
Appendix - Fed-ADE: Adaptive Learning Rate for Federated Post-adaptation under Distribution Shift

Contents

6	A Notation Table	2
7	B Proof of Unbiased Risk Estimation	2
8	C Proof of Theorem 1	3
9	D Proof of Theorem 2	4
10	E Proof of Theorem 3	5
11	E.1. Regret Bound	5
12	E.2. Min-max Optimality	6
13	F. Convergence Analyses	7
14	G Simulation Settings	11
15	G.1. Datasets	11
16	G.2. Online Distribution Shift Modeling	11
17	G.3. Distribution Shift Scenarios	11
18	G.4. Compared Methods	12
19	G.5. Hyperparameters	13
20	G.6. Hardware Specifications	13
21	H Extra Simulation Results	14
22	H.1. Impact of Learning Rate Selection	14
23	H.2. Impact of Distribution of Pre-training Data	15
24	H.3. Why Cosine Similarity? Alternatives and Ablations	15

Table 5. Notation Table

Notation	Description	Notation	Description
\mathbf{x}_G^0	Pre-training data at server	\mathbf{y}_G^0	Pre-training labels at server
$\mathbf{Q}_{G,\mathbf{x}}^0$	Feature distribution of \mathbf{x}_G^0	$\mathbf{Q}_{G,\mathbf{y}}^0$	Label distribution of \mathbf{y}_G^0
\mathcal{I}	Set of class indices	i	Index for class
t	Timestep	T	Total number of timesteps
\mathcal{C}	Set of FL clients	c	Client index
θ	Global model parameters	θ_c	Model parameters of client c
\mathbf{Q}_c^t	Overall data distribution at time t	$\omega(t)$	Weighting function controlling the distribution at time t
$\mathcal{H}(\cdot)$	Softmax prediction of the model	$\mathcal{L}(\cdot)$	Loss function
$\mathcal{F}_c^{t,i}(\theta_c)$	Class-wise risk for class i	$\widehat{\mathcal{F}}_c^{0,i}(\theta_c)$	Empirical risk
\mathbf{M}	Confusion matrix computed using pre-training data	$\mathbf{Q}_{c,\widehat{\mathbf{y}}}^t$	Predicted distribution over \mathbf{x}_c^t
R	Total number of communication rounds	r	Communication round
$\mathcal{C}^{(r)}$	Selected client set in round r	N_c^t	Number of samples of client c at t
ψ_c	Shared layers of client c	ϕ_c	Personalized layers of client c
$\widehat{\psi}$	Aggregated global shared layers	$\{\psi_c, \phi_c\}$	Model split: shared/personalized
$[\eta_{\min}, \eta_{\max}]$	Min/max learning rate bounds	S^t	Distribution dynamics signal
$S_{\text{unc},c}^t$	Uncertainty dynamic signal	$S_{\text{rep},c}^t$	Representation dynamic signal
\mathbf{q}_c^t	Aggregated predictive distribution of current data batch	\mathbf{z}_c^t	Batch-level latent feature vector
$h_{\psi_c}(\cdot)$	Feature extractor	$K_{\cos}, K_h, K_\psi, K_\phi$	Lipschitz constant
B	Loss value upper bound	σ	Min. singular value of \mathbf{M}
Γ	Projection constant	$ \mathcal{I} $	Number of classes
Reg_T	Dynamic regret over T rounds	$\Delta\mathcal{L}^{(0)}$	Initial loss gap for convergence

26 **B. Proof of Unbiased Risk Estimation**

Lemma 5 (Unbiased Risk Estimator). *Given model parameters θ_c independent of the test-time data $(\mathbf{x}_c^t, \mathbf{y}_c^t)$, the estimator $\widehat{\mathcal{F}}_c^t$ in (4) satisfies*

$$\mathbb{E}_{\mathbf{y}_c^t \sim \mathbf{Q}_{c,\mathbf{y}}^t} [\widehat{\mathcal{F}}_c^t(\theta_c)] = \mathcal{F}_c^t(\theta_c),$$

27 *provided that the confusion matrix \mathbf{M} , estimated at the server using sufficient pre-training data $(\mathbf{x}_G^0, \mathbf{y}_G^0)$ ensuring $\widehat{\mathbf{M}} = \mathbf{M}$,*
 28 *and that each client has sufficient local initial data $(\mathbf{x}_c^0, \mathbf{y}_c^0)$ such that $\widehat{\mathcal{F}}_c^0(\theta_c) = \mathcal{F}_c^0(\theta_c)$ for all $i \in \mathcal{I}$.*

Proof. The expected risk decomposes into class-specific components:

$$\mathcal{F}_c^t(\theta_c) = \sum_{i \in \mathcal{I}} [\mathbf{Q}_{\mathbf{y}_c^t}]_i \cdot \mathcal{F}_c^{0,i}(\theta_c)$$

29 where \mathcal{I} denotes the class index set. Through BBSE [40], we approximate $\mathbf{Q}_{\mathbf{y}_c^t} \approx \mathbf{M}^{-1} \mathbf{Q}_{\widehat{\mathbf{y}}_c^t}$ using the confusion matrix \mathbf{M}
 30 from pre-training and empirical predictions $\mathbf{Q}_{\widehat{\mathbf{y}}_c^t}$.

Expanding the estimator's expectation:

$$\mathbb{E}[\widehat{\mathcal{F}}_c^t] = \sum_{i \in \mathcal{I}} [\widehat{\mathbf{M}}^{-1} \mathbb{E}[\mathbf{Q}_{\widehat{\mathbf{y}}_c^t}]]_i \cdot \widehat{\mathcal{F}}_c^{0,i}(\theta_c)$$

Under condition $\mathbb{E}[\mathbf{Q}_{\widehat{\mathbf{y}}_c^t}] = \mathbf{M} \mathbf{Q}_{\mathbf{y}_c^t}$, yielding:

$$\mathbb{E}[\widehat{\mathcal{F}}_c^t] = \sum_{i \in \mathcal{I}} [\mathbf{Q}_{\mathbf{y}_c^t}]_i \cdot \mathcal{F}_c^{0,i}(\theta_c) = \mathcal{F}_c^t(\theta_c)$$

31 The estimation error $|\mathbb{E}[\widehat{\mathcal{F}}_c^t] - \mathcal{F}_c^t|$ is bounded by $\mathcal{O}(1/\sqrt{n_0})$ via concentration inequalities, where a number of pre-training
 32 data $n_0 = |(\mathbf{x}_G^0, \mathbf{y}_G^0)|$. This becomes negligible when $n_0 \geq \Omega(|\mathcal{I}|^2/\epsilon^2)$ for error tolerance ϵ . Practical implementations may
 33 employ singular value thresholding for numerical stability when inverting \mathbf{M} . \square

Complete derivation expanding the expectation operator:

$$\mathbb{E}_{\mathbf{x}_c^t}[\widehat{\mathcal{F}}_c^t] = \sum_{i \in \mathcal{I}} \mathbf{M}^{-1} \mathbb{E}[\mathbf{Q}_{\widehat{\mathbf{y}}_c^t}] \cdot \mathcal{F}_c^{0,i}(\theta_c) = \sum_{i \in \mathcal{I}} [\mathbf{Q}_{\mathbf{y}_c^t}]_i \cdot \mathcal{F}_c^{0,i}(\theta_c)$$

34 confirming the unbiasedness through the relationship $\mathbf{Q}_{\mathbf{y}_c^t} = \mathbf{M}^{-1} \mathbb{E}[\mathbf{Q}_{\widehat{\mathbf{y}}_c^t}]$ established via BBSE methodology.

35 **Remark 6.** The complexity function $\Omega(\cdot)$ in the sample complexity $n_0 \geq \Omega(|Z|^2/\epsilon^2)$ denotes an asymptotic lower bound,
 36 meaning the pre-training dataset size must grow at least quadratically with the number of classes to guarantee estimation
 37 accuracy. This aligns with information-theoretic limits for distribution estimation.

38 C. Proof of Theorem 1

We restate the theorem for completeness. Let $\mathbf{Q}_{c,y}^t$ denote the true predictive distribution marginal at timestep t , and \mathbf{q}_c^t the empirical mean softmax vector defined in (11). Under ϵ -calibration with respect to $\mathbf{Q}_{c,y}^t$, we have $\|\mathbf{q}_c^t - \mathbf{Q}_{c,y}^t\|_2 \leq \epsilon_t$ in expectation. The cumulative surrogate of predictive dynamics is

$$\bar{\mathcal{S}}_{\text{unc}} = \sum_{t=1}^T \mathcal{S}_{\text{unc}}^t = \sum_{t=1}^T (1 - \cos(\mathbf{q}_c^{t-1}, \mathbf{q}_c^t)).$$

39 Our goal is to show that $\bar{\mathcal{S}}_{\text{unc}}$ approximates the cumulative true deviation $\sum_{t=1}^T \|\mathbf{Q}_{c,y}^t - \mathbf{Q}_{c,y}^{t-1}\|_1$ up to an additive error of
 40 order $\mathcal{O}(\sum_t \epsilon_t)$.

Step 1: Lipschitz continuity of cosine distance. The cosine similarity function $\cos(\cdot, \cdot)$ is K_{\cos} -Lipschitz continuous in each of its arguments on the unit sphere, i.e.,

$$|\cos(\mathbf{a}, \mathbf{b}) - \cos(\mathbf{a}', \mathbf{b}')| \leq K_{\cos}(\|\mathbf{a} - \mathbf{a}'\|_2 + \|\mathbf{b} - \mathbf{b}'\|_2).$$

41 Since both \mathbf{q}_c^t and $\mathbf{Q}_{c,y}^t$ are probability vectors normalized to unit norm, this Lipschitz property holds directly.

