# AFL: A Single-Round Analytic Approach for Federated Learning with Pre-trained Models

## Supplementary Material

#### A. Proof of Lemma 1

*Proof.* We prove this Lemma mainly based on the existing MP inverse partition result [4] as follows.

In [4], it has been demonstrated that the MP inverse of any matrix,  $\mathbf{A} = \begin{bmatrix} U & V \end{bmatrix}$  can be written as

$$oldsymbol{\mathcal{A}}^{\dagger} = egin{bmatrix} oldsymbol{U} & oldsymbol{V} \end{bmatrix}^{\dagger} =$$

$$\begin{bmatrix} U^{\dagger} - U^{\dagger}VC^{\dagger} - U^{\dagger}V(I - C^{\dagger}C)KV^{*}U^{\dagger*}U^{\dagger}(I - VC^{\dagger}) \\ V^{\dagger} - V^{\dagger}U\tilde{C}^{\dagger} - V^{\dagger}U(I - \tilde{C}^{\dagger}\tilde{C})\tilde{K}U^{*}V^{\dagger*}V^{\dagger}(I - U\tilde{C}^{\dagger}) \end{bmatrix},$$
(A.1)

where

$$egin{cases} oldsymbol{C} = (oldsymbol{I} - oldsymbol{U}oldsymbol{U}^\dagger)oldsymbol{V} \ oldsymbol{ ilde{C}} = (oldsymbol{I} - oldsymbol{V}oldsymbol{V}^\dagger)oldsymbol{U} \ ,$$

$$\begin{cases} \boldsymbol{K} = \left[ \boldsymbol{I} + (\boldsymbol{I} - \boldsymbol{C}^{\dagger}\boldsymbol{C})\boldsymbol{V}^{*}\boldsymbol{U}^{\dagger}\boldsymbol{V}^{\dagger}\boldsymbol{V}(\boldsymbol{I} - \boldsymbol{C}^{\dagger}\boldsymbol{C}) \right]^{-1} \\ \tilde{\boldsymbol{K}} = \left[ \boldsymbol{I} + (\boldsymbol{I} - \boldsymbol{C}^{\dagger}\boldsymbol{C})\boldsymbol{U}^{*}\boldsymbol{V}^{\dagger}\boldsymbol{V}^{\dagger}\boldsymbol{U}(\boldsymbol{I} - \boldsymbol{C}^{\dagger}\boldsymbol{C}) \right]^{-1} \\ \end{cases}$$
(A.2)

In the case of  $\boldsymbol{X} = \begin{bmatrix} \boldsymbol{X}_u \\ \boldsymbol{X}_v \end{bmatrix}$  with only real numbers, we substitute  $\boldsymbol{U}$  with  $\boldsymbol{X}_u^{\top}$ ,  $\boldsymbol{V}$  with  $\boldsymbol{X}_v^{\top}$ . This rewrites (A.1), (A.2) into

$$\boldsymbol{X}^{\dagger} = \begin{bmatrix} \boldsymbol{X}_{u} \\ \boldsymbol{X}_{v} \end{bmatrix}^{\dagger} = \begin{bmatrix} \boldsymbol{X}_{u}^{\dagger\top} - \boldsymbol{X}_{u}^{\dagger\top} \boldsymbol{X}_{v}^{\top} \boldsymbol{C}^{\dagger} \\ \boldsymbol{X}_{v}^{\dagger\top} - \boldsymbol{X}_{v}^{\dagger\top} \boldsymbol{X}_{u}^{\top} \tilde{\boldsymbol{C}}^{\dagger} \end{bmatrix}^{\top} -$$

$$\boldsymbol{X}_{u}^{\dagger\top} \boldsymbol{X}_{v}^{\top} (\boldsymbol{I} - \boldsymbol{C}^{\dagger} \boldsymbol{C}) \boldsymbol{K} \boldsymbol{X}_{v} \boldsymbol{X}_{u}^{\dagger} \boldsymbol{X}_{u}^{\dagger\top} (\boldsymbol{I} - \boldsymbol{X}_{v}^{\top} \boldsymbol{C}^{\dagger}) \end{bmatrix}^{\top},$$

$$\boldsymbol{X}_{v}^{\dagger\top} \boldsymbol{X}_{u}^{\top} (\boldsymbol{I} - \tilde{\boldsymbol{C}}^{\dagger} \tilde{\boldsymbol{C}}) \tilde{\boldsymbol{K}} \boldsymbol{X}_{u} \boldsymbol{X}_{v}^{\dagger} \boldsymbol{X}_{v}^{\dagger\top} (\boldsymbol{I} - \boldsymbol{X}_{u}^{\top} \tilde{\boldsymbol{C}}^{\dagger}) \end{bmatrix}^{\top},$$
(A.3)

