

## 1509 E Theoretical Analysis

1510 This appendix provides theoretical analysis and derivation details to  
 1511 supplement the methodology section. Please note that the meanings  
 1512 of some symbols differ from those in the main text.

### 1514 E.1 Component I - Data Volume Estimation 1515 Module

1516 We provide a rigorous theoretical analysis of our module from the  
 1517 perspectives of both privacy and utility.

1519 *E.1.1 Privacy Guarantee Analysis.* We begin with the formal defi-  
 1520 nition of Local Differential Privacy (LDP).

1521 **DEFINITION 1 (ε-LOCAL DIFFERENTIAL PRIVACY).** *A randomized  
 1522 algorithm  $\mathcal{A}$  satisfies  $\epsilon$ -LDP if for any two inputs  $x, x'$  in its domain,  
 1523 and for any subset of outputs  $O$  in its range, the following inequality  
 1524 holds:*

$$1526 \Pr[\mathcal{A}(x) \in O] \leq e^\epsilon \cdot \Pr[\mathcal{A}(x') \in O] \quad (28)$$

1527 Based on this definition, we can prove that our proposed module  
 1528 satisfies LDP in every round.

1529 **THEOREM 1 (PRIVACY GUARANTEE OF THE MODULE).** *Our pro-  
 1530 posed data size estimation mechanism ensures that the estimate  $n'_j$   
 1531 reported by client  $j$  in any round  $t$  satisfies  $\epsilon_d$ -LDP.*

1532 **PROOF.** Let the entire client-side procedure be represented by a  
 1533 randomized algorithm  $\mathcal{A}$ , which takes the true data size  $n_j$  as input  
 1534 and produces the reported estimate  $n'_j$  as output. We analyze the  
 1535 privacy properties of this algorithm.

1536 **Sensitivity Analysis:** The first step of the algorithm is to clip  
 1537 the input  $n_j$  (Eq. (??)). For any two possible inputs  $n_j$  and  $n_j^*$ , the  
 1538 L1 distance between their clipped counterparts,  $n_{j,clip}$  and  $n_{j,clip}^*$ ,  
 1539 defines the global sensitivity  $\Delta f$ :

$$1540 \Delta f = \sup_{n_j, n_j^*} |n_{j,clip} - n_{j,clip}^*| \leq (c_d - 1) - 0 = c_d - 1 \quad (29)$$

1541 We use its upper bound,  $c_d$ , as the global sensitivity.

1542 **Laplace Mechanism Analysis:** The core privacy-preserving  
 1543 step is the addition of Laplace noise (Eq. (??)). By the standard  
 1544 definition of the Laplace mechanism, when the noise scale param-  
 1545 eter is set to  $b_d = \Delta f / \epsilon_d = c_d / \epsilon_d$ , the output  $noisy\_d_j^t$  satisfies  
 1546  $\epsilon_d$ -differential privacy with respect to its input  $n_{j,clip}$ . For any two  
 1547 inputs  $n_{j,clip}$  and  $n_{j,clip}^*$ , and any possible output value  $v$ , the ratio  
 1548 of their probability density functions is:

$$1549 \frac{p(v|n_{j,clip})}{p(v|n_{j,clip}^*)} = \frac{\frac{1}{2b_d} \exp(-\frac{|v-n_{j,clip}|}{b_d})}{\frac{1}{2b_d} \exp(-\frac{|v-n_{j,clip}^*|}{b_d})} \quad (30)$$

$$1550 = \exp\left(\frac{|v-n_{j,clip}^*| - |v-n_{j,clip}|}{b_d}\right)$$

1551 By the reverse triangle inequality,  $|v-n_{j,clip}^*| - |v-n_{j,clip}| \leq |n_{j,clip} - n_{j,clip}^*| \leq c_d$ . Therefore:

$$1552 \frac{p(v|n_{j,clip})}{p(v|n_{j,clip}^*)} \leq \exp\left(\frac{|n_{j,clip} - n_{j,clip}^*|}{b_d}\right) \leq \exp\left(\frac{c_d}{c_d/\epsilon_d}\right) = e^{\epsilon_d} \quad (31)$$

1553 This demonstrates that the generation of  $noisy\_d_j^t$  satisfies  $\epsilon_d$ -DP.

1554 **Immunity to Post-processing:** A fundamental property of  
 1555 differential privacy is its immunity to post-processing. In our al-  
 1556 gorithm, the EMA smoothing step (Eq. (??)) and the assignment of  
 1557 the result to  $n'_j$  are computations performed on the value  $noisy\_d_j^t$ ,  
 1558 which already satisfies  $\epsilon_d$ -DP. Since these subsequent operations do  
 1559 not depend on the original private data  $n_j$ , they do not compromise  
 1560 the existing privacy guarantee.  $\square$

1561 In summary, the entire client-side reporting procedure satisfies  
 1562  $\epsilon_d$ -LDP.

1563 **REMARK 1 (CUMULATIVE PRIVACY BUDGET).** *It is important to note  
 1564 that  $\epsilon_d$  represents the privacy budget for a single round. If a client  
 1565 participates in  $T$  consecutive rounds, its total privacy cost, according  
 1566 to the principle of Sequential Composition in differential privacy, is  
 1567 bounded by  $T \cdot \epsilon_d$ . In practice, a sufficiently small  $\epsilon_d$  can be chosen to  
 1568 ensure that the cumulative privacy leakage over the entire FL process  
 1569 remains within an acceptable range.*

1570 **E.1.2 Utility Analysis: Unbiasedness and Order-Preserving Property.**  
 1571 We demonstrate the utility of our estimator by proving that its  
 1572 expectation is unbiased.