Step 2: Bounding the surrogate deviation. By the definition of $\mathcal{S}_{\text{unc}}^t$, we have

$$\begin{aligned} & |\mathcal{S}_{\text{unc}}^t - (1 - \cos(\mathbf{Q}_{c,y}^{t-1}, \mathbf{Q}_{c,y}^t))| \\ & \leq K_{\cos}(\|\mathbf{q}_c^{t-1} - \mathbf{Q}_{c,y}^{t-1}\|_2 + \|\mathbf{q}_c^t - \mathbf{Q}_{c,y}^t\|_2) \leq K_{\cos}(\epsilon_{t-1} + \epsilon_t). \end{aligned}$$

Step 3: Summing over timesteps. Summing over $t = 2, \dots, T$ yields

$$\left| \sum_{t=1}^T \mathcal{S}_{\text{unc}}^t - \sum_{t=1}^T (1 - \cos(\mathbf{Q}_{c,y}^{t-1}, \mathbf{Q}_{c,y}^t)) \right| \leq K_{\cos} \sum_{t=2}^T (\epsilon_t + \epsilon_{t-1}).$$

Step 4: Relating cosine distance to total variation. For normalized probability vectors, the cosine distance upper bounds the ℓ_1 deviation:

$$1 - \cos(\mathbf{Q}_{c,y}^{t-1}, \mathbf{Q}_{c,y}^t) \leq \|\mathbf{Q}_{c,y}^{t-1} - \mathbf{Q}_{c,y}^t\|_1.$$

Combining this with the bound above yields

$$\left| \bar{\mathcal{S}}_{\text{unc}} - \sum_{t=1}^T \|\mathbf{Q}_{c,y}^{t-1} - \mathbf{Q}_{c,y}^t\|_1 \right| \leq K_{\cos} \sum_{t=2}^T (\epsilon_t + \epsilon_{t-1}).$$

43 **D. Proof of Theorem 2**

44 Let \mathbf{z}_c^t denote the empirical ℓ_2 -normalized batch-mean feature vector defined in (13), and $\bar{\mathbf{z}}_c^t = \mathbb{E}[h_{\psi_c}(x)/\|h_{\psi_c}(x)\|_2]$ be its
 45 expectation under the local distribution at time t . We assume $\|\mathbf{z}_c^t - \bar{\mathbf{z}}_c^t\|_2 \leq \epsilon'_t$ in expectation, and that the shared representation
 46 extractor h_{ψ_c} is K_h -Lipschitz.

The cumulative surrogate of representation dynamics is defined as

$$\bar{\mathcal{S}}_{\text{rep}} = \sum_{t=1}^T \mathcal{S}_{\text{rep}}^t = \frac{1}{2} \sum_{t=1}^T (1 - \cos(\mathbf{z}_c^{t-1}, \mathbf{z}_c^t)).$$

47 Since all feature vectors are ℓ_2 -normalized, the cosine distance $\frac{1}{2}(1 - \cos(\mathbf{a}, \mathbf{b}))$ is smooth on the unit sphere, and its deviation
 48 with respect to small perturbations of \mathbf{a} and \mathbf{b} is bounded linearly by their ℓ_2 distance. Using the Lipschitz continuity of h_{ψ_c} ,
 49 the deviation between the empirical mean \mathbf{z}_c^t and its population counterpart $\bar{\mathbf{z}}_c^t$ thus induces a bounded perturbation in the
 50 cosine term.

For each timestep t , we have

$$\left| \mathcal{S}_{\text{rep}}^t - \frac{1}{2}(1 - \cos(\bar{\mathbf{z}}_c^{t-1}, \bar{\mathbf{z}}_c^t)) \right| \leq K_h(\epsilon'_{t-1} + \epsilon'_t).$$

Summing over all $t = 2, \dots, T$ gives

$$\left| \sum_{t=1}^T \mathcal{S}_{\text{rep}}^t - \sum_{t=1}^T \frac{1}{2}(1 - \cos(\bar{\mathbf{z}}_c^{t-1}, \bar{\mathbf{z}}_c^t)) \right| \leq K_h \sum_{t=2}^T (\epsilon'_t + \epsilon'_{t-1}).$$

Therefore, the cumulative surrogate $\bar{\mathcal{S}}_{\text{rep}}$ satisfies

$$\left| \bar{\mathcal{S}}_{\text{rep}} - \sum_{t=1}^T \frac{1}{2} \left(1 - \frac{\langle \bar{\mathbf{z}}_c^{t-1}, \bar{\mathbf{z}}_c^t \rangle}{\|\bar{\mathbf{z}}_c^{t-1}\|_2 \|\bar{\mathbf{z}}_c^t\|_2} \right) \right| \leq K_h \sum_{t=2}^T (\epsilon'_t + \epsilon'_{t-1}),$$

51 which shows that $\bar{\mathcal{S}}_{\text{rep}}$ accurately approximates the temporal trajectory of the expected local feature representations, with error
 52 bounded by the feature-level estimation noise scaled by the Lipschitz constant K_h . \square

53 E. Proof of Theorem 3

54 E.1. Regret Bound

55 We present the proof of Theorem 3 by decomposing the dynamic regret into two terms. Our approach introduces a reference
 56 sequence that evolves in a piecewise-constant manner across predefined intervals. The first term evaluates the algorithm's
 57 performance relative to this sequence, while the second term quantifies how well the reference sequence itself approximates
 58 the optimal comparator. For the following proof, we use $\{\psi_c^t, \phi_c^t\}$ for model parameters to differentiate the model per timestep.

Proof. The regret bound can be split into two terms by introducing a reference sequence $\{\psi_c^t, \phi_c^t\}$ that only changes every τ steps. Specifically, let $\mathcal{J}_m = [(m-1)\tau + 1, m\tau]$ denote the m -th interval. For each interval, the comparator $\{\psi_c^t, \phi_c^t\}$ is chosen as the best fixed decision in that interval, i.e., $\{\psi_c^{\mathcal{J}_m}, \phi_c^{\mathcal{J}_m}\} = \arg \min_{\{\psi_c^t, \phi_c^t\}} \sum_{t \in \mathcal{J}_m} \mathcal{F}_c^t(\{\psi_c^t, \phi_c^t\})$ for $t \in \mathcal{J}_m$. Then,

$$\begin{aligned} & \mathbb{E}_{1:T} \left[\sum_{t=1}^T \mathcal{F}_c^t(\{\psi_c^t, \phi_c^t\}) - \sum_{t=1}^T \mathcal{F}_c^t(\{\psi_c^*, \phi_c^*\}) \right] \\ &= \mathbb{E}_{1:T} \left[\sum_{t=1}^T \mathcal{F}_c^t(\{\psi_c^t, \phi_c^t\}) - \sum_{m=1}^M \sum_{t \in \mathcal{J}_m} \mathcal{F}_c^t(\{\psi_c^{\mathcal{J}_m}, \phi_c^{\mathcal{J}_m}\}) \right] \end{aligned} \quad (a)$$

$$+ \mathbb{E}_{1:T} \left[\sum_{m=1}^M \sum_{t \in \mathcal{J}_m} \mathcal{F}_c^t(\{\psi_c^{\mathcal{J}_m}, \phi_c^{\mathcal{J}_m}\}) - \sum_{t=1}^T \mathcal{F}_c^t(\{\psi_c^*, \phi_c^*\}) \right] \quad (b)$$

59 where $M = \lceil \frac{T}{\tau} \rceil \leq T/\tau + 1$ is the number of intervals. Next, we analyze term (a) and term (b) separately.

Analysis of term (a) This term represents the regret of the algorithm compared to the piecewise-stationary reference sequence. The regret with respect to the expected risk $\mathcal{F}_c^t(\cdot)$ can be related to the unbiased empirical risk estimator $\widehat{\mathcal{F}}_c^t(\cdot)$:

$$\begin{aligned} \text{term (a)} &= \mathbb{E}_{1:T} \left[\sum_{t=1}^T \mathcal{F}_c^t(\{\psi_c^t, \phi_c^t\}) - \sum_{t=1}^T \mathcal{F}_c^t(\{\psi_c^t, \phi_c^t\}) \right] \\ &\leq \mathbb{E}_{1:T} \left[\sum_{t=1}^T \langle \nabla \mathcal{F}_c^t(\{\psi_c^t, \phi_c^t\}), (\{\psi_c^t, \phi_c^t\}) - (\{\psi_c^t, \phi_c^t\}) \rangle \right] \\ &= \mathbb{E}_{1:T} \left[\sum_{t=1}^T \langle \nabla \mathcal{F}_c^t(\{\psi_c^t, \phi_c^t\}) - \nabla \widehat{\mathcal{F}}_c^t(\{\psi_c^t, \phi_c^t\}), (\{\psi_c^t, \phi_c^t\}) - (\{\psi_c^t, \phi_c^t\}) \rangle \right] \\ &\quad + \mathbb{E}_{1:T} \left[\sum_{t=1}^T \langle \nabla \widehat{\mathcal{F}}_c^t(\{\psi_c^t, \phi_c^t\}), (\{\psi_c^t, \phi_c^t\}) - (\{\psi_c^t, \phi_c^t\}) \rangle \right], \end{aligned}$$

60 where the first inequality uses the convexity of $\mathcal{F}_c^t(\cdot)$. The first term above is zero because $\widehat{\mathcal{F}}_c^t$ is an unbiased estimator,
 61 i.e., $\nabla \mathcal{F}_c^t(\{\psi_c^t, \phi_c^t\}) = \mathbb{E}_t[\nabla \widehat{\mathcal{F}}_c^t(\{\psi_c^t, \phi_c^t\}) \mid 1 : t-1]$. For the model sequence $\{\{\psi_c^t, \phi_c^t\}\}_{t=1}^T$ generated by Fed-ADE, the
 62 following lemma holds:

Lemma 7. *Under the assumptions of Theorem 3, Fed-ADE with learning rate $\eta > 0$ in equation (6) satisfies*

$$\sum_{t=1}^T \left\langle \nabla \widehat{\mathcal{F}}_c^t(\{\psi_c^t, \phi_c^t\}), (\{\psi_c^t, \phi_c^t\}) - (\{\psi_c^t, \phi_c^t\}) \right\rangle \leq \frac{2\eta|\mathcal{I}|G^2T}{\sigma^2} + \frac{2\Gamma P_T + \Gamma^2}{2\eta},$$

63 where $P_T = \sum_{t=2}^T \|(\{\psi_c^t, \phi_c^t\}) - (\{\psi_c^{t-1}, \phi_c^{t-1}\})\|_2$ is the total variation of the comparator sequence and $G \triangleq$
 64 $\sup_{(\mathbf{x}_c^t, \mathbf{y}_c^t)} \|\nabla_{\{\psi_c, \phi_c\}} \mathcal{L}(\mathcal{H}(\{\psi_c, \phi_c\}, \mathbf{x}_c^t), \mathbf{y}_c^t)\|_2$ is an upper bound on the gradient norm.

