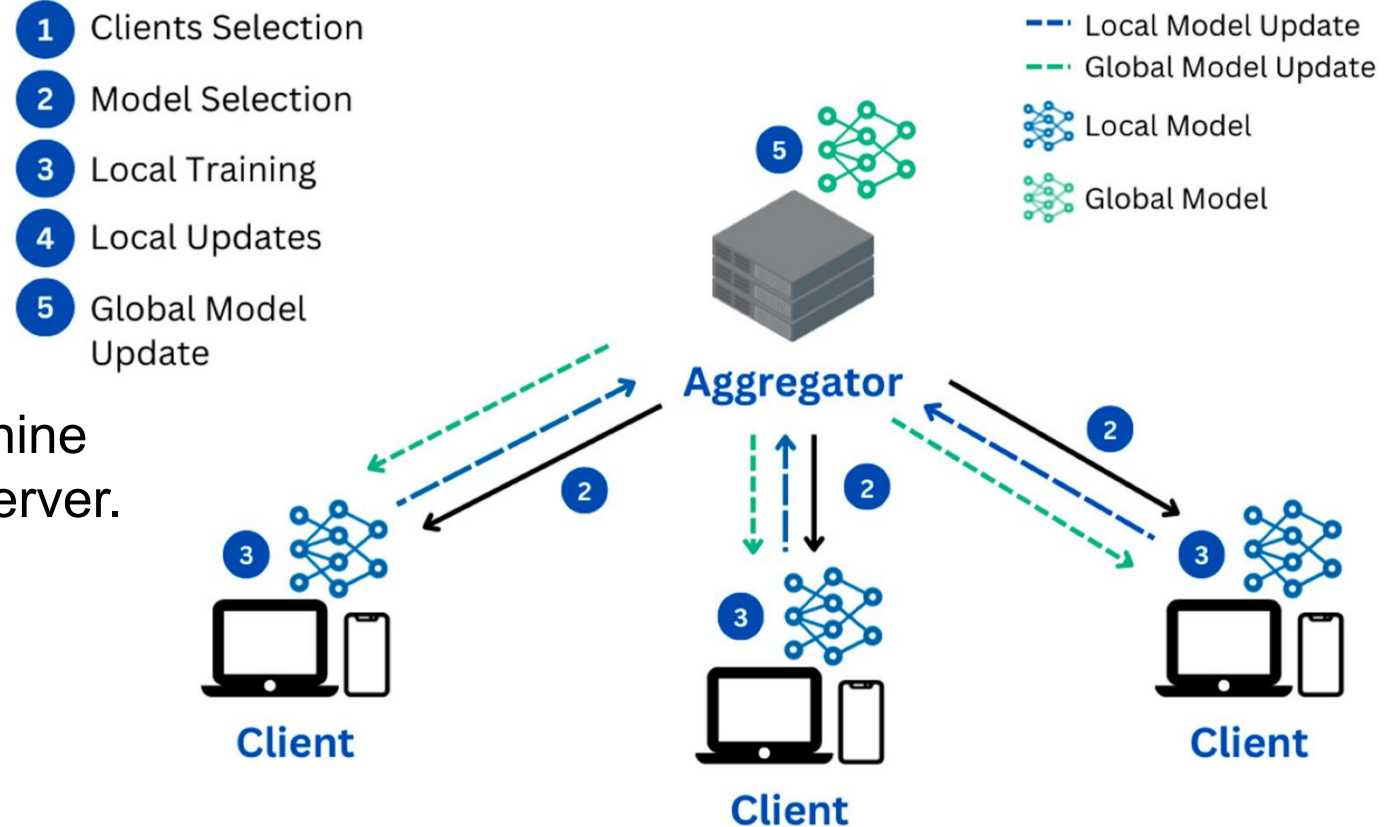
$$\min_{\theta \in \mathbb{R}^d} f(\theta) := \frac{1}{N} \sum_{i=1}^N f_i(\theta)$$

where $f_i(\theta) := \mathbb{E}_{\xi_i \sim \mathcal{D}_i} [F(\theta; \xi_i)]$

- Multiple clients collaboratively train a machine learning model with the help of a central server.
- Each client performs multiple local update based on its local private data
- Server aggregates the global model

Advantages:

- Ensures privacy by avoiding raw data sharing
- Offers scalability and communication efficiency



Challenge 1: Data Heterogeneity

- Data Heterogeneity in FL:

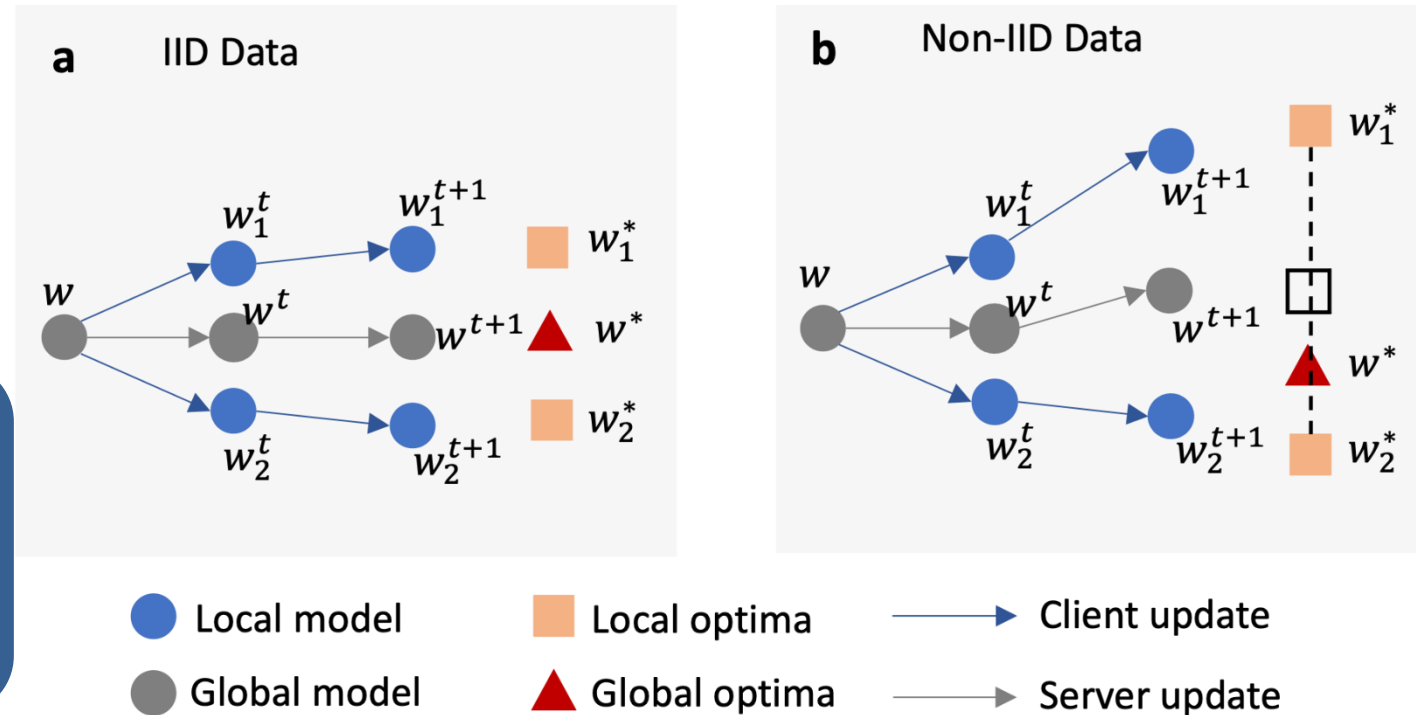
$$f_i(\theta) := \mathbb{E}_{\xi_i \sim \mathcal{D}_i} [F(\theta; \xi_i)]$$

$$\mathcal{D}_i \neq \mathcal{D}_j \text{ for any } i \neq j$$

Bounded Gradient Dissimilarity

$$\frac{1}{N} \sum_{i=1}^N \|\nabla f_i(\theta)\|^2 \leq B \|\nabla f(\theta)\|^2 + \sigma_h^2$$

“Client drift”: Local updates from individual clients diverge significantly from one another and from the global objective



- Quantifying this bound is difficult in FL due to privacy and data constraints.
- Limit the applicability of FL dynamic environments with varying data distributions.

Challenge 2: Problem-Specific Hyperparameter Tuning

• Federated Learning:

$$\min_{\theta \in \mathbb{R}^d} f(\theta) := \frac{1}{N} \sum_{i=1}^N f_i(\theta)$$

where $f_i(\theta) := \mathbb{E}_{\xi_i \sim \mathcal{D}_i} [F(\theta; \xi_i)]$

FedAvg Key Steps:

- K steps local update at client i :

$$g_i^{t,k} = \nabla F(\theta_i^{t,k}; \xi_i^{t,k})$$

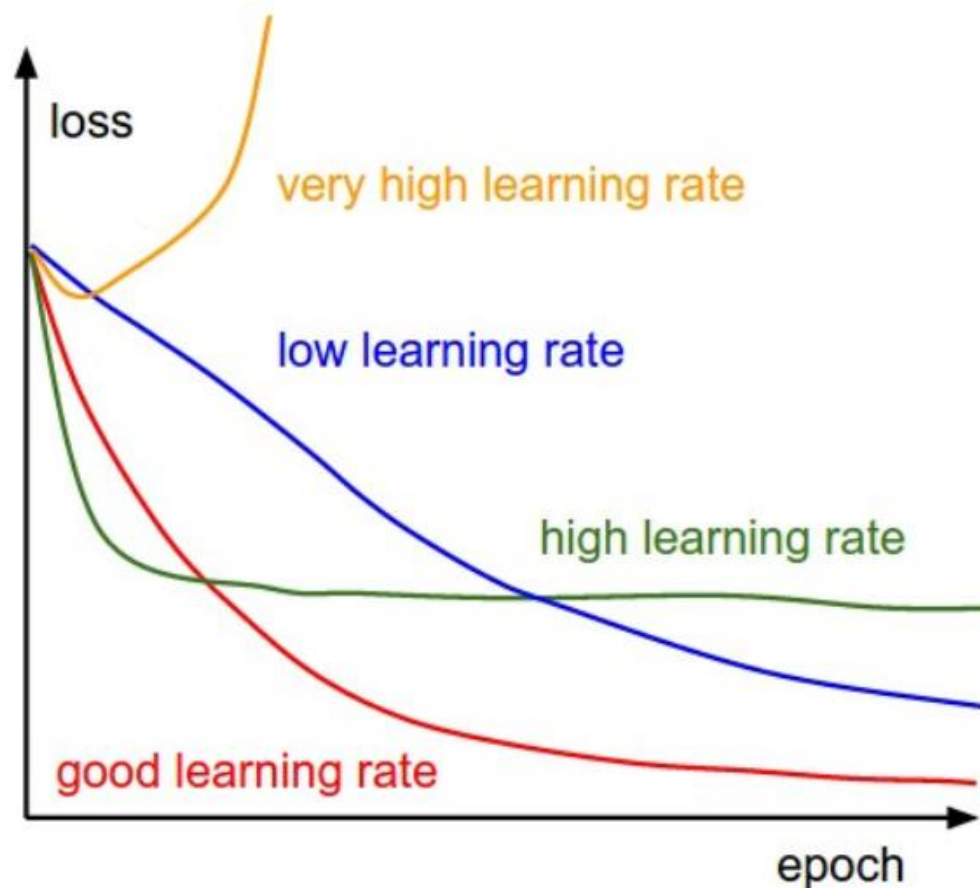
$$\theta_i^{t,k+1} = \theta_i^{t,k} - \eta_t g_i^{t,k}$$