1573 **THEOREM 2 (ASYMPTOTIC UNBIASEDNESS OF THE ESTIMATOR).**  
 1574 *Assuming the smoothing factor  $\gamma_1 \in (0, 1)$  and that the client's true  
 1575 clipped data volume  $n_{j,clip}$  remains stable, the mathematical expecta-  
 1576 tion of the smoothed estimate reported by client  $j$ ,  $\mathbb{E}[s_j^t]$ , asymptoti-  
 1577 cally converges to its true clipped data volume  $n_{j,clip}$ :*

$$1578 \lim_{t \rightarrow \infty} \mathbb{E}[s_j^t] = n_{j,clip} \quad (32)$$

1579 **PROOF.** We begin by analyzing the expectation of the single-  
 1580 round noisy observation  $noisy\_d_j^t$ . The expectation of the Laplace  
 1581 distribution  $\text{Laplace}(0, b_d)$  is zero. Therefore, by the linearity of  
 1582 expectation:

$$1583 \mathbb{E}[noisy\_d_j^t] = \mathbb{E}[n_{j,clip} + noise_j^t] \quad (33)$$

$$1584 = \mathbb{E}[n_{j,clip}] + \mathbb{E}[noise_j^t]$$

$$1585 = n_{j,clip} + 0$$

$$1586 = n_{j,clip}$$

1587 This shows that the noisy observation in each round is, by itself,  
 1588 an unbiased estimator of the true clipped value.

1589 Next, we analyze the expectation of the smoothed value  $s_j^t$ . Ap-  
 1590 plying the linearity of expectation to the equation:

$$1591 \mathbb{E}[s_j^t] = \mathbb{E}[\gamma_1 \cdot s_j^{t-1} + (1 - \gamma_1) \cdot noisy\_d_j^t] \quad (34)$$

$$1592 = \gamma_1 \cdot \mathbb{E}[s_j^{t-1}] + (1 - \gamma_1) \cdot \mathbb{E}[noisy\_d_j^t]$$

$$1593 = \gamma_1 \cdot \mathbb{E}[s_j^{t-1}] + (1 - \gamma_1) \cdot n_{j,clip}$$

1594 This is a first-order linear recurrence relation for  $\mathbb{E}[s_j^t]$ . To find  
 1595 its general solution, we unroll the recurrence, using the initial  
 1596 value  $\mathbb{E}[s_j^0] = n_{j,clip}$ :

$$\begin{aligned}
1625 \quad & \text{condition } E[s_j^0] = E[0] = 0: \\
1626 \quad & \mathbb{E}[s_j^t] = \gamma_1 \mathbb{E}[s_j^{t-1}] + (1 - \gamma_1) n_{j,\text{clip}} \\
1627 \quad & = \gamma_1 \left( \gamma_1 \mathbb{E}[s_j^{t-2}] + (1 - \gamma_1) n_{j,\text{clip}} \right) + (1 - \gamma_1) n_{j,\text{clip}} \\
1628 \quad & = \gamma_1^2 \mathbb{E}[s_j^{t-2}] + \gamma_1 (1 - \gamma_1) n_{j,\text{clip}} + (1 - \gamma_1) n_{j,\text{clip}} \\
1629 \quad & = \dots \\
1630 \quad & = \gamma_1^t \mathbb{E}[s_j^0] + \left( \sum_{i=0}^{t-1} \gamma_1^i \right) (1 - \gamma_1) n_{j,\text{clip}} \\
1631 \quad & = \gamma_1^t \cdot 0 + \left( \frac{1 - \gamma_1^t}{1 - \gamma_1} \right) (1 - \gamma_1) n_{j,\text{clip}} \\
1632 \quad & = (1 - \gamma_1^t) n_{j,\text{clip}}
\end{aligned} \tag{35}$$

We now take the limit of the above expression as  $t \rightarrow \infty$ . Since the smoothing factor  $\gamma_1 \in (0, 1)$ , it follows that  $\lim_{t \rightarrow \infty} \gamma_1^t = 0$ . Therefore:

$$\lim_{t \rightarrow \infty} \mathbb{E}[s_j^t] = \lim_{t \rightarrow \infty} (1 - \gamma_1^t) n_{j,\text{clip}} = (1 - 0) \cdot n_{j,\text{clip}} = n_{j,\text{clip}} \tag{36}$$

This result proves that the expectation of the smoothed estimate asymptotically converges to the true clipped data volume.  $\square$

**REMARK 2 (STATISTICAL ORDER-PRESERVING PROPERTY).** *The Theorem above serves as the theoretical cornerstone for the utility of our module. It reveals that although random noise is introduced in every round, the system effectively filters out the mean effect of this noise through temporal smoothing. Consequently, the set of estimates reported by the clients,  $\{n'_j\}_{j \in \mathcal{M}}$ , will, in statistical expectation, faithfully reflect the true ranking and relative magnitudes of each client's clipped data volume. This ensures that the final score vector  $\mathbf{S}_\text{data}^t$  computed by the server is a meaningful and robust input, providing a solid foundation for subsequent FL algorithms that rely on client contribution assessment.*

## E.2 Theoretical Analysis of Component I - Contribution Estimation Module

**E.2.1 Privacy Guarantee Analysis.** We first present the relevant theorem on the composition of differential privacy, and then leverage it to prove the privacy guarantee of our module.