Since the comparator sequence for term (a) changes only $M-1$ times, $P_T \leq \Gamma(M-1) \leq (\Gamma T)/\tau$. Taking expectations gives

$$\text{term (a)} \leq \frac{2\eta|\mathcal{I}|G^2T}{\sigma^2} + \frac{2\Gamma^2T/\tau + \Gamma^2}{2\eta}.$$

Analysis of term (b) This term accounts for the error from changing the reference sequence. Following [5], we get

$$\text{term (b)} \leq 2\tau \sum_{t=2}^T \sup_{\{\psi_c^t, \phi_c^t\}} |\mathcal{F}_c^t(\{\psi_c^t, \phi_c^t\}) - \mathcal{F}_c^{t-1}(\{\psi_c^t, \phi_c^t\})| \triangleq 2\tau \bar{\mathcal{S}}_c.$$

Combining the two terms, we obtain

$$\begin{aligned} \mathbb{E}_{1:T} \left[\sum_{t=1}^T \mathcal{F}_c^t(\{\psi_c^t, \phi_c^t\}) \right] - \sum_{t=1}^T \mathcal{F}_c^t(\{\psi_c^*, \phi_c^*\}) &\leq \frac{2\eta|\mathcal{I}|G^2T}{\sigma^2} + \frac{2\Gamma^2T/\tau + \Gamma^2}{2\eta} + 2B\tau\bar{\mathcal{S}}_c \\ &\leq \left(\frac{2|\mathcal{I}|G^2}{\sigma^2} + 2B^2 \right) \eta T + \frac{\Gamma^2}{\eta} + 4(\Gamma + 1) \sqrt{\frac{BT\bar{\mathcal{S}}_c}{\eta}}, \end{aligned}$$

65 where we set $\tau = \left\lceil \sqrt{\Gamma^2T/(\eta B\bar{\mathcal{S}}_c)} \right\rceil$ for optimal balance and $B \triangleq \sup_{(\mathbf{x}_c^t, \mathbf{y}_c^t)} |\mathcal{L}(\mathcal{H}(\{\psi_c, \phi_c\}, \mathbf{x}_c^t), \mathbf{y}_c^t)|$ is an upper bound
66 on the loss value, and G represents gradient norm upper bound.

Under the distribution shift assumption, the variation $\bar{\mathcal{S}}_c$ can be further bounded by the change in class and representation priors. Thus,

$$\mathbb{E}[\text{Reg}_T] \leq \left(\frac{2|\mathcal{I}|G^2}{\sigma^2} + 2B^2 \right) \eta T + \frac{\Gamma^2}{\eta} + 4(\Gamma + 1) \sqrt{\frac{BT\bar{\mathcal{S}}_c}{\eta}}.$$

67

□

68 E.2. Min-max Optimality

69 For online convex optimization with general convex losses, [5] showed that the dynamic regret has a lower bound of
70 $\Omega(\bar{\mathcal{S}}_c^{1/3}T^{2/3})$ when only noisy feedback is observed. Here, $\bar{\mathcal{S}}_c$ measures the total variation of the loss functions. The upper
71 bound in Theorem 3 matches this rate, showing that Fed-ADE achieves min-max optimality up to constants.

Lemma 8. *Under the same assumptions as Theorem 3, Fed-ADE with learning rate η satisfies*

$$\mathbb{E}[\text{Reg}_T] \leq 2 \left(\frac{|\mathcal{I}|G^2}{\sigma^2} + B^2 \right) \eta T + \frac{\Gamma^2}{\eta} + 4(\Gamma + 1) \sqrt{\frac{B\bar{\mathcal{S}}_cT}{\eta}} = \mathcal{O} \left(\eta T + \frac{1}{\eta} + \frac{\sqrt{\bar{\mathcal{S}}_cT}}{\eta} \right),$$

72 where $\sigma > 0$ is the minimum singular value of the confusion matrix \mathbf{M} and $\bar{\mathcal{S}}_c$ is the temporal variation of the loss.

73 This result follows by setting $\eta = \Theta(T^{-1/3}\bar{\mathcal{S}}_c^{1/3})$, which yields $\mathcal{O}(\bar{\mathcal{S}}_c^{1/3}T^{2/3})$ regret, matching the lower bound. Similar
74 reasoning shows that the algorithm also achieves $\mathcal{O}(\max\{\bar{\mathcal{S}}_c^{1/3}T^{2/3}, \sqrt{T}\})$ dynamic regret.

75 F. Convergence Analyses

76 We next show that our Fed-ADE provably converges when each per-timestep learning rate $\eta_c^t \in [\eta_{min}, \eta_{max}]$ obeys the same
77 upper bound used in Theorem 12. For the proof, we use three commonly used assumptions [6, 51] as follows.

78 **Assumption 9.** (Lipschitz parameter) For every client c , the loss function \mathcal{L} is continuously differentiable and there exist
79 Lipschitz constants $K_\psi, K_\phi, K_{\psi\phi}, K_{\phi\psi}$ holds that

- 80 • $\nabla_\psi \mathcal{L}(\bar{\psi}^{(r)}, \phi_c^{(r)})$ is K_ψ Lipschitz with respect to ψ and $K_{\psi\phi}$ Lipschitz with respect to $\phi_c^{(r)}$.
- 81 • $\nabla_\phi \mathcal{L}(\bar{\psi}^{(r)}, \phi_c^{(r)})$ is K_ϕ Lipschitz with respect to ϕ and $K_{\phi\psi}$ Lipschitz with respect to $\bar{\psi}^{(r)}$.

82 When the average loss function of total clients is defined as $\mathcal{L}(\bar{\psi}^{(r)}, \phi^{(r)}) = \frac{1}{|\mathcal{C}|} \sum_{c \in \mathcal{C}} \mathbb{E}[\mathcal{L}(\bar{\psi}^{(r)}, \phi_c^{(r)})]$, it has Lipschitz
83 constant K_ψ with respect to $\bar{\psi}^{(r)}$, $K_{\phi\psi}/\sqrt{|\mathcal{C}|}$ with respect to $\phi^{(r)}$, and $K_{\psi\phi}/|\mathcal{C}|$ with respect to any $\phi_c^{(r)}$. Further, there exist
84 χ to measure the relative cross-sensitivity of $\nabla_\phi \mathcal{L}(\bar{\psi}^{(r)}, \phi_c^{(r)})$ and $\nabla_\psi \mathcal{L}(\bar{\psi}^{(r)}, \phi_c^{(r)})$ as follows.

$$\chi = \frac{\max\{K_{\psi\phi}, K_{\phi\psi}\}}{\sqrt{K_\psi K_\phi}} \quad (19)$$

85

Assumption 10. (Bounded variance) The stochastic gradients in the client-side update algorithm are unbiased and have bounded variance as follows.

$$\mathbb{E}[\tilde{\nabla}_\phi \mathcal{L}(\bar{\psi}^{(r)}, \phi_c^{(r)})] = \nabla_\phi \mathcal{L}(\bar{\psi}^{(r)}, \phi_c^{(r)}) \quad (20)$$

$$\mathbb{E}[\tilde{\nabla}_\psi \mathcal{L}(\bar{\psi}^{(r)}, \phi_c^{(r)})] = \nabla_\psi \mathcal{L}(\bar{\psi}^{(r)}, \phi_c^{(r)}) \quad (21)$$

Furthermore, there exist constants σ_ϕ and σ_ψ that meet the following inequation.

$$\mathbb{E}[\|\tilde{\nabla}_\phi \mathcal{L}(\bar{\psi}^{(r)}, \phi_c^{(r)}) - \nabla_\phi \mathcal{L}(\bar{\psi}^{(r)}, \phi_c^{(r)})\|^2] \leq \sigma_\phi^2 \quad (22)$$

$$\mathbb{E}[\|\tilde{\nabla}_\psi \mathcal{L}(\bar{\psi}^{(r)}, \phi_c^{(r)}) - \nabla_\psi \mathcal{L}(\bar{\psi}^{(r)}, \phi_c^{(r)})\|^2] \leq \sigma_\psi^2 \quad (23)$$

86 We can consider $\nabla_\psi \mathcal{L}_c(\bar{\psi}^{(r)}, \phi_c^{(r)})$ as a stochastic partial gradient of average loss function $\mathcal{L}(\bar{\psi}^{(r)}, \phi^{(r)})$ with respect to $\bar{\psi}^{(r)}$.

Assumption 11. (Partial gradient diversity) There exist constants $\delta \geq 0$ and $\rho \geq 0$ such that for $\bar{\psi}^{(r)}$ and all personalized layers $\phi^{(r)} = [\phi_1^{(r)}, \phi_2^{(r)}, \dots]$, which meet the following equations.

$$\frac{1}{|\mathcal{C}|} \sum_{c \in \mathcal{C}} \|\nabla_\psi \mathcal{L}(\bar{\psi}^{(r)}, \phi_c^{(r)}) - \nabla_\psi \mathcal{L}(\bar{\psi}^{(r)}, \phi^{(r)})\|^2 \leq \delta^2 + \rho^2 \|\nabla_\psi \mathcal{L}(\bar{\psi}^{(r)}, \phi^{(r)})\|^2 \quad (24)$$

87 Additionally, $\mathcal{L}(\bar{\psi}, \phi^{(r)})$ is bounded below by $\hat{\mathcal{L}}$, where $\hat{\mathcal{L}}$ meets the equation $\Delta \mathcal{L}^{(0)} = \mathcal{L}(\bar{\psi}^{(0)}, \phi^{(0)}) - \hat{\mathcal{L}}$ at initial communi-
88 cation round.