where

$$\begin{cases} \boldsymbol{C} = (\boldsymbol{I} - \boldsymbol{X}_{u}^{\top} \boldsymbol{X}_{u}^{\dagger \top}) \boldsymbol{X}_{v}^{\top} \\ \tilde{\boldsymbol{C}} = (\boldsymbol{I} - \boldsymbol{X}_{v}^{\top} \boldsymbol{X}_{v}^{\dagger \top}) \boldsymbol{X}_{u}^{\top} \end{cases}, \\ \begin{cases} \boldsymbol{K} = \left[ \boldsymbol{I} + (\boldsymbol{I} - \boldsymbol{C}^{\dagger} \boldsymbol{C}) \boldsymbol{X}_{v} \boldsymbol{X}_{u}^{\dagger} \boldsymbol{X}_{v}^{\dagger \top} \boldsymbol{X}_{v}^{\top} (\boldsymbol{I} - \boldsymbol{C}^{\dagger} \boldsymbol{C}) \right]^{-1} \\ \tilde{\boldsymbol{K}} = \left[ \boldsymbol{I} + (\boldsymbol{I} - \boldsymbol{C}^{\dagger} \boldsymbol{C}) \boldsymbol{X}_{u} \boldsymbol{X}_{v}^{\dagger} \boldsymbol{X}_{v}^{\dagger \top} \boldsymbol{X}_{u}^{\top} (\boldsymbol{I} - \boldsymbol{C}^{\dagger} \boldsymbol{C}) \right]^{-1} \end{cases}. \end{cases}$$
(A.4)

As  $X_u$  and  $X_v$  are of full column ranks, we obtain an alternative formulation of the MP inverse, i.e.,

$$\boldsymbol{X}_{u}^{\dagger} = (\boldsymbol{X}_{u}^{\top} \boldsymbol{X}_{u})^{-1} \boldsymbol{X}_{u}^{\top}, \quad \boldsymbol{X}_{v}^{\dagger} = (\boldsymbol{X}_{v}^{\top} \boldsymbol{X}_{v})^{-1} \boldsymbol{X}_{v}^{\top}.$$
(A.5)

Hence we have

$$C = (I - X_u^{\top} X_u^{\dagger \top}) X_v^{\top}$$
  
=  $(I - X_u^{\top} X_u (X_u^{\top} X_u)^{-1}) X_v^{\top} = 0.$  (A.6)

Similarly,

$$\tilde{\boldsymbol{C}} = (\boldsymbol{I} - \boldsymbol{X}_{v}^{\top} \boldsymbol{X}_{v}^{\dagger \top}) \boldsymbol{X}_{u}^{\top}$$
$$= (\boldsymbol{I} - \boldsymbol{X}_{v}^{\top} \boldsymbol{X}_{v} (\boldsymbol{X}_{v}^{\top} \boldsymbol{X}_{v})^{-1}) \boldsymbol{X}_{u}^{\top} = \boldsymbol{0}.$$
(A.7)

This simplifies  $oldsymbol{K}$  and  $oldsymbol{ ilde{K}}$  as

$$\begin{cases} \boldsymbol{K} = (\boldsymbol{I} + \boldsymbol{X}_{v} \boldsymbol{X}_{u}^{\dagger} \boldsymbol{X}_{u}^{\dagger \top} \boldsymbol{X}_{v}^{\dagger})^{-1} \\ \tilde{\boldsymbol{K}} = (\boldsymbol{I} + \boldsymbol{X}_{u} \boldsymbol{X}_{v}^{\dagger} \boldsymbol{X}_{v}^{\dagger \top} \boldsymbol{X}_{u}^{\dagger})^{-1} \end{cases}$$
(A.8)

According to the Woodbury Matrix Identity, i.e., for conformable matrices  $A \in \mathbb{R}^{n \times n}$ ,  $B \in \mathbb{R}^{n \times m}$ ,  $E \in \mathbb{R}^{m \times m}$ , and  $D \in \mathbb{R}^{m \times n}$ ,

$$(A + BED)^{-1} = A^{-1} - A^{-1}B(E^{-1} + DA^{-1}B)^{-1}DA^{-1},$$
  
(A.9)

we expand K by substituting  $A = I, B = X_v, E = X_u^{\dagger} X_u^{\dagger \top}$ , and  $D = X_v^{\top}$ , leading to

$$\boldsymbol{K} = \boldsymbol{I} - \boldsymbol{X}_{v} (\boldsymbol{X}_{u}^{\top} \boldsymbol{X}_{u} + \boldsymbol{X}_{v}^{\top} \boldsymbol{X}_{v})^{-1} \boldsymbol{X}_{v}^{\top}.$$
(A.10)

Similarly,

$$\tilde{\boldsymbol{K}} = \boldsymbol{I} - \boldsymbol{X}_u (\boldsymbol{X}_u^\top \boldsymbol{X}_u + \boldsymbol{X}_v^\top \boldsymbol{X}_v)^{-1} \boldsymbol{X}_u^\top.$$
(A.11)