**DEFINITION 2 (BASIC COMPOSITION THEOREM).** *Let there be a series of privacy mechanisms  $\mathcal{B}_1, \dots, \mathcal{B}_k$ , where each  $\mathcal{B}_i$  provides  $\epsilon_i$ -DP. The sequence of their outputs,  $(\mathcal{B}_1(D), \dots, \mathcal{B}_k(D))$ , when applied to the same input database  $D$ , satisfies  $(\sum_{i=1}^k \epsilon_i)$ -DP.*

**THEOREM 3 (PRIVACY GUARANTEE OF THE CONTRIBUTION MODULE).** *The contribution estimate  $v_j$  reported by client  $j$  in any round  $t$  satisfies  $\epsilon_{\text{contrib}}$ -LDP.*

**PROOF.** We consider the entire client-side reporting procedure as a randomized algorithm  $\mathcal{B}$ , with the client's private dataset  $D_j$  as input and the public value  $v_j$  as output. This algorithm  $\mathcal{B}$  can be decomposed into two independent sanitization sub-mechanisms and a post-processing step.

- (1) **Privacy of Sub-mechanism  $\mathcal{B}_l$ :** This sub-mechanism processes the initial loss. Its input is  $l_j^t = \mathcal{L}(\mathbf{w}_{\text{start},j}^t; D_j)$ . The clipping operation in Eq. (5) bounds its L1 sensitivity to

$c_l$ . Subsequently, Eq. (6) adds Laplace noise with a scale of  $b_l = c_l/\epsilon_l$ . By definition of the Laplace mechanism, the output of sub-mechanism  $\mathcal{B}_l$ ,  $\text{noisy}_l$ , satisfies  $\epsilon_l$ -LDP.

- (2) **Privacy of Sub-mechanism  $\mathcal{B}_g$ :** This sub-mechanism processes the gradient norm. Its input is  $g_j^t = \nabla \mathcal{L}(\mathbf{w}_{\text{start},j}^t; D_j)$ . The norm clipping in Eq. (5) bounds its L2 sensitivity to  $c_g$ . The Laplace noise added in Eq. (6) has a scale of  $b_g = c_g/\epsilon_g$ . This operation ensures that the output of sub-mechanism  $\mathcal{B}_g$ ,  $\text{noisy}_g$ , satisfies  $\epsilon_g$ -LDP.

- (3) **Applying Composition:** The client effectively reports a tuple  $(\text{noisy}_l, \text{noisy}_g)$  to the server, as the final value  $v_j$  is computed from these two components. Since both  $\text{noisy}_l$  and  $\text{noisy}_g$  are derived from the same private dataset  $D_j$ , we must apply the Basic Composition Theorem from **Definition 2**. Viewing  $\mathcal{B}_l$  and  $\mathcal{B}_g$  as two privacy mechanisms acting on the same data  $D_j$ , the composite mechanism that releases both of their results satisfies  $(\epsilon_l + \epsilon_g)$ -LDP.

- (4) **Post-processing Invariance:** The final reported value  $v_j$  is the result of a deterministic linear transformation (Eq. (7)) on the private tuple  $(\text{noisy}_l, \text{noisy}_g)$ . A fundamental property of differential privacy is its immunity to post-processing, which states that performing any computation on the output of a DP mechanism, without accessing the original private data, does not weaken the privacy guarantee. Therefore, the computation of  $v_j$  inherits the privacy level of its input tuple.

In conclusion, by combining steps 3 and 4, the entire algorithm  $\mathcal{B}$  satisfies  $(\epsilon_l + \epsilon_g)$ -LDP. As per our hyperparameter setting  $\epsilon_{\text{contrib}} = \epsilon_l + \epsilon_g$ , the output  $v_j$  of our module strictly satisfies  $\epsilon_{\text{contrib}}$ -LDP.  $\square$

**E.2.2 Utility and Rationale of Contribution Proxies.** The utility of this module is founded on the rationale that the two chosen proxy metrics can reasonably represent a client's potential contribution. We articulate this rationale through the following propositions.

**PROPOSITION 1 (GRADIENT NORM AS A PROXY FOR POTENTIAL CONTRIBUTION).** *The gradient norm,  $\|g_j^t\|_2 = \|\nabla \mathcal{L}(\mathbf{w}_{\text{start},j}^t; D_j)\|_2$ , quantifies the local rate of change of the loss function, defined by the client's data  $D_j$ , at the current model parameter point. A large gradient norm implies that the loss surface is steep at that point, which intuitively suggests a significant misalignment between the current model  $\mathbf{w}_{\text{start},j}^t$  and the optimal parameters for client  $j$ 's data. Therefore, the client's data holds substantial information for correcting the global model, indicating a high potential contribution.*

**PROPOSITION 2 (INITIAL LOSS AS A PROXY FOR IMPROVEMENT HEADROOM).** *The initial loss,  $l_j^t = \mathcal{L}(\mathbf{w}_{\text{start},j}^t; D_j)$ , directly quantifies the performance of the current model on the client's local data. A high initial loss value indicates poor performance, which is synonymous with a large headroom for improvement. Consequently, such a client has the potential to achieve a more significant loss reduction during its local training, thereby making a greater tangible contribution to the global model's performance enhancement.*