89 **Theorem 12.** The convergence of the Fed-ADE is bounded as follows, where $\Delta \mathcal{L}^{(0)}$ denotes the difference between a bound
90 and initial value of loss $\mathcal{L}(\bar{\psi}, \phi_c)$.²

$$\frac{1}{R} \sum_{r=0}^{R-1} \left(\frac{1}{K_\psi} \mathbb{E}[\Delta_\psi^{(r)}] + \frac{|\mathcal{C}^{(r)}|}{|\mathcal{C}| K_\phi} \mathbb{E}[\Delta_\phi^{(r)}] \right) \leq \frac{\Delta \mathcal{L}^{(0)}}{\eta R} + \eta \sigma_1^2 + \eta^2 \sigma_2^2 \quad (25)$$

²To simplify the inequality in the following assumptions and theorems, we use following shorthands: $\phi^{(r)} = \{\phi_c^{(r)} | c \in \mathcal{C}^{(r)}\}$, $\mathcal{L}(\bar{\psi}^{(r)}, \phi_c^{(r)}) = \mathcal{L}(\{\bar{\psi}^{(r)}, \phi_c^{(r)}\}; \mathbf{x}_c^t)$, $\mathcal{L}(\bar{\psi}^{(r)}, \phi^{(r)}) = \sum_{c \in \mathcal{C}} \frac{N_c}{N} \mathcal{L}(\bar{\psi}^{(r)}, \phi_c^{(r)})$, $\Delta_\psi^{(r)} = \|\nabla_\psi \sum_{c \in \mathcal{C}} \frac{N_c}{N} \mathcal{L}(\bar{\psi}^{(r)}, \phi_c^{(r)})\|^2$, $\Delta_\phi^{(r)} = \frac{1}{|\mathcal{C}|} \sum_{c \in \mathcal{C}} \|\nabla_\phi \mathcal{L}(\bar{\psi}^{(r)}, \phi_c^{(r)})\|^2$.

Here, the constants σ_1^2 and σ_2^2 are defined as follows.

$$\sigma_1^2 = \frac{\delta^2}{K_\psi} \left(1 - \frac{|\mathcal{C}^{(r)}|}{|\mathcal{C}|}\right) + \frac{\sigma_\psi^2}{K_\psi} + \frac{\sigma_\phi^2 (|\mathcal{C}^{(r)}| + \chi^2 (|\mathcal{C}| - |\mathcal{C}^{(r)}|))}{K_\phi |\mathcal{C}|} \quad (26)$$

$$\sigma_2^2 = \left(\frac{\sigma_\psi^2 + \delta^2}{K_\psi} + \frac{\sigma_\phi^2 |\mathcal{C}^{(r)}|}{K_\phi |\mathcal{C}|}\right) (1 - E^{-1}) + \frac{\chi^2 \sigma_\phi^2}{K_\phi} \quad (27)$$

This convergence is bounded when the learning rates are chosen as $\eta_\psi = \eta / (K_\psi E)$ and $\eta_\phi = \eta / (K_\phi E)$, where η satisfies the following inequality.

$$\eta \leq \min \left\{ \frac{1}{24(1 + \rho^2)}, \frac{|\mathcal{C}^{(r)}|}{128\chi^2 (|\mathcal{C}| - |\mathcal{C}^{(r)}|)}, \sqrt{\frac{|\mathcal{C}^{(r)}|}{\chi^2 |\mathcal{C}|}} \right\} \quad (28)$$

Proof. The proof of Theorem 12 begins with the following equation describing the update from round r to round $r + 1$.

$$\mathcal{L}(\psi^{(r+1)}, \phi^{(r+1)}) - \mathcal{L}(\bar{\psi}^{(r)}, \phi^{(r)}) = \mathcal{L}(\bar{\psi}^{(r)}, \phi^{(r+1)}) - \mathcal{L}(\bar{\psi}^{(r)}, \phi^{(r)}) \quad (29)$$

$$+ \mathcal{L}(\psi^{(r+1)}, \phi^{(r+1)}) - \mathcal{L}(\bar{\psi}^{(r)}, \phi^{(r+1)}) \quad (30)$$

91 In the case of the update for the personalized layers in (29), since each client is training with their own data, it can be
 92 demonstrated similarly to the gradient update in conventional deep learning approaches. However, in the case of the update for
 93 the shared layers in (30), as gradient descent is performed for different client models, we aim to establish its convergence.

The smoothness bound for the updating of the shared layers gives an equation as follows.

$$\begin{aligned} \mathcal{L}(\psi^{(r+1)}, \phi^{(r+1)}) - \mathcal{L}(\bar{\psi}^{(r)}, \phi^{(r+1)}) &\leq \left\langle \nabla_\psi \mathcal{L}(\bar{\psi}^{(r)}, \phi^{(r+1)}), \psi^{(r+1)} - \bar{\psi}^{(r)} \right\rangle \\ &+ \frac{K_\psi}{2} \|\psi^{(r+1)} - \bar{\psi}^{(r)}\|^2 \end{aligned} \quad (31)$$

According to Lipschitz continuity, we can express the gradient as shown in the (31), when $\langle A, B \rangle$ means $A^T B$. The proof proceeds by demonstrating that for all communication rounds $0 < r < R$, the total sum of this gradient does not diverge but rather remains bounded, as is well known. However, a key distinction between the conventional deep learning approach and the personalized federated learning method is that only a subset of clients among the total clients participate in training during a single round. In these conditions, the virtual full participants $\tilde{\phi}^{(r+1)}$ are used to move all the dependencies of the model update on the participants $\mathcal{C}^{(r)} \in \mathcal{C}$ [51]. When applying the virtual full participants to the (31), we can get an equation as follows.

$$\begin{aligned} \left\langle \nabla_\psi \mathcal{L}(\bar{\psi}^{(r)}, \phi^{(r+1)}), \psi^{(r+1)} - \bar{\psi}^{(r)} \right\rangle &= \left\langle \nabla_\psi \mathcal{L}(\bar{\psi}^{(r)}, \tilde{\phi}^{(r+1)}), \psi^{(r+1)} - \psi^{(r)} \right\rangle \\ &+ \left\langle \nabla_\psi \mathcal{L}(\bar{\psi}^{(r)}, \tilde{\phi}^{(r+1)}) - \nabla_\psi \mathcal{L}(\bar{\psi}^{(r)}, \phi^{(r+1)}), \psi^{(r+1)} - \bar{\psi}^{(r)} \right\rangle \end{aligned} \quad (32)$$

Applying Young's inequality and Lipschitz inequality to the (32), the (31) can be expressed as follows.

$$\begin{aligned} \mathcal{L}(\psi^{(r+1)}, \phi^{(r+1)}) - \mathcal{L}(\bar{\psi}^{(r)}, \phi^{(r+1)}) &\leq \left\langle \nabla_\psi \mathcal{L}(\bar{\psi}^{(r)}, \phi^{(r+1)}), \psi^{(r+1)} - \bar{\psi}^{(r)} \right\rangle \\ &+ \frac{1}{2K_\psi} \|\nabla_\psi \mathcal{L}(\bar{\psi}^{(r)}, \tilde{\phi}^{(r+1)}) - \nabla_\psi \mathcal{L}(\bar{\psi}^{(r)}, \phi^{(r+1)})\|^2 \\ &+ K_\psi \|\psi^{(r+1)} - \bar{\psi}^{(r)}\|^2 \end{aligned} \quad (33)$$

This allows us to take an expectation. Therefore, the final bounded form can be described as follows with Assumption 11, where the $\tilde{\Delta}_\psi^{(r)}$ is the analog of $\Delta_\psi^{(r)}$ with the virtual variable $\tilde{\phi}^{(r+1)}$.

$$\frac{1}{R} \sum_{r=0}^{R-1} \left(\frac{\eta_\psi E}{8} \mathbb{E}[\tilde{\Delta}_\psi^{(r)}] + \frac{\eta_\phi E c}{16|\mathcal{C}|} \mathbb{E}[\Delta_\phi^{(r)}] \right) \leq \frac{\Delta \mathcal{L}^{(0)}}{R} + \mathcal{O}(\eta_\psi^2 + \eta_\phi^2). \quad (34)$$

94 **Bound of shared layers update** When using (34), the update of the shared layers can be bounded as follows.

$$\begin{aligned} & \mathcal{L}(\psi^{(r+1)}, \phi^{(r+1)}) - \mathcal{L}(\psi^{(r)}, \phi^{(r+1)}) \\ & \leq \left\langle \nabla_p \mathcal{L}(\psi^{(r)}, \tilde{\phi}^{(r+1)}), \psi^{(r+1)} - \psi^{(r)} \right\rangle K_\psi \|\psi^{(r+1)} - \psi^{(r)}\|^2 + \frac{\chi^2 K_\phi}{2|\mathcal{C}|} \sum_{c \in \mathcal{C}} \|\tilde{\phi}_c^{(r+1)} - \phi_c^{(r+1)}\|^2 \end{aligned} \quad (35)$$

Thus, the dependence of the second line term on $\phi^{(r+1)}$ is successfully eliminated. To bind the third line term, expectation can be used as follows.

$$\begin{aligned} & \mathbb{E} \left[\mathcal{L}(\psi^{(r+1)}, \phi^{(r+1)}) - \mathcal{L}(\psi^{(r)}, \phi^{(r+1)}) \right] \\ & \leq -\frac{\eta_\psi E}{4} \mathbb{E}[\tilde{\Delta}_\psi^{(r)}] + \frac{2\eta_\psi K_\psi^2}{|\mathcal{C}|} \sum_{c \in \mathcal{C}} \sum_{e=0}^E \mathbb{E} \|\tilde{\psi}_{c,e}^{(r)} - \psi^{(r)}\|^2 + 4\eta_\phi^2 E^2 K_\phi \sigma_\phi^2 \chi^2 \left(1 - \frac{|\mathcal{C}^{(r)}|}{|\mathcal{C}|}\right) \\ & \quad + \frac{K_\psi \eta_\psi^2 E^2}{|\mathcal{C}^{(r)}|} (\sigma_\psi^2 + 3\delta^2 (1 - \frac{|\mathcal{C}^{(r)}|}{|\mathcal{C}|})) + 8\eta_\psi^2 E^2 K_\phi \chi^2 \left(1 - \frac{|\mathcal{C}^{(r)}|}{|\mathcal{C}|}\right) \Delta_\phi^{(r)} \end{aligned} \quad (36)$$