Thus,

$$\boldsymbol{X}^{\dagger} = \begin{bmatrix} \boldsymbol{X}_{u}^{\dagger\top} - \boldsymbol{X}_{u}^{\dagger\top} \boldsymbol{X}_{v}^{\top} \boldsymbol{K} \boldsymbol{X}_{v} \boldsymbol{X}_{u}^{\dagger} \boldsymbol{X}_{u}^{\dagger\top} \\ \boldsymbol{X}_{v}^{\dagger\top} - \boldsymbol{X}_{v}^{\dagger\top} \boldsymbol{X}_{u}^{\top} \boldsymbol{\tilde{K}} \boldsymbol{X}_{u} \boldsymbol{X}_{v}^{\dagger} \boldsymbol{X}_{v}^{\dagger\top} \end{bmatrix}^{\top}.$$
 (A.12)

Let  $oldsymbol{X}^{\dagger} = egin{bmatrix} oldsymbol{ar{U}} & oldsymbol{ar{V}} \end{bmatrix}$  , we have

$$\bar{\boldsymbol{U}} = \left(\boldsymbol{X}_{u}^{\dagger\top} - \boldsymbol{X}_{u}^{\dagger\top}\boldsymbol{X}_{v}^{\top}\boldsymbol{K}\boldsymbol{X}_{v}\boldsymbol{X}_{u}^{\dagger}\boldsymbol{X}_{u}^{\dagger\top}\right)^{\top} \\
= \boldsymbol{X}_{u}^{\dagger} - \boldsymbol{X}_{u}^{\dagger}\boldsymbol{X}_{u}^{\dagger\top}\boldsymbol{X}_{v}^{\top}\boldsymbol{K}^{\top}\boldsymbol{X}_{v}\boldsymbol{X}_{u}^{\dagger} \\
\bar{\boldsymbol{V}} = \left(\boldsymbol{X}_{v}^{\dagger\top} - \boldsymbol{X}_{v}^{\dagger\top}\boldsymbol{X}_{u}^{\top}\tilde{\boldsymbol{K}}\boldsymbol{X}_{u}\boldsymbol{X}_{v}^{\dagger}\boldsymbol{X}_{v}^{\dagger\top}\right)^{\top} \\
= \boldsymbol{X}_{v}^{\dagger} - \boldsymbol{X}_{v}^{\dagger}\boldsymbol{X}_{v}^{\dagger\top}\boldsymbol{X}_{u}^{\top}\tilde{\boldsymbol{K}}^{\top}\boldsymbol{X}_{u}\boldsymbol{X}_{v}^{\dagger} \qquad (A.13)$$

Substitute K and  $\tilde{K}$  with (A.10) and (A.11), we may rewrite (A.13) into

$$\begin{split} \bar{\boldsymbol{U}} &= \boldsymbol{X}_{u}^{\dagger} - \boldsymbol{X}_{u}^{\dagger} \boldsymbol{X}_{u}^{\dagger \top} \boldsymbol{X}_{v}^{\top} \left( \boldsymbol{I} - \boldsymbol{X}_{v} (\boldsymbol{X}_{u}^{\top} \boldsymbol{X}_{u} + \boldsymbol{X}_{v}^{\top} \boldsymbol{X}_{v})^{-1} \boldsymbol{X}_{v}^{\top} \right)^{\top} \boldsymbol{X}_{v} \boldsymbol{X}_{u}, \\ \bar{\boldsymbol{V}} &= \boldsymbol{X}_{v}^{\dagger} - \boldsymbol{X}_{v}^{\dagger} \boldsymbol{X}_{v}^{\dagger \top} \boldsymbol{X}_{u}^{\top} \left( \boldsymbol{I} - \boldsymbol{X}_{u} (\boldsymbol{X}_{u}^{\top} \boldsymbol{X}_{u} + \boldsymbol{X}_{v}^{\top} \boldsymbol{X}_{v})^{-1} \boldsymbol{X}_{u}^{\top} \right)^{\top} \boldsymbol{X}_{u} \boldsymbol{X}_{v}^{\dagger}. \end{split}$$

$$(A.14)$$

That is,

Let

$$\begin{cases} \boldsymbol{C}_{u} = \boldsymbol{X}_{u}^{\top} \boldsymbol{X}_{u} \\ \boldsymbol{C}_{v} = \boldsymbol{X}_{v}^{\top} \boldsymbol{X}_{v} \end{cases}, \text{ and } \begin{cases} \boldsymbol{R}_{u} = \boldsymbol{R}_{u}^{-1} \\ \boldsymbol{R}_{v} = \boldsymbol{R}_{v}^{-1} \end{cases}.$$
(A.16)