**REMARK 3 (HOLISTIC CONTRIBUTION METRIC).** *Relying on either proxy alone could lead to sub-optimal selections. For instance, anomalous data could produce an extremely large gradient norm,*

1741 while low-quality data might also result in a very high initial loss.  
 1742 Our module, by taking a weighted sum of the two noisy proxies,  
 1743  $v_j = \eta_l \cdot \text{noisy\_}l_j + \eta_g \cdot \text{noisy\_}g_j$ , aims to strike a balance. It seeks  
 1744 to identify clients that not only present a discrepancy with the current  
 1745 model (proxied by the gradient norm) but also on which the  
 1746 model can effectively learn from this discrepancy to achieve significant  
 1747 performance gains (proxied by the initial loss). The weighting  
 1748 hyperparameters  $\eta_l$  and  $\eta_g$  provide the necessary flexibility to tune  
 1749 this trade-off between the two dimensions.

1750

### 1751 E.3 Component II - Hedonic Game-Based Stable 1752 Coalition Formation

1753 The core goal of this section is to prove that our proposed online  
 1754 estimators  $\hat{\mu}_e$  and  $\hat{\sigma}^2$  are effective approximations of the prior  
 1755 parameters  $\mu_e$  and  $\sigma^2$  in Donahue K. theory, thereby proving the  
 1756 rationality of the dynamic game threshold  $\hat{T}$ .

1757

#### 1758 E.3.1 The validity of $\hat{\mu}_e$ .

1759

1760 ASSUMPTION 1 (THEORETICAL BASIS AND ERROR DECOMPOSITION). According to the work of Donahue K. et al., there is a set  
 1761 of  $\mathcal{M}$ . The real data of each client  $j$  is determined by an unknown  
 1762  $D$ -dimensional parameter vector  $\theta_j^*$  and a noise variance  $\sigma_{\epsilon,j}^2$ . These  
 1763 parameters are drawn independently and identically (IID) from a  
 1764 higher-level prior distribution. We define the following two core the-  
 1765 oretical prior parameters, the Mean Irreducible Error and Inherent  
 1766 Heterogeneity:

1767

$$\mu_e := \mathbb{E}_j[\sigma_{\epsilon,j}^2], \quad \sigma^2 := \text{Var}_j(\theta_j^*). \quad (37)$$

1768

1769 ASSUMPTION 2 (LOCAL MODEL CONVERGENCE). We assume that  
 1770 in an effective FL process, for any client  $j$ , its local model  $\theta_j^t$  ob-  
 1771 tained after  $t$  rounds of training will converge probabilistically or  
 1772 mean-squared to the optimal model  $\theta_j^*$  on its local data distribution  
 1773  $D_j$ . This assumption means that as training progresses, the model  
 1774 will fully learn generalizable patterns in the local data. We assume  
 1775 that  $\lim_{t \rightarrow \infty} \theta_j^t = \theta_j^*$  is a reasonable simplification under ideal con-  
 1776 vergence conditions, which captures the core dynamics of the local  
 1777 model constantly approaching the optimal representation of its data  
 1778 distribution.

1779

1780 THEOREM 4 (ASYMPTOTICALLY UNBIASED ESTIMATION). Under  
 1781 the condition that Assumption 2 holds, our estimator  $\hat{\mu}_e^t$  is an asymp-  
 1782 tootically unbiased estimator of the theoretical parameter  $\mu_e$ , that is:

1783

$$\lim_{t \rightarrow \infty} \mathbb{E}[\hat{\mu}_e^t] = \mu_e \quad (38)$$

1784

1785 PROOF. According to the setting of Donahue K. et al.,  $\mu_e$  is the  
 1786 expected irreducible error of the system average. For any client  $j$ ,  
 1787 the expectation of its loss function  $\mathbb{E}[L_j]$  can be decomposed. Let  
 1788  $f_j^*$  be the best possible function (i.e., Bayesian optimal classifier)  
 1789 on the data distribution  $D_j$  of client  $j$ ,  $\theta_j^t$  be the model parameters  
 1790 after  $t$  rounds of training, and its corresponding function is  $f(\theta_j^t, \cdot)$ .

1791

1792 Then the expected loss of  $t$  rounds is:

$$\begin{aligned} \mathbb{E}[L_j^t] &= \mathbb{E}_{(x,y) \sim D_j} \left[ (f(\theta_j^t; x) - y)^2 \right] \\ &= \mathbb{E}_{(x,y) \sim D_j} \left[ (f(\theta_j^t; x) - f_j^*(x) + f_j^*(x) - y)^2 \right] \\ &= \mathbb{E}_x \left[ (f(\theta_j^t; x) - f_j^*(x))^2 \right] + \mathbb{E}_{(x,y)} \left[ (f_j^*(x) - y)^2 \right] \\ &= \text{Error}_{\text{reducible}}(\theta_j^t) + \text{Error}_{\text{irreducible},j}, \end{aligned} \quad (39)$$

1793 where the cross term  $\mathbb{E}[(f(\theta_j^t; x) - f_j^*(x))(f_j^*(x) - y)]$  is zero since  
 1794  $f_j^*(x)$  is the conditional expectation  $\mathbb{E}[y|x]$  given  $x$ .  $\text{Error}_{\text{irreducible},j} =$   
 1795  $\mathbb{E}[(f_j^*(x) - y)^2] = \sigma_{\epsilon,j}^2$  is the noise variance inherent to client  $j$ ,  
 1796 and the theoretical parameter  $\mu_e = \mathbb{E}_j[\sigma_{\epsilon,j}^2]$ .