Note that $24K_\psi \eta_\psi E(1 + \rho^2) \leq 1$ is used to simplify the coefficients of some of the terms above. The term in the third line is referred to as client shift in the literature. This term can be described using the virtual variable $\tilde{\phi}^{(r+1)}$ as follows.

$$\begin{aligned} & \frac{2\eta_\psi K_\psi^2}{|\mathcal{C}|} \sum_{c \in \mathcal{C}} \sum_{e=0}^E \mathbb{E} \|\tilde{\psi}_{c,e}^{(r)} - \psi^{(r)}\|^2 \\ & \leq \frac{16\eta_\psi^3 K_\psi^2 E(E-1)}{|\mathcal{C}|} (\delta^2 + \rho^2 \mathbb{E} \|\nabla_\psi \mathcal{L}(\psi^{(r)}, \tilde{\phi}^{(r+1)})\|^2) + 8\eta_\psi^3 K_\psi^2 E^2 (E-1) \sigma_\psi^2 \end{aligned} \quad (37)$$

By inputting (37) into the previous inequality of expectation (36), the inequality can be described as follows.

$$\begin{aligned} & \mathbb{E}[\mathcal{L}(\psi^{(r+1)}, \phi^{(r+1)}) - \mathcal{L}(\psi^{(r)}, \phi^{(r+1)})] \\ & \leq -\frac{\eta_\psi E}{8} \mathbb{E}[\tilde{\Delta}_\psi^{(r)}] + \frac{K_\psi \eta_\psi^2 E^2}{|\mathcal{C}^{(r)}|} (\sigma_\psi^2 + 2\delta^2 (1 - \frac{|\mathcal{C}^{(r)}|}{|\mathcal{C}|})) + 4\eta_\phi^2 E^2 K_\phi \sigma_\phi^2 \chi^2 \left(1 - \frac{|\mathcal{C}^{(r)}|}{|\mathcal{C}|}\right) \\ & \quad + 8\eta_\phi^2 E^2 K_\phi \chi^2 \left(1 - \frac{|\mathcal{C}^{(r)}|}{|\mathcal{C}|}\right) \Delta_\phi^{(r)} + 8\eta_\psi^2 K_\psi^3 E^2 \left(1 - \frac{|\mathcal{C}^{(r)}|}{|\mathcal{C}|}\right) (\sigma_\psi^2 + 2\delta_\psi) \end{aligned} \quad (38)$$

95 **Bound of personalized layers** The above analysis for updating the shared layers can be applied to updating the personalized
 96 layers. This analysis provides the following, which displays the bound of the update of the personalized layers. It can be
 97 simplified some coefficients using $128\eta_\phi E K_\phi \chi^2 (\frac{|\mathcal{C}|}{|\mathcal{C}^{(r)}|} - 1) \leq 1$.

$$\begin{aligned} & \mathbb{E} \left[\mathcal{L}(\phi^{(r+1)}, \psi^{(r+1)}) - \mathcal{L}(\psi^{(r)}, \phi^{(r)}) \right] \\ & \leq -\frac{\eta_\psi E}{8} \mathbb{E}[\tilde{\Delta}_\psi^{(r)}] - \frac{\eta_\phi E |\mathcal{C}^{(r)}|}{16|\mathcal{C}|} \mathbb{E}[\Delta_\phi^{(r)}] + 4\eta_\phi^2 K_\phi E^2 \sigma_\phi^2 \left(\frac{|\mathcal{C}^{(r)}|}{|\mathcal{C}|} + \chi^2 \left(1 - \frac{|\mathcal{C}^{(r)}|}{|\mathcal{C}|}\right) \right) \\ & \quad + \frac{\eta_\psi^2 K_\psi E^2}{|\mathcal{C}^{(r)}|} (\sigma_\psi^2 + 2\delta^2 (1 - \frac{|\mathcal{C}^{(r)}|}{|\mathcal{C}|})) + 8\eta_\psi^2 K_\psi^2 E^2 (E-1) (\sigma_\psi^2 + 2\delta^2) + \frac{|\mathcal{C}^{(r)}|}{|\mathcal{C}|} 4\eta_\phi^3 K_\phi^2 E^2 (E-1) \sigma_\phi^2 \end{aligned} \quad (39)$$

98 The summation of inequality (39) for considering the communication round is given below.

$$\begin{aligned}
& \frac{1}{R} \sum_{r=0}^{R-1} \left(\frac{\eta_\psi E}{8} \mathbb{E}[\tilde{\Delta}_\psi^{(r)}] + \frac{\eta_\phi E |\mathcal{C}^{(r)}|}{|\mathcal{C}|} \mathbb{E}[\Delta_\phi^{(r)}] \right) \\
& \leq \frac{\Delta \mathcal{L}^{(0)}}{R} + 4\eta_\phi^2 K_\phi E^2 \sigma_\phi^2 \left(\frac{|\mathcal{C}^{(r)}|}{|\mathcal{C}|} + \chi^2 \left(1 - \frac{|\mathcal{C}^{(r)}|}{|\mathcal{C}|} \right) \right) + \frac{\eta_\psi^2 K_\psi E^2}{|\mathcal{C}^{(r)}|} \left(\sigma_\psi^2 + 2\delta^2 \left(1 - \frac{|\mathcal{C}^{(r)}|}{|\mathcal{C}|} \right) \right) \\
& \quad + 8\eta_\psi^3 K_\psi^2 E^2 (E-1) (\sigma_\psi^2 + 2\delta^2) + \frac{|\mathcal{C}^{(r)}|}{|\mathcal{C}|} 4\eta_\phi^3 K_\phi^2 E^2 (E-1) \sigma_\phi^2
\end{aligned} \tag{40}$$

99 This is the bound about virtual variable $\tilde{\phi}^{(r+1)}$. To describe this bound with the participants $\phi^{(r+1)}$, the assumption of the
100 Lipschitz parameter is used as follows.

$$\begin{aligned}
& \mathbb{E} \left\| \nabla_\psi \mathcal{L}(\psi^{(r)}, \phi^{(r)}) - \nabla_\psi \mathcal{L}(\psi^{(r)}, \tilde{\phi}^{(r+1)}) \right\|^2 \\
& \leq \frac{1}{|\mathcal{C}|} \sum_{c \in \mathcal{C}} \mathbb{E} \left\| \nabla_\psi \mathcal{L}(\psi^{(r)}, \phi_c^{(r)}) - \nabla_\psi \mathcal{L}(\psi^{(r)}, \tilde{\phi}_c^{(r+1)}) \right\|^2 \\
& \leq \frac{\chi^2 K_\phi K_\psi}{|\mathcal{C}|} \sum_{c \in \mathcal{C}} \mathbb{E} \left\| \tilde{\phi}_c^{(r+1)} - \phi_c^{(r)} \right\|^2 \\
& \leq \frac{\chi^2 K_\phi K_\psi}{|\mathcal{C}|} \sum_{c \in \mathcal{C}} \left(16\eta_\phi^2 E^2 \left\| \nabla_p \mathcal{L}(\psi^{(r)}, \phi_c^{(r)}) \right\|^2 + 8\eta_\phi^2 E^2 \sigma_\phi^2 \right) = 8\eta_\phi^2 E^2 \chi^2 K_\phi K_\psi \left(\sigma_\phi^2 + 2\mathbb{E}[\Delta_\phi^{(r)}] \right)
\end{aligned} \tag{41}$$

Using the Cauchy–Schwarz inequality described as follows.

$$\left\| \nabla_\psi \mathcal{L}(\psi^{(r)}, \phi^{(r)}) \right\|^2 \leq 2 \left\| \nabla_\psi \mathcal{L}(\psi^{(r)}, \phi^{(r)}) - \nabla_\psi \mathcal{L}(\psi^{(r)}, \tilde{\phi}^{(r+1)}) \right\|^2 + 2 \left\| \nabla_\psi \mathcal{L}(\psi^{(r)}, \tilde{\phi}^{(r+1)}) \right\|^2 \tag{42}$$

By utilizing inequality (42), we obtain the following result.

$$\mathbb{E}[\Delta_\phi^r] \leq \mathbb{E}[\tilde{\Delta}_\psi^{(r)}] + 16\eta_\phi^2 E^2 \chi^2 K_\phi K_\psi (\sigma_\phi^2 + 2\mathbb{E}[\Delta_\phi^{(r)}]) \tag{43}$$

Through inequality (43), we can get the following inequality.

$$\frac{\eta_\psi E}{16} \mathbb{E}[\Delta_\psi^{(r)}] + \frac{\eta_\phi E |\mathcal{C}^{(r)}|}{32|\mathcal{C}|} \mathbb{E}[\Delta_\phi^{(r)}] \leq \frac{\eta_\psi E}{8} \mathbb{E}[\tilde{\Delta}_\psi^{(r)}] + \frac{\eta_\phi E |\mathcal{C}^{(r)}|}{16|\mathcal{C}|} \mathbb{E}[\Delta_\phi^{(r)}] + \eta_\psi \eta_\phi^2 E^3 \sigma_\phi^2 \chi^2 K_\phi K_\psi \tag{44}$$

101 By summing inequality (44) over the total number of communication rounds and applying the parameter settings $\eta_\phi =$
102 $\eta/(K_\phi E)$ and $\eta_\psi = \eta/(K_\psi E)$, the proof is completed. \square

103 To ensure that the convergence analysis above remains valid when using an adaptive learning rate in (9), we require
104 that the value of the adaptive learning rate η_c across all clients and rounds is bounded by the same upper bound used in
105 Theorem 12. Since our proposed label shift adaptive learning rate can not have value outside of the bound $[\eta_{min}, \eta_{max}]$ (where
106 $\eta_c^t \in [\eta_{min}, \eta_{max}]$), by appropriately choosing $[\eta_{min}, \eta_{max}] \leq \min \left\{ \frac{1}{24(1+\rho^2)}, \frac{|\mathcal{C}^{(r)}|}{128\chi^2(|\mathcal{C}|-|\mathcal{C}^{(r)}|)}, \sqrt{\frac{|\mathcal{C}^{(r)}|}{\chi^2|\mathcal{C}|}} \right\}$, we guarantee that
107 all learning rates used in the adaptive scheme satisfy the conditions required for Theorem 12.