- Global model aggregation at server:

$$g^t = \frac{1}{NK} \sum_{i=1}^N (\theta^t - \theta_i^{t,K})$$

$$\theta^{t+1} = \theta^t - \gamma g^t$$

Stepsize setting is crucial



Challenge 2: Problem-Specific Hyperparameter Tuning

- Federated Learning:

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^d} f(\boldsymbol{\theta}) := \frac{1}{N} \sum_{i=1}^N f_i(\boldsymbol{\theta})$$

where $f_i(\boldsymbol{\theta}) := \mathbb{E}_{\boldsymbol{\xi}_i \sim \mathcal{D}_i} [F(\boldsymbol{\theta}; \boldsymbol{\xi}_i)]$



- Problem-Specific Constants:

- L : Smoothness constant
- σ^2 : Stochastic gradient variance
- B, σ_h^2 : coefficients on gradient dissimilarity bound
- $\Delta := f(\boldsymbol{\theta}_0) - f^*$: Initial suboptimality gap

Assumptions

 L -Smoothness

$$\|\nabla F(\boldsymbol{\theta}; \boldsymbol{\xi}_i) - \nabla F(\boldsymbol{\delta}; \boldsymbol{\xi}_i)\| \leq L \|\boldsymbol{\theta} - \boldsymbol{\delta}\|$$

Stochastic Gradient Variance

$$\mathbb{E}_{\boldsymbol{\xi}_i} \|\nabla F(\boldsymbol{\theta}; \boldsymbol{\xi}_i) - \nabla f_i(\boldsymbol{\theta})\|^2 \leq \sigma^2$$

Bounded Gradient Dissimilarity

$$\frac{1}{N} \sum_{i=1}^N \|\nabla f_i(\boldsymbol{\theta})\|^2 \leq B \|\nabla f(\boldsymbol{\theta})\|^2 + \sigma_h^2$$

Challenge 1: Problem-Specific Hyperparameter Tuning

- Federated Learning:

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^d} f(\boldsymbol{\theta}) := \frac{1}{N} \sum_{i=1}^N f_i(\boldsymbol{\theta})$$

where $f_i(\boldsymbol{\theta})$



- Estimating those constants in FL is difficult due to data privacy restrictions and computational complexity.
- Problem-specific tuning limits the applicability of those FL approaches in dynamic environments (e.g., IoT, edge devices).

- L : Smoothness constant
- σ^2 : Stochastic gradient variance
- L, σ_h^2 : coefficients on gradient dissimilarity bound
- $\Delta := f(\boldsymbol{\theta}_0) - f^*$: Initial suboptimality gap

Assumptions

L-Smoothness

$$\|\nabla F(\boldsymbol{\theta}; \mathcal{D}) - \nabla F(\boldsymbol{\delta}; \mathcal{D})\| \leq L \|\boldsymbol{\theta} - \boldsymbol{\delta}\|$$

$$\|\nabla F(\boldsymbol{\theta}; \mathcal{D}) - \nabla f_i(\boldsymbol{\theta})\|^2 \leq \sigma^2$$

Bounded Gradient Dissimilarity

$$\frac{1}{N} \sum_{i=1}^N \|\nabla f_i(\boldsymbol{\theta})\|^2 \leq B \|\nabla f(\boldsymbol{\theta})\|^2 + \sigma_h^2$$

Our Method: Tuning-Free Asynchronuous Federated Learning (AFL)

- Federated Learning:

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^d} f(\boldsymbol{\theta}) := \frac{1}{N} \sum_{i=1}^N f_i(\boldsymbol{\theta})$$

where $f_i(\boldsymbol{\theta}) := \mathbb{E}_{\boldsymbol{\xi}_i \sim \mathcal{D}_i} [F(\boldsymbol{\theta}; \boldsymbol{\xi}_i)]$



Eliminating the requirement on data heterogeneity bounds

Independent of all problem-specific parameters, enabling **tuning-free**

Assumptions **L -Smoothness**

$$\|\nabla F(\boldsymbol{\theta}; \boldsymbol{\xi}_i) - \nabla F(\boldsymbol{\delta}; \boldsymbol{\xi}_i)\| \leq L \|\boldsymbol{\theta} - \boldsymbol{\delta}\|$$

Stochastic Gradient Variance

$$\mathbb{E}_{\boldsymbol{\xi}_i} \|\nabla F(\boldsymbol{\theta}; \boldsymbol{\xi}_i) - \nabla f_i(\boldsymbol{\theta})\|^2 \leq \sigma^2$$

Bounded Gradient Dissimilarity

$$\frac{1}{N} \sum_{i=1}^N \|\nabla f_i(\boldsymbol{\theta})\|^2 \leq B \|\nabla f(\boldsymbol{\theta})\|^2 + \sigma_h^2$$

Outline

1 Background
and
Motivation



2. Algorithm
Design



3. Simulation
Results

Handling data heterogeneity: **Momentum + Control variates**

- Local update at client i :