1797 ASSUMPTION 3 (CONVERGENCE OF ESTIMATORS - 1). We assume  
 1798 that in an effective FL process, the model parameters  $\theta_j^t$  of client  $j$   
 1799 converge to their local optimal solution  $\theta_j^*$ , so that the reducible error  
 1800 converges to a minimum value, ideally zero:

$$\lim_{t \rightarrow \infty} \text{Error}_{\text{reducible}}(\theta_j^t) = \lim_{t \rightarrow \infty} \mathbb{E}[(f_j(\theta_j^t, x) - f_j^*(x))^2] = 0. \quad (40)$$

1801 This means that the model has been fully learned in non-convex  
 1802 optimization, and the error caused by model inaccuracy cannot be  
 1803 significantly reduced through further training.

1804 Based on this assumption, we can have:

1805 COROLLARY 1. Substituting this limit into the loss decomposition,  
 1806 we obtain that the expectation of the local loss will converge to the  
 1807 local irreducible error in the later stages of training:

$$\lim_{t \rightarrow \infty} \mathbb{E}[L_j^t] = \lim_{t \rightarrow \infty} \left( \text{Error}_{\text{reducible}}(\theta_j^t) + \sigma_{\epsilon,j}^2 \right) = \sigma_{\epsilon,j}^2. \quad (41)$$

1808 The instantaneous observation value  $\mu_{e,\text{obs}}^t$  we calculated in  
 1809 round  $t$  is the mean of  $L_j^{t-1}$  on the current participating client  
 1810 set  $\mathcal{N}^{t-1}$ :

$$\mu_{e,\text{obs}}^t = \frac{1}{|\mathcal{N}^{t-1}|} \sum_{j \in \mathcal{N}^{t-1}} L_j^{t-1}. \quad (42)$$

1811 If we assume that the client sampling process is unbiased, that  
 1812 is,  $E_{\mathcal{N}^{t-1}}[\cdot] = \mathbb{E}_j[\cdot]$ , then its expectation is:

$$\mathbb{E}[u_{e,\text{obs}}^t] = \mathbb{E} \left[ \frac{1}{|\mathcal{N}^{t-1}|} \sum_{j \in \mathcal{N}^{t-1}} L_j^{t-1} \right] = \mathbb{E}_j[L_j^{t-1}]. \quad (43)$$

1813 Therefore, when  $t \rightarrow \infty$ :

$$\lim_{t \rightarrow \infty} \mathbb{E}[u_{e,\text{obs}}^t] = \mathbb{E}_j \left[ \lim_{t-1 \rightarrow \infty} \mathbb{E}[L_j^{t-1}] \right] = \mathbb{E}_j[\sigma_{\epsilon,j}^2] = \mu_e. \quad (44)$$

1814 The EMA update mechanism  $\hat{\mu}_e^t = \gamma_2 \cdot \hat{\mu}_e^{t-1} + (1 - \gamma_2) \cdot \mu_{e,\text{obs}}^t$   
 1815 can be expanded into a weighted sum of all historical observations:

$$\hat{\mu}_e^t = (1 - \gamma_2) \sum_{i=0}^{t-1} \gamma_2^i \mu_{e,\text{obs}}^{t-i} + \gamma_2^t \hat{\mu}_e^0. \quad (45)$$

1816 When  $t \rightarrow \infty$ ,  $\hat{\mu}_e^t$  will converge to the expected value of  $\mu_{e,\text{obs}}^t$ . In  
 1817 the steady state,  $\mathbb{E}[\hat{\mu}_e^t] \approx \mathbb{E}[\mu_{e,\text{obs}}^t]$ . Therefore, our final estimator  
 1818  $\hat{\mu}_e$  is an asymptotically unbiased and consistent estimate of the  
 1819 theoretical parameter  $\mu_e$ , and its validity is guaranteed.  $\square$

### 1857 E.3.2 The validity of $\widehat{\sigma}^2$ .

1858 **THEOREM 5 (THEORETICAL BASIS AND ESTIMATOR).** *Under the*  
 1859 *condition that Assumption 2 holds, estimator  $\widehat{\sigma}^2$  is an asymptotically*  
 1860 *unbiased estimator of the theoretical parameter  $\sigma^2$ , that is,*  
 1861  $\lim_{t \rightarrow \infty} \mathbb{E}[\widehat{\sigma}^2] = \sigma^2$ .

1862 **PROOF.** The theoretical parameter  $\sigma^2$  is the variance between  
 1863 the optimal local models  $\theta_j^*$ :

$$1864 \sigma^2 = \text{Var}_j(\theta_j^*) = \mathbb{E}_j[\|\theta_j^* - \mathbb{E}_j[\theta_j^*]\|^2]. \quad (46)$$

1865 Instantaneous observation  $\sigma_{\text{obs},t}^2$  is the unbiased sample variance of  
 1866 the model  $\theta_j^{t-1}$  in round  $t-1$ , and its expectation is:

$$1867 \mathbb{E}[\sigma_{\text{obs},t}^2] = \mathbb{E}\left[\frac{1}{|\mathcal{N}^{t-1}| - 1} \sum_{j \in \mathcal{N}^{t-1}} \|\theta_j^{t-1} - \bar{\theta}_{\text{obs}}^{t-1}\|^2\right] = \text{Var}_j(\theta_j^{t-1}) \quad (47)$$