we have

$$\begin{cases} \bar{\boldsymbol{U}} = \left[ \boldsymbol{I} - \boldsymbol{R}_{u}\boldsymbol{C}_{v} - \boldsymbol{R}_{u}\boldsymbol{C}_{v}(\boldsymbol{C}_{u} + \boldsymbol{C}_{v})^{-1}\boldsymbol{C}_{v} \right] \boldsymbol{X}_{u}^{\dagger} \\ \bar{\boldsymbol{V}} = \left[ \boldsymbol{I} - \boldsymbol{R}_{v}\boldsymbol{C}_{u} - \boldsymbol{R}_{v}\boldsymbol{C}_{u}(\boldsymbol{C}_{u} + \boldsymbol{C}_{v})^{-1}\boldsymbol{C}_{u} \right] \boldsymbol{X}_{v}^{\dagger} \end{cases}$$
(A.17)

Thus,

$$\boldsymbol{X}^{\dagger} = \begin{bmatrix} \boldsymbol{X}_{u} \\ \boldsymbol{X}_{v} \end{bmatrix}^{\dagger} = \begin{bmatrix} \bar{\boldsymbol{U}} & \bar{\boldsymbol{V}} \end{bmatrix},$$
 (A.18)

which completes the proof.

# **B.** Proof of Theorem 1

Proof. As indicated in Lemma 1, we have

$$\boldsymbol{X}^{\dagger} = \begin{bmatrix} \boldsymbol{\bar{U}} & \boldsymbol{\bar{V}} \end{bmatrix} \tag{A.19}$$

where

$$\begin{cases} \bar{\boldsymbol{U}} = \left[\boldsymbol{I} - \boldsymbol{R}_{u}\boldsymbol{C}_{v} - \boldsymbol{R}_{u}\boldsymbol{C}_{v}(\boldsymbol{C}_{u} + \boldsymbol{C}_{v})^{-1}\boldsymbol{C}_{v}\right]\boldsymbol{X}_{u}^{\dagger} \\ \bar{\boldsymbol{V}} = \left[\boldsymbol{I} - \boldsymbol{R}_{v}\boldsymbol{C}_{u} - \boldsymbol{R}_{v}\boldsymbol{C}_{u}(\boldsymbol{C}_{u} + \boldsymbol{C}_{v})^{-1}\boldsymbol{C}_{u}\right]\boldsymbol{X}_{v}^{\dagger} , \\ \end{cases}$$
(A.20)

and

$$\begin{cases} \boldsymbol{R}_{u} = (\boldsymbol{X}_{u}^{\top}\boldsymbol{X}_{u})^{-1} = \boldsymbol{X}_{u}^{\dagger}\boldsymbol{X}_{u}^{\dagger\top} \\ \boldsymbol{R}_{v} = (\boldsymbol{X}_{v}^{\top}\boldsymbol{X}_{v})^{-1} = \boldsymbol{X}_{v}^{\dagger}\boldsymbol{X}_{v}^{\dagger\top} \\ \begin{cases} \boldsymbol{C}_{u} = \boldsymbol{R}_{u}^{-1} = \boldsymbol{X}_{u}^{\top}\boldsymbol{X}_{u} \\ \boldsymbol{C}_{v} = \boldsymbol{R}_{v}^{-1} = \boldsymbol{X}_{v}^{\top}\boldsymbol{X}_{v} \end{cases}$$
(A.21)

Hence,

$$W = X^{\dagger} Y = \begin{bmatrix} \bar{U} & \bar{V} \end{bmatrix} \begin{bmatrix} Y_u \\ Y_v \end{bmatrix}$$
$$= \bar{U} Y_u + \bar{V} Y_v. \qquad (A.22)$$

By substituting  $\bar{U}$  and  $\bar{V}$  with those in (A.20), we rewrite (A.22) into

$$\begin{split} \hat{\boldsymbol{W}} &= \left[\boldsymbol{I} - \boldsymbol{R}_{u}\boldsymbol{C}_{v} - \boldsymbol{R}_{u}\boldsymbol{C}_{v}(\boldsymbol{C}_{u} + \boldsymbol{C}_{v})^{-1}\boldsymbol{C}_{v}\right]\boldsymbol{X}_{u}^{\dagger}\boldsymbol{Y}_{u} \\ &+ \left[\boldsymbol{I} - \boldsymbol{R}_{v}\boldsymbol{C}_{u} - \boldsymbol{R}_{v}\boldsymbol{C}_{u}(\boldsymbol{C}_{u} + \boldsymbol{C}_{v})^{-1}\boldsymbol{C}_{u}\right]\boldsymbol{X}_{v}^{\dagger}\boldsymbol{Y}_{v}. \end{split}$$
(A.23)