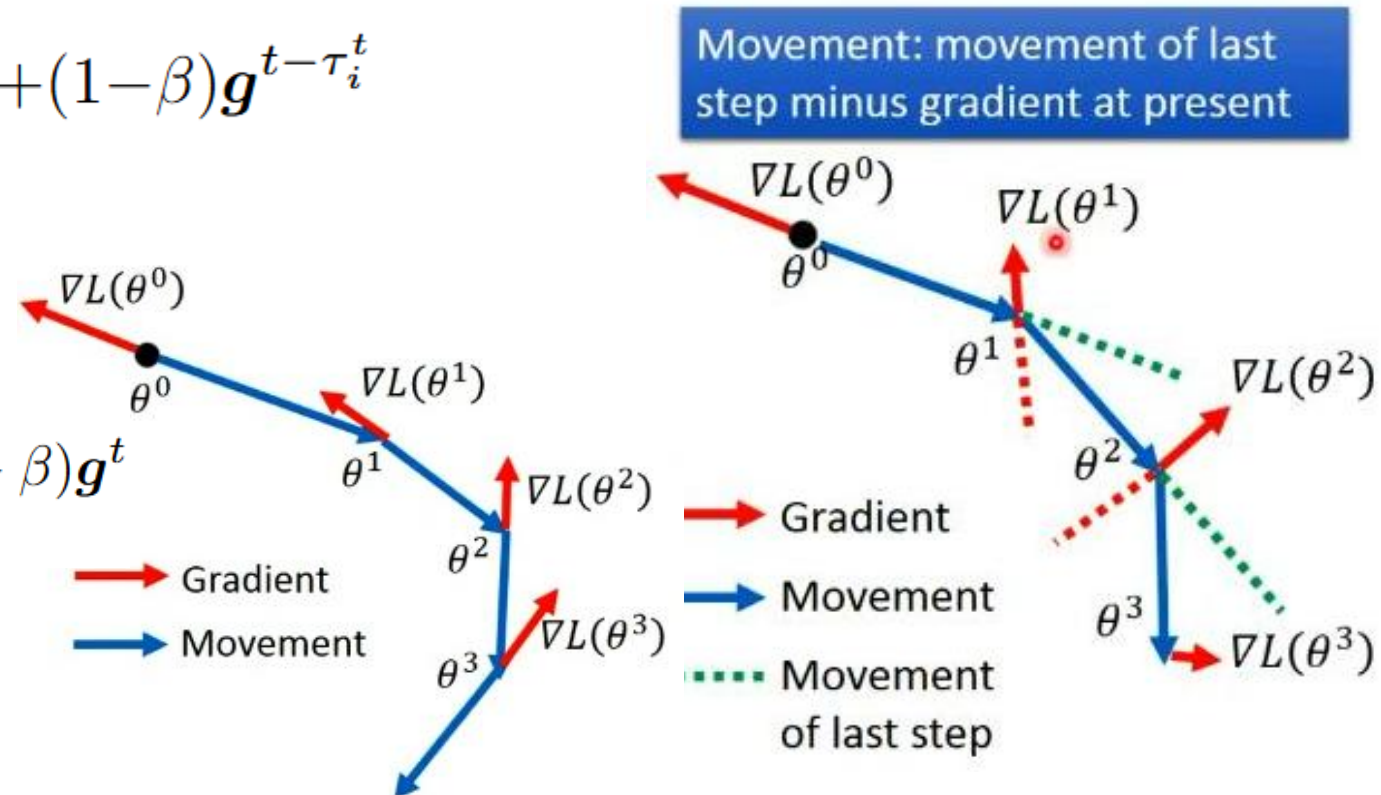
$$\mathbf{g}_i^{t,k} = \beta \left(\nabla F \left(\boldsymbol{\theta}_i^{t,k}; \boldsymbol{\xi}_i^{t,k} \right) - \tilde{\mathbf{c}}_i^t + \mathbf{c}^{t-\tau_i^t} \right) + (1-\beta) \mathbf{g}^{t-\tau_i^t}$$

$$\boldsymbol{\theta}_i^{t,k+1} = \boldsymbol{\theta}_i^{t,k} - \eta \mathbf{g}_i^{t,k}$$

- Global aggregation at server:

$$\mathbf{g}^{t+1} = \beta \left(\frac{1}{S} \sum_{i \in \mathcal{S}_t} (\mathbf{c}_i^{t+1} - \mathbf{c}_i^t) + \mathbf{c}^t \right) + (1-\beta) \mathbf{g}^t$$

$$\mathbf{c}^{t+1} = \mathbf{c}^t + \frac{1}{N} \sum_{i \in \mathcal{S}_t} (\mathbf{c}_i^{t+1} - \mathbf{c}_i^t)$$



Momentum: Accumulating past gradients across iterations and clients

Algorithm 1 MasFL: Procedures at Central Server

- 1: **Require:** Initial model θ^0 , control variates $c_i^0 = \frac{1}{K} \sum_{k=0}^{K-1} \nabla F(\theta^0; \xi_i^{-1,k})$ for any i , $c^0 = \frac{1}{N} \sum_i c_i^0$, momentum $g^0 = c^0$, global learning rate γ , local learning rate η , and momentum parameter β
- 2: **for** $t = 0, \dots, T - 1$ **do**
- 3: Randomly selected a set of clients \mathcal{S}_t
- 4: Update

$$c_i^{t+1} = \begin{cases} \tilde{c}_i^{t+1}, & \text{if } i \in \mathcal{S}_t \\ c_i^t, & \text{otherwise} \end{cases}$$