1868 **ASSUMPTION 4 (CONVERGENCE OF ESTIMATORS - 2).** *According to*  
 1869 *Assumption 2, each local model  $\theta_j^t$  converges to its own ideal model*  
 1870  $\theta_j^*$ :

$$1871 \lim_{t \rightarrow \infty} \theta_j^t = \theta_j^*. \quad (48)$$

1872 **COROLLARY 2 (CONVERGENT TRANSITIVITY).** *Variance  $\text{Var}(\cdot)$  is a*  
 1873 *continuous functional. According to the Continuous Mapping Theorem,*  
 1874 *the convergence of the variable can be transferred to its continuous*  
 1875 *function. According to Assumption 1, the definition of  $\text{Var}_j(\theta_j^*)$  is  $\sigma^2$ .*  
 1876 *Therefore:*

$$1877 \lim_{t \rightarrow \infty} \text{Var}_j(\theta_j^{t-1}) = \text{Var}_j\left(\lim_{t \rightarrow \infty} \theta_j^{t-1}\right) = \text{Var}_j(\theta_j^*) = \sigma^2 \quad (49)$$

1878 In Sec 4.1, our instantaneous heterogeneity observation  $\sigma_{\text{obs},t}^2$  is:

$$1879 \sigma_{\text{obs},t}^2 = \frac{1}{|\mathcal{N}^{t-1}| - 1} \sum_{j \in \mathcal{N}^{t-1}} \|\theta_j^{t-1} - \bar{\theta}_{\text{obs}}^{t-1}\|^2. \quad (50)$$

1880 This is a standard unbiased sample variance estimator, whose  
 1881 expectation is equal to the population variance:

$$1882 \mathbb{E}[\sigma_{\text{obs},t}^2] = \text{Var}_j(\theta_j^{t-1}), \quad (51)$$

1883 assuming client sampling is unbiased.

1884 Therefore, when  $t$  is large enough:

$$1885 \lim_{t \rightarrow \infty} \mathbb{E}[\sigma_{\text{obs},t}^2] = \lim_{t-1 \rightarrow \infty} \text{Var}_j(\theta_j^{t-1}) = \sigma^2. \quad (52)$$

1886 Similarly, EMA estimator

$$1887 \widehat{\sigma}^2 = \gamma_2 \cdot \widehat{\sigma}^{t-1} + (1 - \gamma_2) \cdot \sigma_{\text{obs},t}^2 \quad (53)$$

1888 will converge to the expected value of  $\sigma_{\text{obs},t}^2$ .

1889 Therefore:

$$1890 \lim_{t \rightarrow \infty} \widehat{\sigma}^2 = \lim_{t \rightarrow \infty} \mathbb{E}[\sigma_{\text{obs},t}^2] = \sigma^2, \quad (54)$$

1891 and this proves that  $\widehat{\sigma}^2$  is an asymptotically unbiased estimator of  
 1892 the theoretical parameter  $\sigma^2$ .  $\square$

1893 **REMARK 4 (WARM-UP AND EMA CONDITIONS).** *The above proofs*  
 1894 *all rely on the limit condition of  $t \rightarrow \infty$ . In the early stages of limited*  
 1895 *training, the convergence assumption does not actually hold, and  $L_j^t$*   
 1896 *and  $\theta_j^t$  contain huge noise and deviations.*

1897 The above proves the asymptotic properties of the estimator. In  
 1898 a finite number of training rounds, the initial estimates  $\widehat{\mu}_e^0$  and  $\widehat{\sigma}^2$   
 1899 obtained in the warm-up phase by batch averaging  $T_{\text{warmup}}$  rounds  
 1900 are, according to the law of large numbers, closer to their expected  
 1901 initial estimates than any single-round instantaneous values. This  
 1902 provides a robust "anchor point" for EMA, preventing it from being  
 1903 "biased" by extreme noise values in the early stages, thereby  
 1904 accelerating the convergence of the estimator to its true value. In  
 1905 EMA smoothing,  $\mu_e$  and  $\sigma^2$  are static properties that describe the  
 1906 data distribution, while  $\mu_{e,\text{obs}}^t$  and  $\sigma_{\text{obs},t}^2$  are dynamic observations  
 1907 affected by client sampling and training randomness. Through the  
 1908 smoothing factor  $\gamma_2$ , the system controls the trade-off between the  
 1909 estimated bias and variance, ensuring that throughout the dynamic  
 1910 process,  $\hat{T}^t = \widehat{\mu}_e^t / \widehat{\sigma}^2$  is a stable proxy that gradually approaches  
 1911 the true value.

1912 The mathematical logic of Sec 4.2 is rigorous. It transforms the  
 1913 two theoretical priors  $(\mu_e, \sigma^2)$  that are not directly observable into  
 1914 empirical quantities  $(\widehat{\mu}_e^t, \widehat{\sigma}^2)$  that can be dynamically estimated  
 1915 through  $L_j^t$  and  $\theta_j^t$  in the FL process. Its effectiveness is guaranteed  
 1916 by the basic convergence assumptions of machine learning and  
 1917 statistical estimation theory.

## 1918 E.4 Component III - Multi-Granularity Aggregation with Soft Update

1919 This section aims to prove that the aggregation strategy we designed  
 1920 for the coalitions is reasonable and effective, and the soft update  
 1921 mechanism is a necessary supplement to ensure the long-term con-  
 1922 vergence of the system. We omit the Uniform and Coarse-grained  
 1923 here, as related work has provided sufficient proof already.