As  $\hat{W}_u = X_u^{\dagger} Y_u$  and  $\hat{W}_v = X_v^{\dagger} Y_v$ , (A.23) can be rewritten as

$$\hat{\boldsymbol{W}} = \begin{bmatrix} \boldsymbol{I} - \boldsymbol{R}_u \boldsymbol{C}_v - \boldsymbol{R}_u \boldsymbol{C}_v (\boldsymbol{C}_u + \boldsymbol{C}_v)^{-1} \boldsymbol{C}_v \end{bmatrix} \hat{\boldsymbol{W}}_u \\ + \begin{bmatrix} \boldsymbol{I} - \boldsymbol{R}_v \boldsymbol{C}_u - \boldsymbol{R}_v \boldsymbol{C}_u (\boldsymbol{C}_u + \boldsymbol{C}_v)^{-1} \boldsymbol{C}_u \end{bmatrix} \hat{\boldsymbol{W}}_v.$$
(A.24)

That is,

$$\hat{\boldsymbol{W}} = \boldsymbol{\mathcal{W}}_u \hat{\boldsymbol{W}}_u + \boldsymbol{\mathcal{W}}_v \hat{\boldsymbol{W}}_v, \qquad (A.25)$$

where

$$\begin{cases} \mathcal{W}_{u} = I - R_{u}C_{v} - R_{u}C_{v}(C_{u} + C_{v})^{-1}C_{v} \\ \mathcal{W}_{v} = I - R_{v}C_{u} - R_{v}C_{u}(C_{u} + C_{v})^{-1}C_{u} \end{cases}, \\ \begin{cases} C_{u} = X_{u}^{\top}X_{u} \\ C_{v} = X_{v}^{\top}X_{v} \end{cases} \text{ and } \begin{cases} R_{u} = C_{u}^{-1} \\ R_{v} = C_{v}^{-1} \end{cases}. \end{cases}$$
(A.26)

## C. Proof of Theorem 2

*Proof.* First we consider the aggregation of two clients. Directly substituting  $\hat{W}_u$  as  $\hat{W}_u^r$  and changing  $C_u$ ,  $C_v$  to  $C_u^r = (X_u^\top X_u + \gamma I)$ ,  $C_v^r = (X_v^\top X_v + \gamma I)$  in Theorem B, we have

$$\hat{\boldsymbol{W}}^{\mathrm{r}} = \boldsymbol{\mathcal{W}}_{u}^{\mathrm{r}} \hat{\boldsymbol{W}}_{u}^{\mathrm{r}} + \boldsymbol{\mathcal{W}}_{v}^{\mathrm{r}} \hat{\boldsymbol{W}}_{v}^{\mathrm{r}}, \qquad (A.27)$$

where

$$\begin{cases} \boldsymbol{\mathcal{W}}_{u}^{\mathrm{r}} = \boldsymbol{I} - \boldsymbol{R}_{u}^{\mathrm{r}} \boldsymbol{C}_{v}^{\mathrm{r}} - \boldsymbol{R}_{u}^{\mathrm{r}} \boldsymbol{C}_{v}^{\mathrm{r}} (\boldsymbol{C}_{u}^{\mathrm{r}} + \boldsymbol{C}_{v}^{\mathrm{r}})^{-1} \boldsymbol{C}_{v}^{\mathrm{r}} \\ \boldsymbol{\mathcal{W}}_{v}^{\mathrm{r}} = \boldsymbol{I} - \boldsymbol{R}_{v}^{\mathrm{r}} \boldsymbol{C}_{u}^{\mathrm{r}} - \boldsymbol{R}_{v}^{\mathrm{r}} \boldsymbol{C}_{u}^{\mathrm{r}} (\boldsymbol{C}_{u}^{\mathrm{r}} + \boldsymbol{C}_{v}^{\mathrm{r}})^{-1} \boldsymbol{C}_{u}^{\mathrm{r}} \\ \end{cases} \\ \begin{cases} \boldsymbol{C}_{u}^{\mathrm{r}} = (\boldsymbol{X}_{u}^{\top} \boldsymbol{X}_{u} + \gamma \boldsymbol{I}) \\ \boldsymbol{C}_{v}^{\mathrm{r}} = (\boldsymbol{X}_{v}^{\top} \boldsymbol{X}_{v} + \gamma \boldsymbol{I}) \end{cases} \text{ and } \begin{cases} \boldsymbol{R}_{u}^{\mathrm{r}} = \boldsymbol{C}_{u}^{\mathrm{r}-1} \\ \boldsymbol{R}_{v}^{\mathrm{r}} = \boldsymbol{C}_{v}^{\mathrm{r}-1} \end{cases} . \end{cases} \text{ (A.28)} \\ \end{cases} \\ \text{Since } \hat{\boldsymbol{W}}_{u}^{\mathrm{r}} = (\boldsymbol{X}_{u}^{\top} \boldsymbol{X}_{u} + \gamma \boldsymbol{I})^{-1} \boldsymbol{X}_{u}^{\top} \boldsymbol{Y}_{u}, \text{ then} \end{cases}$$