- 5: Aggregate momentum

$$g^{t+1} = \beta \left(\frac{1}{S} \sum_{i \in \mathcal{S}_t} (c_i^{t+1} - c_i^t) + c^t \right) + (1 - \beta)g^t$$

- 6: Update global model $\theta^{t+1} = \theta^t - \gamma g^{t+1}$
- 7: Aggregate control variate

$$c^{t+1} = c^t + \frac{1}{N} \sum_{i \in \mathcal{S}_t} (c_i^{t+1} - c_i^t)$$

- 8: Send θ^{t+1} and $\beta c^{t+1} + (1 - \beta)g^{t+1}$ to all clients
- 9: **end for**

Algorithm 2 MasFL: Procedures at Client i

- 1: Receive $\theta^{t-\tau_i^t}$ and $\beta c^{t-\tau_i^t} + (1 - \beta)g^{t-\tau_i^t}$ from server.
Set $\theta_i^{t,0} = \theta^{t-\tau_i^t}$
- 2: **for** $k = 0, \dots, K - 1$ **do**
- 3: Compute

$$g_i^{t,k} = \beta \left(\nabla F(\theta_i^{t,k}; \xi_i^{t,k}) - \tilde{c}_i^t + c^{t-\tau_i^t} \right) + (1 - \beta)g^{t-\tau_i^t}$$

- 4: Update local model $\theta_i^{t,k+1} = \theta_i^{t,k} - \eta g_i^{t,k}$
- 5: **end for**
- 6: Send $\tilde{c}_i^{t+1} = \frac{1}{K} \sum_{k=0}^{K-1} \nabla F(\theta_i^{t,k}; \xi_i^{t,k})$ to the server

Key Techniques: Normalized Gradient Descent

- Traditional gradient descent:

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \eta \nabla f(\mathbf{x}_t)$$

- The stepsize constraint is:

$$\text{e.g. } \eta \leq 1/L$$

The step size η must be small enough to account for the gradient's magnitude, which is governed by L . Large gradients can lead to overshooting and instability.

- Normalized gradient descent:

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \eta \frac{\nabla f(\mathbf{x}_t)}{\|\nabla f(\mathbf{x}_t)\|}$$

The gradient is normalized, so the step size η no longer depends on how large the gradient is (irrespective of L). Consistent step sizes allows NGD to navigate flat regions, steep regions, and saddle points more effectively.



Algorithm 3 AdaMasFL: Procedures at Central Server

- 1: **Require:** Initial model θ^0 , control variates $c_i^0 = \frac{1}{K} \sum_{k=0}^{K-1} \nabla F(\theta^0; \xi_i^{-1,k})$ for any i , $c^0 = \frac{1}{N} \sum_i c_i^0$, momentum $g^0 = c^0$, global learning rate γ , local learning rate η , and momentum parameter β .
- 2: **for** $t = 0, \dots, T - 1$ **do**
- 3: Randomly selected a set of clients \mathcal{S}_t
- 4: Update

$$c_i^{t+1} = \begin{cases} \tilde{c}_i^{t+1}, & \text{if } i \in \mathcal{S}_t \\ c_i^t, & \text{otherwise} \end{cases}$$

- 5: Aggregate local updates $\bar{g}^t = \frac{1}{S} \sum_{i \in \mathcal{S}_t} \Delta_i^t$
- 6: Update global model $\theta^{t+1} = \theta^t - \gamma \bar{g}^t$
- 7: Aggregate momentum

$$g^{t+1} = \beta \left(\frac{1}{S} \sum_{i \in \mathcal{S}_t} (c_i^{t+1} - c_i^t) + c^t \right) + (1 - \beta) g^t$$

- 8: Aggregate control variate

$$c^{t+1} = c^t + \frac{1}{N} \sum_{i \in \mathcal{S}_t} (c_i^{t+1} - c_i^t)$$

- 9: Download θ^{t+1} and $\beta c^{t+1} + (1 - \beta) g^{t+1}$ to all clients
- 10: **end for**

Algorithm 4 AdaMasFL: Procedures at Client i

- 1: Receive $\theta^{t-\tau_i^t}$ and $\beta c^{t-\tau_i^t} + (1 - \beta) g^{t-\tau_i^t}$ from server.
Set $\theta_i^{t,0} = \theta^{t-\tau_i^t}$
- 2: **for** $k = 0, \dots, K - 1$ **do**
- 3: Compute

$$g_i^{t,k} = \beta \left(\nabla F(\theta_i^{t,k}; \xi_i^{t,k}) - \tilde{c}_i^t + c^{t-\tau_i^t} \right) + (1 - \beta) g^{t-\tau_i^t}$$

- 4: Update local model

$$\theta_i^{t,k+1} = \theta_i^{t,k} - \eta \frac{g_i^{t,k}}{\|g_i^{t,k}\|}$$

- 5: **end for**
- 6: Send $\Delta_i^t = \frac{1}{\eta K} (\theta^{t-\tau_i^t} - \theta_i^{t,K})$ and $\tilde{c}_i^{t+1} = \frac{1}{K} \sum_{k=0}^{K-1} \nabla F(\theta_i^{t,k}; \xi_i^{t,k})$ to the server

Convergence of MasFL

Theorem 1. Suppose that Assumptions 1 and 2 holds. Let $\{\theta^t\}_{t=1}^T$ be the global iterates generated by MasFL. Set $\beta = \sqrt{\frac{SK}{T}}$ and $\gamma = \frac{1}{4L} \sqrt{\frac{SK}{T}}$. Define $a := \tau_{\max}^2 \beta^2 + 20e^2 \eta^2 K^2 L^2$. If the condition $1 - 4a - \sqrt{\frac{SK}{T}} \geq 0$ is satisfied, i.e., $\eta \leq \frac{\sqrt{T - \sqrt{SKT} - 4SK\tau_{\max}^2}}{4\sqrt{5}eKL\sqrt{T}}$, then it holds that

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} \|\nabla f(\theta^t)\|^2 \leq \mathcal{O} \left(\frac{L\Delta + \sigma^2}{\sqrt{SKT}} + \frac{\sigma^2}{T} + \frac{L^2 \sigma^2}{NKT} \right)$$

for sufficiently large T , where e denotes Euler's number and $\Delta := f(\theta^0) - f^*$ represents the initial optimality gap with $f^* = \min_{\theta} f(\theta) > -\infty$.