1924 **E.4.1 Theoretical Basis of Fine-grained Aggregation Strategy.** The  
 1925 optimal weight  $v_{ji}$  of its Fine-grained federation is derived by min-  
 1926 imizing the expected MSE of player  $j$ . The solution (as shown in  
 1927 Lemma 7.1 of Donahue K.) is of a complex form and depends on  $\mu_e$ ,  
 1928  $\sigma^2$ , and  $n_i$  of all clients. The underlying logic is that client  $j$  should  
 1929 give higher weights to clients whose models  $\theta_i$  are "similar" to its  
 1930 own  $\theta_j$ , in order to minimize the introduced bias. We demonstrate  
 1931 that our Fine-grained aggregation strategy is more adaptable to  
 1932 dynamic systems in the following proof.

1933 **COROLLARY 3 (GRADIENT DIRECTION AS PROXY FOR DATA DIS-**  
 1934 **TRIBUTION).** *The direction of the model update vector  $\Delta\theta_j$  is a valid*  
 1935 *approximation of the average gradient direction  $\mathbb{E}_{x \sim D_j}[\nabla_{\theta} L_j(\theta)]$  of*  
 1936 *the local data distribution  $D_j$  of client  $j$  at the current model point*  
 1937  *$w_{\text{start},j}$ . Therefore, the cosine similarity between update vectors is an*  
 1938 *effective measure of local task similarity.*

1939 **PROOF.** For  $E$ -step local updates (learning rate is  $\eta$ ),  $\Delta\theta_j^t$  can be  
 1940 written as:

$$1941 \Delta\theta_j^t = \theta_j^t - w_{\text{start},j}^t = \sum_{e=1}^E \Delta\theta_{j,e}^t, \quad (55)$$

1942 where  $\Delta\theta_{j,e}^t$  is the update at step  $e$ . For SGD,

$$1943 \Delta\theta_{j,e}^t = -\eta \cdot g_j(\theta_{j,e-1}^t), \quad (56)$$

1944 where  $g_j$  is the gradient over a mini-batch. When  $\eta$  is small enough  
 1945 and  $E$  is not too large,  $\theta$  does not change much during the local  
 1946 update, and  $\theta_{j,e-1}^t \approx w_{\text{start},j}^t$ .

1973 Therefore, in FL, the model update  $\Delta\theta_j = \theta_j^t - w_{\text{start},j}^t$  can be  
 1974 viewed as an approximation of the average gradient direction driven  
 1975 by client  $j$ 's local data  $D_j$  in the model parameter space. Therefore,  
 1976  $\text{sim}(\Delta\theta_i, \Delta\theta_j)$  directly measures the consistency of the goals of the  
 1977 two client local learning tasks:  
 1978

$$1979 \Delta\theta_j^t \approx \sum_{e=1}^E -\eta \cdot \mathbb{E}_{\text{batch}_e \sim D_j} [g_j(w_{\text{start},j}^t)] = -E \cdot \eta \cdot \mathbb{E}_{x \sim D_j} [g_j(w_{\text{start},j}^t)]. \quad (57)$$

1982 This approximation shows that the direction of  $\Delta\theta_j^t$  is approxi-  
 1983 mately collinear with the local expected gradient direction  $\nabla_{\theta} L_j(w_{\text{start},j}^t)$ .  
 1984 Since the expected gradient is uniquely determined by the data dis-  
 1985 tribution  $D_j$ ,  $\text{sim}(\Delta\theta_i, \Delta\theta_j)$  is an approximation of  $\text{sim}(\nabla L_i, \nabla L_j)$ ,  
 1986 and is therefore a reasonable measure of the similarity of the tasks  
 1987 of  $D_i$  and  $D_j$  at the current point.

1988 Because our weight calculation formula is:

$$1989 \text{vec}_{j,k} = (1 + \Gamma_{k,j}^t) / \sum_{i \in C_f^t} (1 + \Gamma_{i,j}^t), \quad (58)$$

1992 where  $\Gamma_{k,j}^t = \text{sim}(\Delta\theta_j^t, \Delta\theta_k^t)$ .

1993 If the update direction of client  $k$  is highly consistent with that of  
 1994  $j$  ( $\Gamma_{k,j}^t \rightarrow 1$ ),  $\text{vec}_{j,k}$  gets a high weight. This is completely consistent  
 1995 with the spirit of "cooperating with similar parties" in Donahue K.'s  
 1996 theory. If the update direction of client  $k$  is orthogonal or opposite  
 1997 to that of  $j$  ( $\Gamma_{k,j}^t \leq 0$ ),  $\text{vec}_{j,k}$  gets a low weight. This also achieves the  
 1998 goal of "avoiding cooperation with those who conflict with tasks".  
 1999

□

2000 Based on the above, our Fine-grained aggregation strategy re-  
 2001 shapes Donahue K.'s theory in terms of implementation by using  
 2002 the cosine similarity of model updates as a similarity metric. It is  
 2003 not only consistent with the original work in terms of logical goals,  
 2004 but also has more advantages in mathematics: 1. Adaptability: No  
 2005 need to rely on global, static prior parameters  $\mu_e, \sigma^2$ , but adaptively  
 2006 adjust according to the dynamic learning behavior  $\Delta\theta$  of each round.  
 2007 2. Robustness: Directly operating "knowledge increment" rather  
 2008 than "model state" can better handle the situation where model pa-  
 2009 rameters vary greatly in non-IID scenarios, avoiding catastrophic  
 2010 forgetting or parameter conflicts.  
 2011