$$\mathcal{W}_{u}^{\mathrm{r}}\hat{W}_{u}^{\mathrm{r}} = [I - R_{u}^{\mathrm{r}}C_{v}^{\mathrm{r}} - R_{u}^{\mathrm{r}}C_{v}^{\mathrm{r}}(C_{u}^{\mathrm{r}} + C_{v}^{\mathrm{r}})^{-1}C_{v}^{\mathrm{r}}]R_{u}^{\mathrm{r}}X_{u}^{\top}Y_{u}$$
$$= [R_{u}^{\mathrm{r}} - R_{u}^{\mathrm{r}}C_{v}^{\mathrm{r}}R_{u}^{\mathrm{r}} - R_{u}^{\mathrm{r}}C_{v}^{\mathrm{r}}(C_{u}^{\mathrm{r}} + C_{v}^{\mathrm{r}})^{-1}C_{v}^{\mathrm{r}}R_{u}^{\mathrm{r}}]X_{u}^{\top}Y_{u}.$$
(A.29)

According to the Woodbury Matrix Identity in (A.9), let B = I, D = I, we have

$$(\mathbf{A} + \mathbf{E})^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1} (\mathbf{A}^{-1} + \mathbf{E}^{-1})^{-1} \mathbf{A}^{-1}.$$
 (A.30)

Then we have

$$A^{-1}(A^{-1} + E^{-1})^{-1}A^{-1} = A^{-1} - (A + E)^{-1}$$
. (A.31)

Swapping  $A^{-1}$  with  $C_v^{r}$  and  $E^{-1}$  with  $C_u^{r}$ , we have

$$C_v^{\rm r} (C_v^{\rm r} + C_u^{\rm r})^{-1} C_v^{\rm r} = C_v^{\rm r} - (R_u^{\rm r} + R_v^{\rm r})^{-1}.$$
 (A.32)

Similarly,

$$\boldsymbol{R}_{u}^{r}(\boldsymbol{R}_{u}^{r}+\boldsymbol{R}_{v}^{r})^{-1}\boldsymbol{R}_{r}^{r}=\boldsymbol{R}_{u}^{r}-(\boldsymbol{C}_{u}^{r}+\boldsymbol{C}_{v}^{r})^{-1}.$$
 (A.33)

By substituting (A.32) into (A.29),

$$\mathcal{W}_{u}^{r}\hat{W}_{u}^{r} = [R_{u}^{r} - R_{u}^{r}C_{v}^{r}R_{u}^{r} + R_{u}^{r}C_{v}^{r}R_{u}^{r} - R_{u}^{r}(R_{u}^{r} + R_{v}^{r})^{-1}R_{u}^{r}]X_{u}^{\top}Y_{u}.$$
(A.34)

Then

$$\boldsymbol{\mathcal{W}}_{u}^{\mathrm{r}}\boldsymbol{\hat{W}}_{u}^{\mathrm{r}} = [\boldsymbol{R}_{u}^{\mathrm{r}} - \boldsymbol{R}_{u}^{\mathrm{r}}(\boldsymbol{R}_{u}^{\mathrm{r}} + \boldsymbol{R}_{v}^{\mathrm{r}})^{-1}\boldsymbol{R}_{u}^{\mathrm{r}}]\boldsymbol{X}_{u}^{\top}\boldsymbol{Y}_{u}. \quad (A.35)$$

Further substituting (A.33) into (A.35), we have

$$\mathcal{W}_{u}^{\mathbf{r}} \hat{\mathcal{W}}_{u}^{\mathbf{r}} = [\mathbf{R}_{u}^{\mathbf{r}} - \mathbf{R}_{u}^{\mathbf{r}} (\mathbf{R}_{u}^{\mathbf{r}} + \mathbf{R}_{v}^{\mathbf{r}})^{-1} \mathbf{R}_{u}^{\mathbf{r}}] \mathbf{X}_{u}^{\top} \mathbf{Y}_{u}$$

$$= [\mathbf{R}_{u}^{\mathbf{r}} - \mathbf{R}_{u}^{\mathbf{r}} + (\mathbf{C}_{u}^{\mathbf{r}} + \mathbf{C}_{v}^{\mathbf{r}})^{-1}] \mathbf{X}_{u}^{\top} \mathbf{Y}_{u}$$

$$= (\mathbf{C}_{u}^{\mathbf{r}} + \mathbf{C}_{v}^{\mathbf{r}})^{-1} \mathbf{X}_{u}^{\top} \mathbf{Y}_{u}$$

$$= (\mathbf{C}_{u} + \mathbf{C}_{v} + 2\gamma \mathbf{I})^{-1} \mathbf{X}_{u}^{\top} \mathbf{Y}_{u}. \quad (A.36)$$