108 G. Simulation Settings

109 G.1. Datasets

110 We evaluate on various benchmarks including five image-based benchmarks- Tiny ImageNet [32], CIFAR-10 [29], CIFAR-
111 10-C [19], CIFAR-100 [29], and CIFAR-100-C [19], and a text-based benchmark LAMA [50]. The following is a detailed
112 description of the datasets used in this work.

- 113 • **Tiny ImageNet [32]:** Tiny ImageNet is a subset of the ImageNet dataset, containing 200 object classes, each with 500
114 training images, 50 validation images, and 50 test images. All images are RGB and resized to 64×64 pixels.
- 115 • **CIFAR-10 [29]:** The CIFAR-10 dataset consists of 60,000 RGB images of size 32×32 across 10 classes, with 50,000
116 images for training and 10,000 for testing. We applied standard preprocessing including normalization with dataset-wide
117 mean and standard deviation.
- 118 • **CIFAR-10-C [19]:** CIFAR-10-C is an extension of the CIFAR-10 dataset designed to evaluate model robustness against
119 common corruptions. It contains 950,000 images created by applying 19 different types of corruption—such as noise, blur,
120 weather effects, and digital distortions—at 5 severity levels to the original CIFAR-10 test set of 10,000 images.
- 121 • **CIFAR-100 [29]:** CIFAR-100 is a fine-grained variant of CIFAR-10 containing 100 object classes. It consists of 60,000
122 color images at 32×32 pixels, split into 50,000 training and 10,000 test images. The 100 classes are grouped into 20
123 superclasses for hierarchical classification research.
- 124 • **CIFAR-100-C [19]:** CIFAR-100C applies the same 19 corruption types with 5 severity levels from CIFAR-10-C onto
125 the CIFAR-100 dataset. It consists of 950,000 corrupted images used as a benchmark for analyzing model robustness on
126 fine-grained image classification tasks.
- 127 • **LAMA [50]:** LLanguage Model Analysis is a benchmark dataset designed to evaluate the factual and commonsense
128 knowledge acquired by language models. In our experiments, we focus on the T-REx subset of LAMA, which consists
129 of (subject, relation, object) triples derived from Wikidata. This subset consists of 41 relational tasks, with 100 samples
130 randomly sampled per task. To simulate heterogeneous client environments, we assign disjoint subsets of tasks to different
131 clients.

132 G.2. Online Distribution Shift Modeling

133 To emulate online distribution shift in (1), each client’s data distribution moves from a common initial uniform vector \mathbf{Q}_c^0
134 toward a client-specific target \mathbf{Q}_c^T according to a time-dependent weight $\omega(t)$. We evaluate under four prototypical distribution
135 shift schedules—linear (Lin.), sine (Sin.), square (Squ.), and Bernoulli (Ber.), each driving a client’s distribution from a
136 uniform start toward a client-specific target [10, 45, 56, 65, 69]. Following is a detailed description of four prototypical
137 label-shift schedules according to different $\omega(t)$ functions, and Figure 3 shows how the parameter $\omega(t)$ changes over time for
138 various shift types.

- 139 • **Linear shift (Lin.):** $\omega(t) = \frac{t}{T}$, where t is the current timestep and T is the total number of timesteps. This provides a smooth
140 and gradual transition from \mathbf{Q}_c^0 to \mathbf{Q}_c^T over time.
- 141 • **Sine shift (Sin.):** $\omega(t) = \sin\left(\frac{\pi t}{\sqrt{T}}\right)$. This introduces cyclical variations, capturing periodic shifts in the data distribution.
- 142 • **Square Shift (Squ.):** $\omega(t)$ alternates between 0 and 1 every $\frac{\sqrt{T}}{2}$ timestep, resulting in abrupt, periodic changes in the class
143 distribution.
- 144 • **Bernoulli shift (Bern.):** $\omega(t)$ retains its previous value $\omega(t-1)$ with a probability of $\frac{1}{\sqrt{T}}$, or flips to $1 - \omega(t-1)$.
145 This configuration captures stochastic variations, ensuring that class priors change randomly while maintaining a scale
146 proportional to \sqrt{T} .

147 For \mathbf{Q}_c^T , we used the Dirichlet distribution to consider heterogeneous data distributions among clients. The distribution
148 factor α of the Dirichlet distribution controls the level of heterogeneity among clients. Smaller α leads more heterogeneous
149 distribution, and Larger α leads to Independent and identically distributed data system. In this study, we set $Dir(\alpha = 0.1)$.
150 Figure 4 visualizes the heterogeneity among clients according to Dirichlet distribution factor α .

151 G.3. Distribution Shift Scenarios

152 In all experiments, each target distribution of client \mathbf{Q}_{y^T} in (1) is defined according to the specific type of distribution shift
153 being simulated. We summarize below the construction and datasets used for each shift scenario.

- 154 • **Label Shift.** Label shift simulations are conducted on Tiny ImageNet [32], CIFAR-10 [29], and LAMA [50], where each
155 client’s class prior evolves over time according to one of four temporal schedules (Lin., Sin, Squ., Ber.) which follows (1).
- 156 • **Covariate Shift.** The target feature distribution $\mathbf{Q}_{c,x}^T$ reflects corruption-induced changes to the input domain. Each

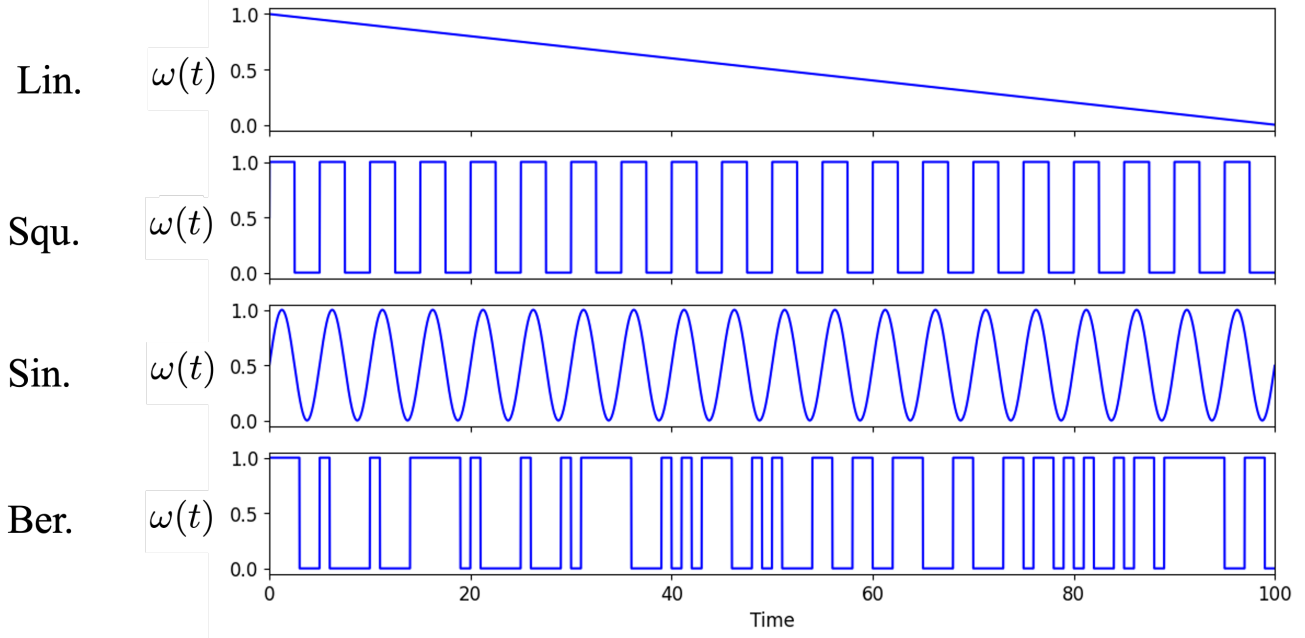


Figure 3. Visualization of how the parameter $\omega(t)$ changes over time for various shift types. The Lin. presents a steady and continuous increase or decrease, reflecting a gradual and predictable change. Both the Squ. and Sin. display periodic patterns. The Squ. alternates sharply between two values at regular intervals, while the Sin. oscillates smoothly in a wave-like manner. The Ber. introduces randomness at each timestep, resulting in a stochastic and less predictable trajectory for $\omega(t)$.

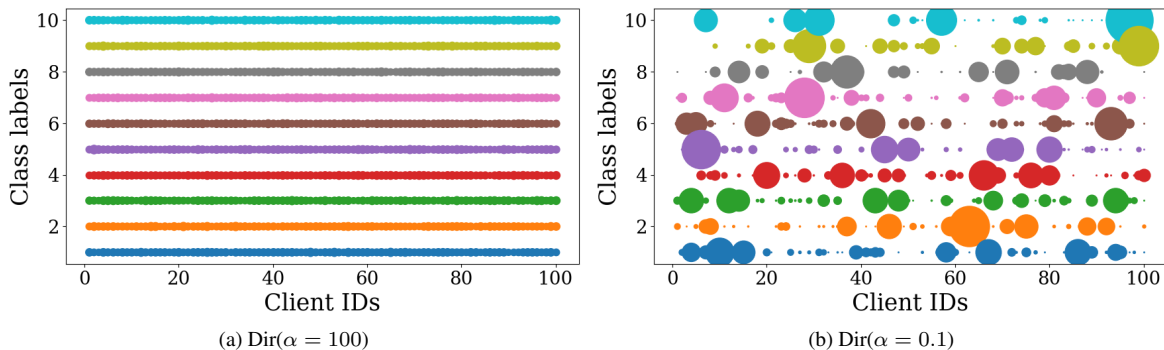


Figure 4. Visualization of heterogeneity among clients according to Dirichlet distribution factor α in CIFAR-10 dataset. The x -axis represents the ID of the client, while the y -axis represents the class of the data. The color of the circles varies with each class type, and the size of the circles indicates the size of the data. This level of heterogeneity is not dependent on the data set but on α values.