2012 **REMARK 5 (ADVANTAGES OF WEIGHTING PROPERTY).** *Donahue*  
 2013 *K.'s theory is based on the similarity of model parameters  $\theta_i$ , which*  
 2014 *can be misleading in non-IID settings (e.g., two models with similar*  
 2015 *functions may have parameters far apart). Our method is based on*  
 2016 *gradient direction and focuses on the consistency of learning objec-*  
 2017 *ctives, which is a more robust similarity measure in the non-convex*  
 2018 *optimization landscape.*

2019 **E.4.2 Effectiveness of Soft Update Mechanism.** Since the theoreti-  
 2020 cal basis of Donahue K. is limited and based on a single-round  
 2021 static game, it does not involve the evolution path of the model in  
 2022 multiple rounds of iterations. Therefore, we introduce a soft update  
 2023 mechanism after multi-granularity aggregation.

2024 **THEOREM 6 (BOUNDED DIVERGENCE VIA SOFT UPDATE).** *In fully*  
 2025 *personalized FL, introducing a global reference model  $w_{\text{ref}}$  as the*  
 2026 *starting point for training some clients can provide an upper bound*  
 2027 *for the variance  $\text{Var}_j(\theta_j^t)$  between client models, thereby prevent-*  
 2028 *ing the system from diverging due to knowledge drift.*

2029 **PROOF. Divergence risk of unconstrained systems:** In a fully  
 2030 personalized FL framework, without any global constraints, the  
 2031 model evolution of each client  $j$ :  $\theta_j^{t+1} = \text{LocalTrain}_j(\theta_j^t)$  is an inde-  
 2032 pendent Markov chain. Under non-IID, the steady-state distribu-  
 2033 tions of these chains may be different, resulting in  $\lim_{t \rightarrow \infty} \text{Var}_j(\theta_j^t)$   
 2034 being non-zero or even divergent, causing the entire system to fail  
 2035 to converge to a good performance level. This can be formalized as:

$$2036 \lim_{t \rightarrow \infty} \text{Var}_j(\|\theta_j^t - \bar{\theta}^t\|_2^2) \rightarrow \infty \text{ or const} > 0, \quad (59)$$

2037 where  $\bar{\theta}^t$  is form of global average model.

2038 **Centroid property of  $w_{\text{ref}}^t$ :** According to the definition of  $w_{\text{ref}}^t$ ,  
 2039  $w_{\text{ref}}^t = E_{j \in N^t, n'_j} [\theta_j^t]$  (before aggregation). This shows that  $w_{\text{ref}}^t$  is  
 2040 the weighted centroid of the current active client group in the  
 2041 parameter space. It is defined as Eq. (23) where:

$$2042 p_i^t = \left( \sum_{j \in C_i^t} n'_j \right) / \left( \sum_{k \in N^t} n'_k \right), \quad (60)$$

$$2043 \theta_{i,\text{rep}}^t = \frac{\sum_{j \in C_i^t} n'_j \theta_j^t}{\sum_{j \in C_i^t} n'_j}. \quad (61)$$

2044 **Regularization Effect of Soft Update:** The global reference  
 2045 model  $w_{\text{ref}}^t$  we introduced is designed to solve this problem. The  
 2046 soft update mechanism sets  $w_{\text{start},k}^{t+1} = w_{\text{ref}}^t$  for a portion of clients  
 2047  $k$  (especially those that have just joined or have not participated  
 2048 for a long time). This is equivalent to "pulling back" the models of  
 2049 these clients to the knowledge center of round  $t$  before the start of  
 2050 round  $t + 1$ .

2051 We can regard the evolution of the entire system as a controlled  
 2052 random process, as a variant of a conditional reset or proximal  
 2053 update. At each time step  $t$ , the soft update operation effectively  
 2054 reduces the model variance. Let  $V^t = \text{Var}_j(\theta_j^t)$ . At the beginning  
 2055 of  $t + 1$  round, due to soft update,  $\mathbb{E}[V_{\text{start}}^{t+1}] < V^t$ . Although local  
 2056 training will increase the variance again, the periodic pullback  
 2057 operation ensures that the variance does not grow unbounded.

2058 More formally, we can model this process as  $V^{t+1} \leq (1 - \alpha)V^t + C$ ,  
 2059 where  $\alpha > 0$  is the variance reduction rate brought by soft update  
 2060 and  $C$  is the variance increase introduced by local training. As  
 2061 long as  $C$  is bounded (which holds under the bounded gradient  
 2062 assumption), this iterative process ensures that  $V^t$  will eventually  
 2063 converge to a bounded steady-state value, thus preventing the  
 2064 system from diverging.

□

2065 Based on the above, the multi-granularity aggregation strategy  
 2066 proposed by our method is a rigorous inheritance and optimization  
 2067 of Donahue K.'s theory, especially the similarity measure in the Fine-  
 2068 grained strategy is more robust mathematically. The soft update  
 2069 mechanism fundamentally solves the convergence shortcomings of  
 2070 static theory in dynamic applications, providing a solid theoretical  
 2071 guarantee for the stable operation of our entire framework.

2072