Similarly,

$$\boldsymbol{\mathcal{W}}_{v}^{\mathrm{r}}\boldsymbol{\hat{W}}_{v}^{\mathrm{r}} = (\boldsymbol{C}_{u} + \boldsymbol{C}_{v} + 2\gamma\boldsymbol{I})^{-1}\boldsymbol{X}_{v}^{\top}\boldsymbol{Y}_{v}. \qquad (A.37)$$

Thus equation (A.27) can be converted to

$$\hat{\boldsymbol{W}}^{\mathrm{r}} = (\boldsymbol{C}_u + \boldsymbol{C}_v + 2\gamma \boldsymbol{I})^{-1} (\boldsymbol{X}_u^{\top} \boldsymbol{Y}_u + \boldsymbol{X}_v^{\top} \boldsymbol{Y}_v). \quad (A.38)$$

Since  $\hat{W} = X^{\dagger}Y$  and  $X = \begin{bmatrix} X_u \\ X_v \end{bmatrix}$ ,  $Y = \begin{bmatrix} Y_u \\ Y_v \end{bmatrix}$ , with full-column rank of X,

$$\hat{\boldsymbol{W}} = \boldsymbol{X}^{\dagger} \boldsymbol{Y} = (\boldsymbol{X}^{\mathsf{T}} \boldsymbol{X})^{-1} \boldsymbol{X}^{\mathsf{T}} \boldsymbol{Y}$$
(A.39)  
$$= (\begin{bmatrix} \boldsymbol{X}_{u}^{\mathsf{T}} & \boldsymbol{X}_{v}^{\mathsf{T}} \end{bmatrix} \begin{bmatrix} \boldsymbol{X}_{u} \\ \boldsymbol{X}_{v} \end{bmatrix})^{-1} \begin{bmatrix} \boldsymbol{X}_{u}^{\mathsf{T}} & \boldsymbol{X}_{v}^{\mathsf{T}} \end{bmatrix} \begin{bmatrix} \boldsymbol{Y}_{u} \\ \boldsymbol{Y}_{v} \end{bmatrix})$$
$$= (\boldsymbol{X}_{u}^{\mathsf{T}} \boldsymbol{X}_{u} + \boldsymbol{X}_{v}^{\mathsf{T}} \boldsymbol{X}_{v})^{-1} (\boldsymbol{X}_{u}^{\mathsf{T}} \boldsymbol{Y}_{u} + \boldsymbol{X}_{v}^{\mathsf{T}} \boldsymbol{Y}_{v})$$
$$= (\boldsymbol{C}_{u} + \boldsymbol{C}_{v})^{-1} (\boldsymbol{X}_{u}^{\mathsf{T}} \boldsymbol{Y}_{u} + \boldsymbol{X}_{v}^{\mathsf{T}} \boldsymbol{Y}_{v}).$$

By comparing with (A.27) and (A.39), we can obtain the relation between  $\hat{W}$  and  $\hat{W}^{r}$  as follows.

$$\hat{\boldsymbol{W}}^{\mathrm{r}} = (\boldsymbol{C}_{u}^{\mathrm{r}} + \boldsymbol{C}_{v}^{\mathrm{r}})^{-1} (\boldsymbol{C}_{u} + \boldsymbol{C}_{v}) \hat{\boldsymbol{W}}.$$
 (A.40)

By extending to the multi-client scenario, we have

$$\hat{\boldsymbol{W}}_{\text{agg},k}^{\text{r}} = (\boldsymbol{C}_{\text{agg},k}^{\text{r}})^{-1} \boldsymbol{C}_{\text{agg},k} \hat{\boldsymbol{W}}_{\text{agg},k}, \qquad (A.41)$$

where

$$C_{\text{agg},k}^{\text{r}} = C_{\text{agg},k} + k\gamma I = \sum_{i}^{k} C_{i}^{\text{r}},$$
$$C_{i}^{\text{r}} = X_{i}^{\top} X_{i} + \gamma I, \qquad (A.42)$$

which complete the proof.

### D. Validating AA Laws on Dummy Dataset

Here we validate the AA laws in AFL, whether the aggregated weight  $\hat{W}_{agg,K}$  equals to  $\hat{W}$  trained on a centralized dataset. This is done by measuring the deviation (i.e.,  $\Delta W = \|\hat{W} - \hat{W}_{agg,K}\|_1$ ) between the joint-trained weight and the aggregated one on a dummy dataset.

**Dummy Dataset.** We randomly generate a 512dimension and 10,000-sample dummy dataset. This dataset has 10 classes, with each class containing an identical number of samples. The samples in the dummy dataset are randomly but evenly distributed to K clients (we set K = 2, 10, 20, 50, 100, 200).