157 client is assigned a distinct and fixed corruption type (e.g., Gaussian noise, blur), while the corruption severity level shifts
 158 dynamically over time according to one of four temporal schedules. We use the CIFAR-10 [29] \rightarrow CIFAR-10-C [19] and
 159 CIFAR-100 [29] \rightarrow CIFAR-100-C [19] benchmarks to simulate these covariate-shift environments.

160 G.4. Compared Methods

161 We evaluate our proposed Fed-ADE against six existing unsupervised online adaptation baselines:

- 162 • **FTH** [58]: FTH maintains an average of estimated label distributions across all previous timesteps. By aggregating the
 163 entire history of predictions, FTH assumes that past distributional information is informative for current adaptation, making
 164 it robust to gradual, monotonic shifts but potentially less responsive to abrupt changes.
- 165 • **ATLAS** [2]: ATLAS maintains an ensemble of base learners, each optimized with a distinct learning rate. Using an
 166 exponential weighting scheme, ATLAS dynamically aggregates these learners according to their recent performance,

- 167 enabling robust adaptation to various patterns and magnitudes of label shift.
- 168 • **UDA [18]**: UDA is a technique where a model is trained on labeled data from a source domain and adapted to work well
- 169 on a target domain that has no labels, even though the data distributions are different. The main idea is to reduce the gap
- 170 between the source and target domains by aligning their feature representations or adjusting the model so it can generalize to
- 171 the target data, all without using any target domain labels.
- 172 • **UNIDA [63]**: UNiDA is a method designed to handle domain adaptation scenarios. This assumes the same class set across
- 173 domains, which aims to adapt models by distinguishing between common, source-private, and target-private classes, often
- 174 using prototype-based or confidence-based techniques to avoid misclassifying unknown target classes as known ones.
- 175 • **Fed-POE [43]**: Fed-POE is a federated online adaptation framework in which each client combines local fine-tuning with
- 176 an ensemble of periodically aggregated global models. While this approach supports client personalization and leverages
- 177 federated knowledge, it employs a fixed learning rate and incurs computational overhead from managing multiple model
- 178 instances.
- 179 • **FedCCFA [9]**: FedCCFA is a federated learning framework for concept shift, where clients may experience different
- 180 distributions. It performs class-level classifier clustering to share aggregated class-wise classifiers among similar clients, and
- 181 aligns feature spaces using clustered feature anchors with an entropy-based adaptive weighting.

182 G.5. Hyperparameters

183 In this section, we summarize hyperparameters used in the simulations.

Table 6. Hyperparameters for Tiny ImageNet, CIFAR-10, and CIFAR-100.

Hyperparameter	Tiny ImageNet	CIFAR-10	CIFAR-100
Number of clients	100	100	100
Communication rounds (R)	10	10	10
Local epochs	4	4	4
Optimizer	SGD	SGD	SGD
Batch size	128	32	32
Participant rate	10%	10%	10%
Model architecture	ResNet18	CNN w/ 3 residual blocks	CNN w/ 3 residual blocks
Shared layers (ψ_c)	Representation layers	Representation layers	Representation layers
Lowest learning rate	5×10^{-6}	5×10^{-6}	5×10^{-6}
Highest learning rate	10^{-4}	10^{-4}	10^{-4}

Table 7. Hyperparameters for CIFAR-10-C, CIFAR-100-C, and LAMA

Hyperparameter	CIFAR-10-C	CIFAR-100-C	LAMA
Number of clients	100	100	5
Communication rounds (R)	10	10	5
Local epochs	4	4	1
Optimizer	SGD	SGD	AdamW
Batch size	32	32	4
Participant rate	10%	10%	100%
Model architecture	CNN w/ 3 residual blocks	CNN w/ 3 residual blocks	Llama-3.2-3B
Shared layers (ψ_c)	Representation layers	Representation layers	LoRA on q/k/v (layers 0-13)
Lowest learning rate	5×10^{-6}	5×10^{-6}	1×10^{-3}
Highest learning rate	1×10^{-4}	1×10^{-4}	1×10^{-2}

184 G.6. Hardware Specifications

Table 8. Hardware specifications used for all simulations.

Component	Specification
CPU	AMD Ryzen 9 7950X, 16-Core, 32 Threads
GPU	NVIDIA GeForce RTX 3090X2
RAM	256 GB DDR5
Storage	1.8 TB NVMe SSD

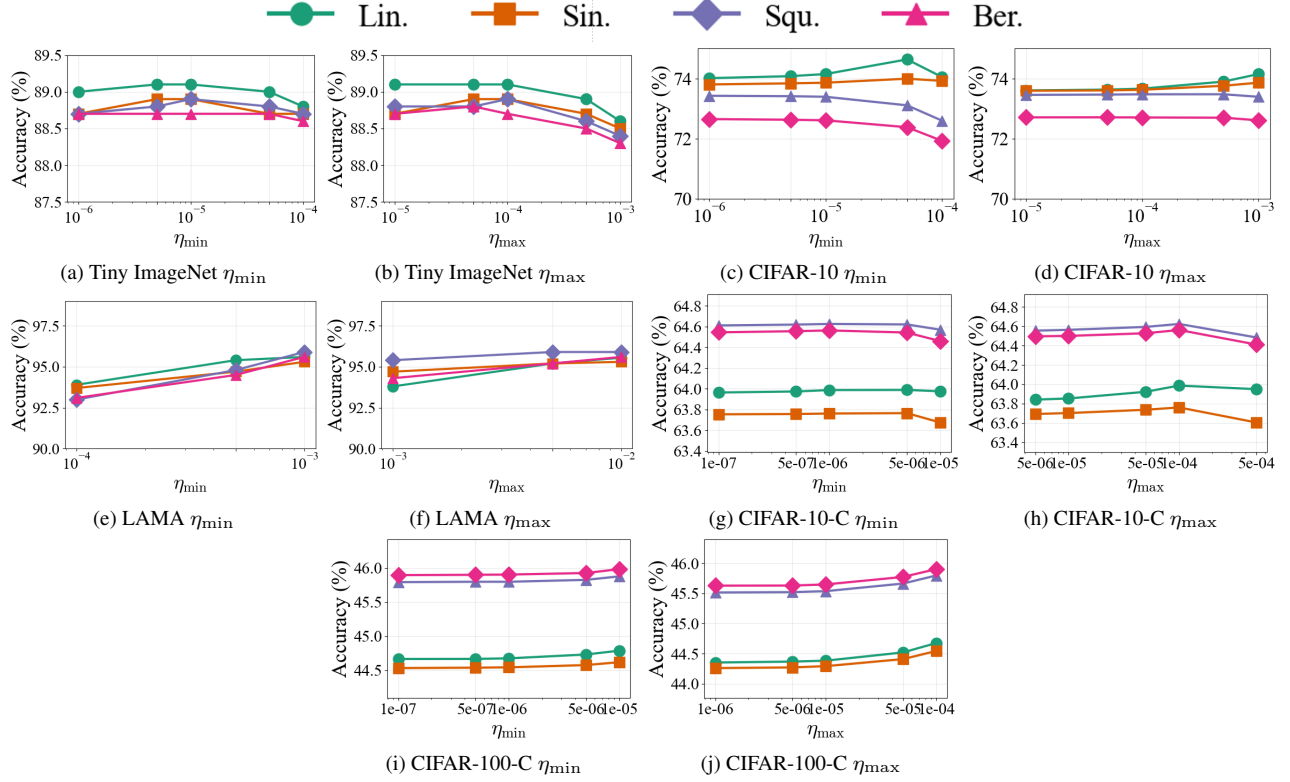


Figure 5. Impact of learning rate bounds on Fed-ADE under label shift scenarios.

185 H. Extra Simulation Results

186 H.1. Impact of Learning Rate Selection

187 Figure 5 examines the sensitivity of Fed-ADE to the choice of learning rate bounds under all scenarios and benchmarks.
 188 For image benchmarks, we vary η_{\min} while fixing $\eta_{\max} = 10^{-4}$, and vary η_{\max} while fixing $\eta_{\min} = 5 \times 10^{-6}$. For text
 189 benchmark LAMA, we vary η_{\min} while fixing $\eta_{\max} = 10^{-2}$, and vary η_{\max} while fixing $\eta_{\min} = 10^{-4}$. Across all cases,
 190 Fed-ADE consistently maintains high accuracy within a broad range of hyperparameter choices. While extreme values lead to
 191 slight performance degradation, the overall trend shows that our dynamics-driven adaptation is far less sensitive to the exact
 192 learning rate bounds compared to fixed-rate baselines. This robustness highlights that Fed-ADE reduces the need for extensive
 193 hyperparameter tuning, an essential property in federated and online learning where search budgets are limited.

Table 9. Performance comparison of different online distribution shift adaptation methods across datasets and shift types with pretrained model from Gaussian distribution (average accuracy (%) and average wall time (sec.))