Assumptions

 L -Smoothness

$$\|\nabla F(\theta; \xi_i) - \nabla F(\delta; \xi_i)\| \leq L \|\theta - \delta\|$$

Stochastic Gradient Variance

$$\mathbb{E}_{\xi_i} \|\nabla F(\theta; \xi_i) - \nabla f_i(\theta)\|^2 \leq \sigma^2$$

Convergence of AdaMasFL

Theorem 2. Suppose that Assumptions 1 and 2 holds. Let $\{\theta^t\}_{t=1}^T$ be the global iterates generated by AdaMasFL. Set $\gamma = \frac{(SK)^{1/4}}{T^{3/4}}$, $\eta = \frac{1}{K\sqrt{T}}$, and $\beta = \sqrt{\frac{SK}{T}}$, then it holds that

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} \|\nabla f(\theta^t)\| \leq \mathcal{O} \left(\frac{\Delta + L + \sigma + \sqrt{L\sigma}}{(SKT)^{\frac{1}{4}}} + \sigma \sqrt{\frac{SK}{T}} + \frac{\sqrt{L\sigma\tau_{\max}}}{T^{\frac{3}{8}}(SK)^{\frac{1}{8}}} \right. \\ \left. + \bar{\tau}\sigma\sqrt{\frac{S}{T}} + \tau_{\max}L\frac{(SK)^{\frac{1}{4}}}{T^{\frac{3}{4}}} + \bar{\tau}\sqrt{L\sigma\tau_{\max}}\frac{(SK)^{\frac{3}{8}}}{T^{\frac{7}{8}}} \right)$$

for sufficiently large T .

Tuning-Free AFL

All hyperparameters (η : local learning rate, γ : global learning rate, and β : momentum parameter) in Algorithm AdaMasFL are **explicated determined by system-predefined constants**: S (the number of participation clients), K (local update times), T (iteration times)

Table 1: Comparisons of AFL algorithms for handling heterogeneous data.

(Convergence Rate = The convergence rate of different algorithms in terms of $\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} \|\nabla f(\theta^t)\|$; Additional Assumptions = Additional assumptions aside from Assumptions 2.1–2.3; BDH = Bounded data heterogeneity define in (1); BG = Bounded gradient that $\|\nabla f_i(\theta)\| \leq G, \forall i, \theta$.)