**Results.** As indicated in Table A.1, without the RI process, the deviation is negligible for K = 2, 10, but it grows with an increasing K and could become unacceptable (e.g.,  $3.67 \times 10^{12}$  for K = 200). This is because the full-column rank assumption might not hold anymore for large K. By adopting the RI process, the deviations become negligible (around  $10^{-10}$ ) for various K values as shown in the second row of Table A.1. The RI process introduces  $\gamma$  (we adopt  $\gamma = 1$  in this case, but any value would suffice) to satisfy full-column rank condition, which is later removed in (16) to restore the AA law's optimality. This experiment has well demonstrated AFL's invariance to data partitioning with empirical evidence. The codes for the dummy data validation can be found in the file  $App_Dummy.ipynb$  in the released code.

#### **E.** Implementation Details of Experiments

For the compared methods, we set the local epoch to 1 and all the clients are selected to participate each round after local training. The batch size is set to 64 and we employ SGD optimizer with learning rate of 0.05. The number of global communication rounds is set to be 500 since there is little or no performance gain with more rounds. We report the average and standard deviation of best top-1 accuracy in three runs. All the experiments are conducted on a NVIDIA GeForce RTX 4090 GPU. All the compared methods except FedDisco are implemented with PFLlib [45] and FedDisco is implemented upon FedAvg via official codes.

For the specified hyperparameters in compared methods, we tune the parameters via grid search. For FedProx [18], we tune the hyperparameter  $\mu$  from {0.0001, 0.001, 0.01, 0.1}. For MOON [16], we tune the hyperparameter  $\mu$  from {0.1, 1, 5, 10}. For FedDyn [1], we tune the hyperparameter  $\alpha$  from {0.001, 0.01, 0.1, 1.0}. For FedNTD [15], we tune the hyperparameters  $\tau$  from {0.1, 0.5, 1.0, 2.0}. For FedDisco [43], we tune the hyperparameters *a* from {0.01, 0.05, 0.1, 0.5} and *b* from {0.005, 0.01, 0.05, 0.1}. The best parameters we adopted are,  $\mu = 0.001$  in FedProx [18],  $\mu = 1$  in MOON [16],  $\alpha = 1.0$  in FedDyn [1],  $\tau = 0.5, \beta = 1.0$  in FedNTD [15] and a = 0.05, b = 0.01 in FedDisco [43].

#### F. Necessity of FL with Pre-trained Model

To validate the necessity of FL with pre-trained backbone, we train the models locally without aggregation under the setting of  $\alpha$ =0.1, K=100. We report the average and maximum test accuracy of local training among all the clients. As shown in Table. A.2, without aggregation, the results of local training (12.04% and 16.36%) fall behind FedAvg (56.57%)

Table A.1. Deviation  $\Delta W$  between the joint-trained weight and the aggregated one (average of 5 runs).

Difference	K = 2	K = 10	K = 20	K = 50	K = 100	K = 200
$\Delta W$ (w/o RI)	$7.83 \times 10^{-14}$	$1.76 \times 10^{-12}$	$9.86 \times 10^{-1}$	5.90	$5.93 \times 10^4$	$3.67 \times 10^{12}$
$\Delta W$ (w/ RI)	$4.94 \times 10^{-14}$	$1.74 \times 10^{-12}$	$5.09 \times 10^{-10}$	$8.45 \times 10^{-10}$	$7.57 \times 10^{-10}$	$7.81 \times 10^{-10}$

and AFL (58.56%) with large gaps. Training local models without FL could suffer from the data heterogeneity and the collaboration among clients is beneficial. This pattern is also validated in previous studies with pre-trained backbone [7].

Table A.2. Comparison between FL teniques including FedAvg and AFL with local training when utilizing pre-trained backbone.

Methods	Local Max	Local Avg	FedAvg	AFL
Acc.(%)	16.36	12.04	56.57	58.56

## G. Comparative Study with Single-Round FL

In this section, we provide a comparative study with another single-round FL technique FedFisher [11] to further validate our proposed AFL. We compare AFL with FedFisher under the setting of  $\alpha = 0.1, K = 50$  (larger K will lead to out-of-memory in FedFisher) with the same pre-trained ResNet-18 provided by the repository of FedFisher. As shown in Table A.3, AFL outperforms FedFisher with considerable distance (35.87% v.s. 19.31%). FedFisher utilizes iterative gradient-descent to aggregate local weights and MSE loss is established by the difference of global and local weights to preventing drifting between them during aggregation. However, this technique could still suffer from the data heterogeneity, while AFL formulates the AA law to achieve the invariance to data partitioning, enabling outperforming result when compared with FedFisher.

Table A.3. Comparative study between AFL and FedFisher.

Methods	FedFisher	AFL
Acc.(%)	19.31	35.87