Dataset	Shift	Localized Learning				Federated Learning					Fed-ADE
		FTH	ATLAS	UNIDA	UDA	Fed-POE	FedCCFA	FixLR(Low)	FixLR(Mid)	FixLR(High)	
(i) Label Shift Scenarios											
Tiny ImageNet	Lin.	74.1±2.3	70.7±2.1	75.1±2.1	63.1±2.3	81.5±1.2	81.7±1.8	81.5±1.1	81.2±1.3	78.1±1.5	84.5±1.1
	Sin.	73.7±2.4	71.1±1.9	74.8±2.5	61.3±2.1	82.1±1.2	80.9±1.1	81.9±1.7	81.0±2.1	77.5±2.1	84.7±0.7
	Squ.	72.6±2.4	70.9±1.8	74.6±2.2	61.7±2.1	81.5±1.3	78.8±1.9	82.4±1.2	80.9±1.6	77.7±1.3	83.9±0.9
	Ber.	72.1±2.3	71.0±2.0	74.7±2.1	61.1±1.7	80.9±1.6	79.1±1.1	82.1±0.9	80.8±1.7	77.8±1.5	84.2±0.4
CIFAR-10	Lin.	28.4±1.7	22.6±3.3	17.2±4.1	27.4±4.2	66.1±1.2	53.6±2.1	64.8±2.2	65.7±2.1	51.3±3.1	68.0±1.7
	Sin.	28.9±2.1	23.1±3.1	18.3±4.5	27.1±3.7	66.7±1.2	53.3±2.1	64.7±1.9	65.1±1.8	51.6±4.1	67.7±1.7
	Squ.	25.9±1.7	21.6±2.6	15.8±3.7	25.1±2.2	67.1±1.5	54.3±1.3	64.8±1.2	65.7±1.6	52.3±3.5	66.4±1.6
	Ber.	26.1±2.1	21.7±2.9	16.1±4.1	24.5±2.5	67.0±1.3	54.1±1.3	64.3±1.7	65.4±1.4	52.6±2.8	67.1±1.3
(ii) Covariate Shift Scenarios											
CIFAR-10	Lin.	15.4±0.4	10.1±0.8	39.1±0.5	37.8±1.1	41.6±0.4	35.2±0.5	58.0±0.6	58.1±0.6	36.7±2.4	59.5±0.5
	Sin.	15.2±0.5	11.9±0.7	39.2±0.4	37.2±0.9	41.7±0.6	33.1±0.4	57.1±0.5	57.8±0.7	35.3±1.7	59.3±0.5
CIFAR-10-C	Squ.	13.5±0.5	13.1±0.6	38.3±0.4	36.8±1.2	43.1±0.5	32.1±0.7	56.9±0.9	57.3±1.1	35.6±2.1	60.4±0.5
	Ber.	14.8±0.5	12.2±1.1	38.7±0.4	36.6±1.1	43.2±0.5	33.0±0.5	56.7±0.7	67.1±1.0	46.1±2.7	60.3±0.4
CIFAR-100	Lin.	4.2±1.7	2.4±0.2	29.1±0.1	31.3±0.3	20.2±0.4	20.7±0.6	42.0±0.4	40.8±0.2	37.1±1.1	43.1±0.5
	Sin.	5.6±2.2	2.5±0.1	29.5±0.3	30.7±0.8	20.3±0.5	20.7±0.6	41.9±0.4	40.9±0.3	36.6±1.3	44.4±0.4
CIFAR-100-C	Squ.	5.7±2.1	3.2±0.3	27.1±0.4	29.1±0.7	22.9±0.4	18.6±0.8	41.6±0.3	41.2±0.5	37.1±1.5	45.1±0.5
	Ber.	6.1±1.8	3.1±0.4	27.8±0.4	29.6±0.8	21.8±0.8	19.3±0.8	40.3±0.5	41.4±0.5	34.1±1.4	44.8±0.6

H.2. Impact of Distribution of Pre-training Data

In practical federated deployments, the exact distribution of the pre-training data on the server is often unknown to clients. To examine the robustness of Fed-ADE under such realistic conditions, we evaluate post-adaptation performance on image benchmarks when the pre-training data follows two non-uniform distributions: *Gaussian* and *Exponential Decay*. As shown in Table 9, and Table 10, Fed-ADE maintains consistently strong performance across all distribution settings and shift types, with only a minor variation compared to the uniform case (cf. Table 1). These results demonstrate that our dynamics-based adaptation effectively generalizes even when the pre-training distribution differs from the assumed prior, validating its robustness under practical deployment scenarios.

H.3. Why Cosine Similarity? Alternatives and Ablations

In this paper, our Fed-ADE must be (i) online & per-client (no extra communication), (ii) label-free (only pseudo-labels/representations), (iii) robust to pseudo-label noise and transient shifts, (iv) lightweight in compute and memory, and (v) bounded and interpretable so that we can stably map it to a learning rate in $[\eta_{\min}, \eta_{\max}]$ across heterogeneous clients.

Cosine similarity meets these requirements

- **Bounded:** S_c^t is bounded in $[0,1]$, so it leads to a stable mapping to learning across clients and timestep
- **Efficient:** Linear-time, no optimization/probabilistic inference
- **Robust:** Scale-invariant and less sensitive to sparsity/outliers in pseudo-label histograms; temperature rescaling or mild miscalibration does not change directions.

Alternatives considered We implemented three commonly suggested alternatives and found consistent drawbacks in our setting.

(1) KL divergence on batch histograms: *Pros:* classical information-theoretic distance. *Cons:* asymmetric and *unbounded*; becomes unstable/infinite under zeros/sparsity—common with noisy pseudo-labels on small batches. Requires smoothing/regularization and hyperparameters (e.g., ϵ -floor), which are highly dataset/client dependent and harm cross-client consistency. Mapping an unbounded statistic to a stable learning rate is delicate.

(2) Wasserstein distance: Requires a ground metric on classes or high-dim features; even with Sinkhorn regularization, it needs iterative optimization (per batch pair) and careful tuning. *Pros:* captures geometric discrepancies. *Cons:* higher compute

Table 10. Performance comparison of different online distribution shift adaptation methods across datasets and shift types with pretrained model from exponential decay distribution(average accuracy (%) and average wall time (sec.))

Dataset	Shift	Localized Learning				Federated Learning					Fed-ADE
		FTH	ATLAS	UNIDA	UDA	Fed-POE	FedCCFA	FixLR(Low)	FixLR(Mid)	FixLR(High)	
(i) Label Shift Scenarios											
Tiny ImageNet	Lin.	78.2±1.0	70.2±2.1	77.3±0.8	69.0±0.6	79.9±0.3	82.1±0.4	81.8±0.5	82.5±0.3	81.4±0.2	83.1±0.4
	Sin.	77.9±0.8	71.4±1.6	77.1±0.6	69.4±0.5	80.4±0.6	81.5±0.4	82.1±0.2	82.7±0.5	81.3±0.3	82.7±0.5
	Squ.	77.2±0.8	71.1±1.8	76.9±0.6	69.4±0.5	80.9±0.5	81.2±0.3	82.3±0.4	82.1±0.4	80.1±0.3	82.6±0.3
	Ber.	78.2±1.1	70.7±1.6	76.6±0.7	69.7±0.4	80.4±0.5	80.4±0.3	81.1±0.4	82.1±0.2	80.4±0.3	81.8±0.4
CIFAR-10	Lin.	27.7±1.4	29.0±2.0	20.2±1.3	27.6±1.8	65.1±1.0	63.8±0.9	64.1±1.4	64.8±1.3	60.1±1.3	66.2±1.7
	Sin.	28.4±1.2	28.3±2.0	20.3±1.1	28.1±1.6	64.7±0.6	63.1±1.3	63.8±1.2	64.3±1.4	60.3±0.8	65.9±1.7
	Squ.	27.1±1.2	28.5±1.8	21.1±0.9	24.8±1.6	64.6±0.9	62.8±1.1	64.5±1.1	64.1±1.6	61.1±0.9	66.4±1.5
	Ber.	26.6±1.5	27.1±1.9	20.7±1.1	25.1±2.1	64.3±1.1	62.3±0.7	64.7±1.4	63.8±1.5	60.8±1.2	67.1±1.3
(ii) Covariate Shift Scenarios											
CIFAR-10	Lin.	23.7±0.2	13.9±0.2	43.3±0.6	45.7±0.8	44.5±0.8	43.1±0.5	63.4±0.2	63.9±0.3	40.6±2.1	58.3±0.5
	Sin.	22.9±0.2	13.2±0.2	44.8±0.1	43.5±0.3	44.7±1.2	43.3±0.4	63.5±0.3	63.7±0.3	38.1±3.1	58.2±0.5
CIFAR-10-C	Squ.	23.8±0.1	14.1±0.4	42.2±0.1	43.3±0.2	48.5±1.2	41.6±0.7	62.7±2.1	64.5±2.1	39.6±2.8	59.0±0.5
	Ber.	23.6±0.2	14.2±0.3	42.7±0.2	42.3±0.1	48.7±0.9	42.0±0.5	64.1±2.1	64.4±2.1	40.8±2.2	59.0±0.4
CIFAR-100	Lin.	5.1±1.8	6.4±0.4	24.4±0.5	27.6±0.5	19.1±2.1	17.9±1.1	40.8±0.5	39.1±0.6	37.7±1.2	42.7±0.5
	Sin.	5.3±1.7	5.1±0.4	24.3±0.6	28.1±0.5	20.1±1.9	18.3±1.1	40.7±0.6	38.6±0.7	37.6±1.2	42.4±0.6
CIFAR-100-C	Squ.	4.8±1.5	4.7±0.5	24.1±0.4	28.3±0.7	19.8±2.1	18.1±0.9	39.4±0.6	38.2±0.7	38.1±1.1	42.1±0.4
	Ber.	4.3±1.3	5.1±0.3	23.5±0.4	27.8±0.6	19.7±2.3	16.9±1.0	40.1±0.4	37.7±0.6	37.8±1.2	42.0±0.5

219 and memory; sensitive to transient fluctuations; scale depends on ground metric \Rightarrow learning rate scaling across clients becomes
 220 ill-conditioned. In high-dim feature space, pairwise costs/transport exacerbate overhead.
 221 **(3) Bayesian change-point detection:** *Pros:* probabilistic treatment of regime switches. *Cons:* assumes i.i.d. within segments
 222 and requires prior modeling; posterior updates/inference add computation; fragile to batch-level noise (false alarms/delays).
 223 Typically designed for detection (binary), not for per-step *continuous* shift magnitudes needed to drive learning rate.

224 **Performance comparison of cosine similarity vs. alternatives** We compared cosine similarity with KL, Wasserstein, and
 225 Bayesian CPD under four dynamic shift patterns (Lin./Sin./Squ./Ber.) on CIFAR-10. Cosine consistently achieved the best
 226 accuracy:

Measure	Lin.	Sin.	Squ.	Ber.
Fed-ADE (Cosine similarity)	73.8 ± 0.6	73.6 ± 0.5	72.2 ± 1.6	72.9 ± 2.2
KL divergence	63.2 ± 3.8	62.4 ± 2.3	61.1 ± 3.4	61.4 ± 3.6
Wasserstein	70.8 ± 0.6	69.9 ± 1.1	68.8 ± 1.3	66.5 ± 2.1
Bayesian CPD	71.5 ± 0.5	70.8 ± 0.8	69.7 ± 2.4	69.0 ± 1.1

228 Gains of the cosine similarity method are largest under label noise and transient shifts, where bounded, direction-only
 229 comparison stabilizes learning rate scaling across clients.
 230