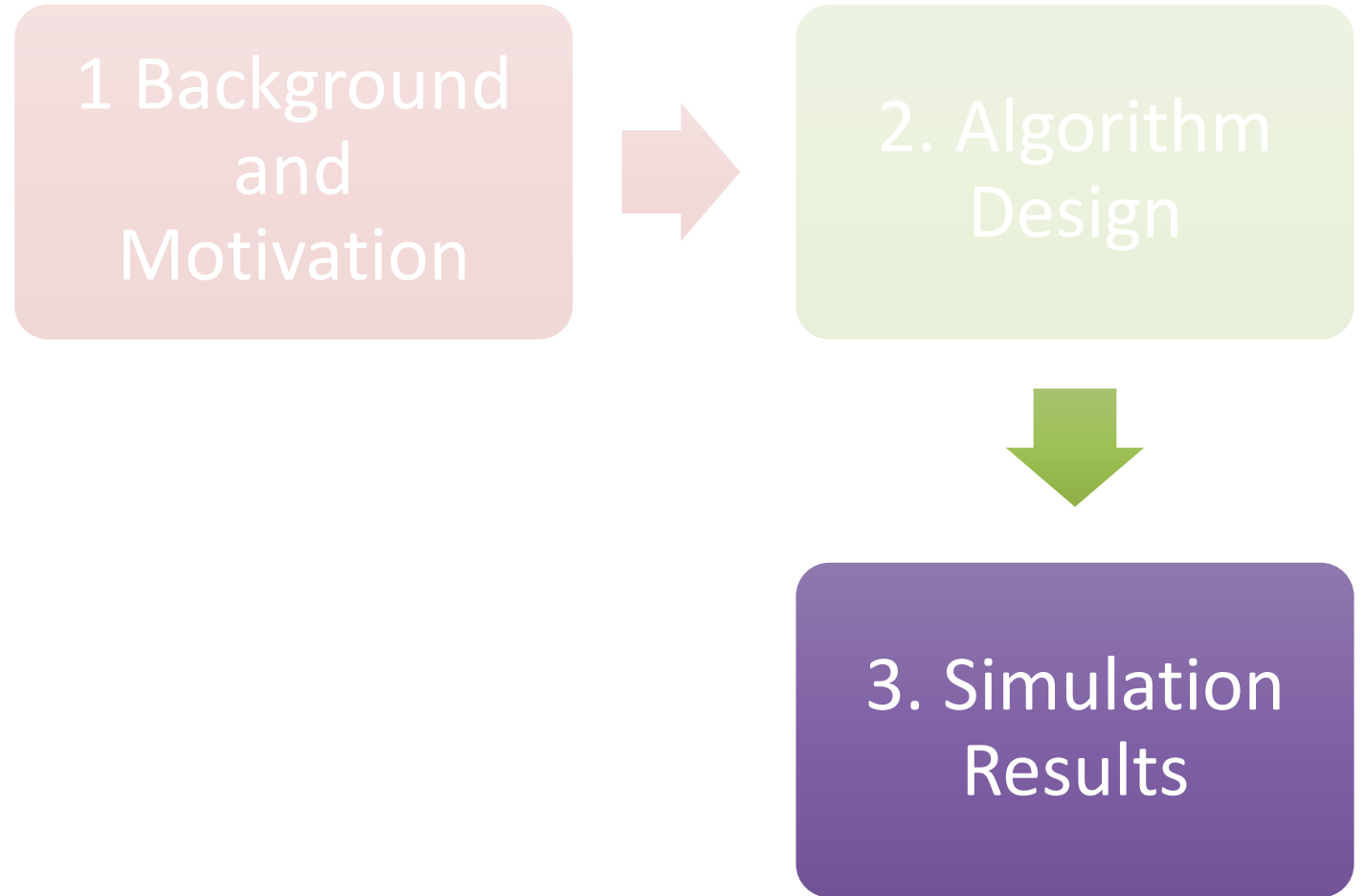
Algorithms	Convergence Rate ¹	Additional Assumptions	Stepsize Restrictions	Stepsize-related Problem-Parameters
FedBuff (Nguyen et al., 2022)	$\mathcal{O}\left(\frac{K^{1/4}\sigma_g}{(ST)^{1/4}} + \sqrt{\frac{\tau_{\max}\bar{\tau}}{T}}\right)$	BDH	$\eta\gamma \leq \frac{1}{4K\tau_{\max}^{3/2}}$	τ_{\max}
CA ² FL (Wang et al., 2023)	$\mathcal{O}\left(\frac{1}{(SKT)^{1/4}} + \sqrt{\frac{\tau_{\max}}{T}}\right)$	BDH	$\eta\gamma \leq \frac{S}{36K\tau_{\max}^2L^2}, \eta \leq \frac{1}{36K\sqrt{\tau_{\max}L}}$	τ_{\max}, L
DeFedAvg-nIID (Wang et al., 2024a)	$\mathcal{O}\left(\frac{1}{(SKT)^{1/4}} + \frac{1}{\sqrt{KT}}\right)$	BDH, BG	$\eta\gamma \leq \frac{1}{4LK\tau_{\max}}, \eta \leq \frac{1}{4\sqrt{3}LK}$	τ_{\max}, L
FADAS (Wang et al., 2024c)	$\mathcal{O}\left(\frac{1}{(ST)^{1/4}} + \sqrt{\frac{\tau_{\max}\bar{\tau}}{T}}\right)$	BDH, BG	$\eta\gamma \leq \min\left\{\frac{\epsilon^2 S(N-1)}{180C_G^2 N(N-S)\tau_{\max}^2 KL}, \frac{\sqrt{\epsilon^3 S(N-1)}}{12\sqrt{C_G N(N-S)\tau_{\max}^2 KL}}\right\}, \eta \leq \frac{\sqrt{\epsilon}}{\sqrt{360C_G\tau_{\max}^2 KL}}^2$	τ_{\max}, L, G
MasFL	$\mathcal{O}\left(\frac{\sqrt{\kappa}}{(SKT)^{1/4}} + \sqrt{\frac{\kappa}{T}}\right)^3$	–	$\beta = \sqrt{\frac{SK}{T}}, \gamma = \frac{\beta}{4L}, \eta \leq \frac{\sqrt{T - \sqrt{SKT} - 4SK\tau_{\max}^2}}{4\sqrt{5eKL}\sqrt{T}}$	τ_{\max}, L
AdaMasFL	$\mathcal{O}\left(\frac{1}{(SKT)^{1/4}} + \frac{\sqrt{\tau_{\max}}}{T^{3/8}}\right)$	–	$\beta = \sqrt{\frac{SK}{T}}, \gamma = \frac{(SK)^{1/4}}{T^{3/4}}, \eta = \frac{1}{K\sqrt{T}}$	–

¹ For the convergence rate defined in terms of $\frac{1}{T} \sum_{t=0}^{T-1} \|\nabla f(\theta^t)\|^2$, we can readily obtained the corresponding rate with respect to $\frac{1}{T} \sum_{t=0}^{T-1} \|\nabla f(\theta^t)\|$ by taking square root on both sides of the associated bound. This operator is verified by the following fact: $\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} \|\nabla f(\theta^t)\| = \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} \sqrt{\|\nabla f(\theta^t)\|^2} \leq \frac{1}{T} \sum_{t=0}^{T-1} \sqrt{\mathbb{E} \|\nabla f(\theta^t)\|^2} \leq \sqrt{\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} \|\nabla f(\theta^t)\|^2}$, where the first and second inequalities utilizes Jensen's inequality as the square root function is concave.

² $C_G := \eta KG + \epsilon$ and $\epsilon > 0$ is an adaptive optimization parameter.

³ $\kappa := \frac{3-4a}{1-4a-\sqrt{SK/T}}$ and $a := \tau_{\max}^2\beta^2 + 20e^2\eta^2K^2L^2$.

Outline



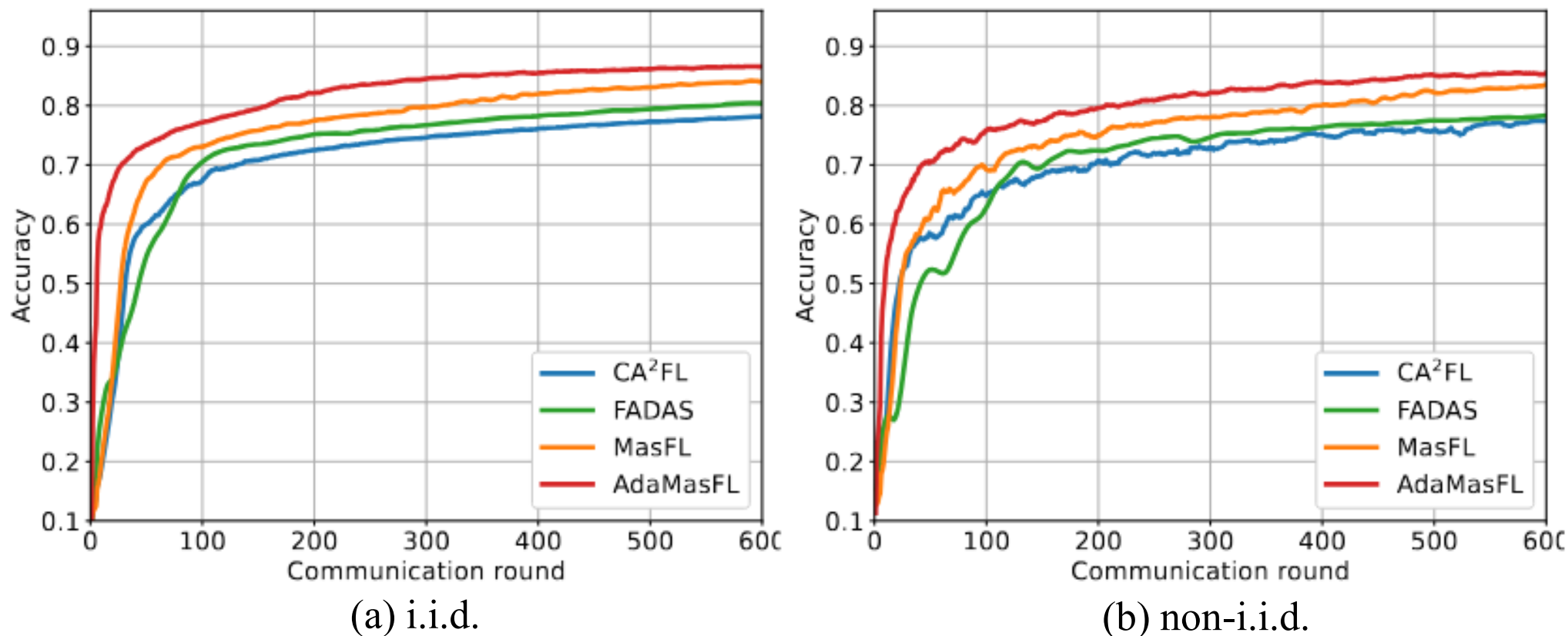


Figure 1. Test accuracy versus the number of communication rounds on the FMNIST dataset (CNN, Dir(0.5)).

Our approach achieves start-of-the-art performance while eliminating the tedious process of stepsize tuning

The stepsizes of AdaMasFL are set directly based on the guidance of Theorem 2. The stepsize of all baselines are perfectly tuned by grid search.

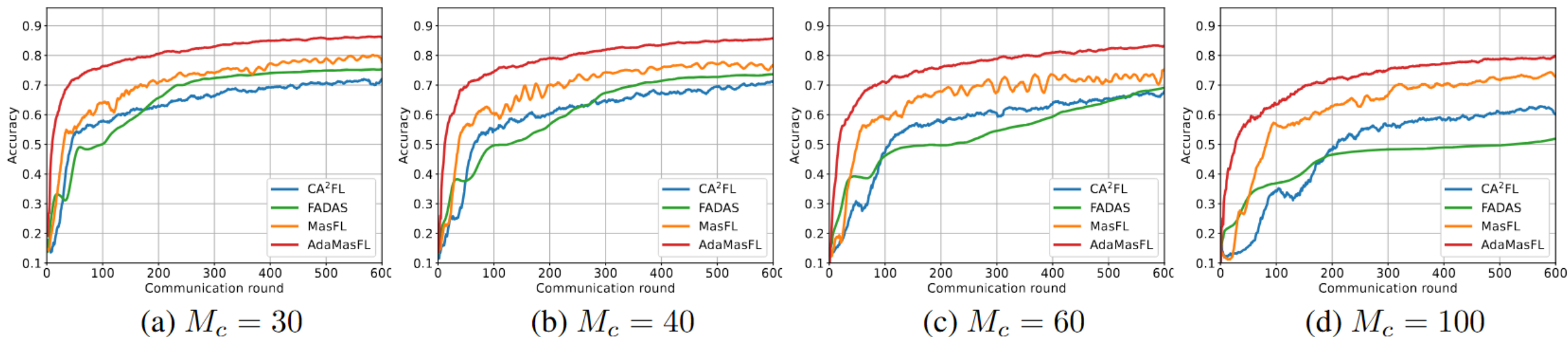


Figure 2. Test accuracy on the non-i.i.d. FMNIST dataset under varying levels of asynchrony

- M_c denotes the number of clients performing local updates concurrently. a larger M_c enables more frequent global aggregations, resulting in greater asynchronous delays.
- In Figure 2, we fix the learning rate settings of all algorithms to those used in Figure 1 and evaluate their performance as M_c increases.

AdaMasFL demonstrates exceptional robustness to varying levels of asynchrony, maintaining nearly consistent performance despite increasing delays

THANKS

Thanks!