

483 **A Appendix**

484 **A.1 Theoretical Proofs**

485 **Notations of Convolutional Operations.** In our paper, we express convolution operation as
 486 $\mathbf{Z}_u = \alpha \mathbf{X}_p \mathbf{D}_u$. More explicitly, the formulation writes as $\mathbf{Z}_u = \alpha \mathbf{X}_p \star \mathbf{D}_u$, where \star is the
 487 convolutional operation. By converting convolutional kernel \mathbf{D}_u into a Toeplitz matrix, we can
 488 replace the convolution operation $\mathbf{X}_p \star \mathbf{D}_u$ with matrix multiplication $\mathbf{X}_p \mathbf{D}_u$. We also modify α by
 489 $I_{h,w} \otimes \alpha$, where \otimes is Kronecker product, to enable the matrix multiplication $\alpha \mathbf{X}_p \mathbf{D}_u$.

490 **Proposition A.1.** *Suppose \mathbf{D}_u and \mathbf{D}_v are two different sets of filter atoms for a convolutional layer*
 491 *with the common atom coefficients α , we can upper bound the changes in the corresponding features*
 492 *$\mathbf{Z}_u, \mathbf{Z}_v$ with atom changes,*

$$\|\mathbf{Z}_u - \mathbf{Z}_v\|_F \leq (\|\alpha\|_F \lambda) \sqrt{|\mathcal{B}|} \cdot \|(\mathbf{D}_u - \mathbf{D}_v)\|_F, \quad \text{with } \lambda = \sup_{b \in \mathcal{B}} \|\mathbf{X}\|_{F, N_b}, \quad (8)$$

493 *Proof.* Recall the decomposed convolution can be expressed as,

$$\mathbf{Z} = \sum_{i=1}^m \alpha_i \langle \mathbf{X}, \mathbf{D}[i] \rangle_{N_b} \quad (9)$$

494 $\forall b$ we have,

$$\begin{aligned} |\mathbf{Z}_u(b) - \mathbf{Z}_v(b)| &= \left| \sum_{i=1}^m \alpha_i \langle \mathbf{X}, \mathbf{D}_u[i] \rangle_{N_b} - \sum_{i=1}^m \alpha_i \langle \mathbf{X}, \mathbf{D}_v[i] \rangle_{N_b} \right| \\ &\leq \|\alpha\|_F \left(\sum_{i=1}^m |\langle \mathbf{X}, (\mathbf{D}_u[i] - \mathbf{D}_v[i]) \rangle_{N_b}|^2 \right)^{1/2}. \end{aligned} \quad (10)$$

495 By Cauchy-Schwarz inequality,

$$\begin{aligned} |\langle \mathbf{X}, (\mathbf{D}_u[i] - \mathbf{D}_v[i]) \rangle_{N_b}| &\leq \|\mathbf{X}\|_{F, N_b} \cdot \|\mathbf{D}_u[i] - \mathbf{D}_v[i]\|_{F, N_b} \\ &\leq \lambda \cdot \|\mathbf{D}_u[i] - \mathbf{D}_v[i]\|_{F, N_b} \end{aligned} \quad (11)$$

496 we have that

$$\begin{aligned} \sum_{b \in \mathcal{B}} |\mathbf{Z}_u(b) - \mathbf{Z}_v(b)|^2 &\leq \|\alpha\|_F^2 \sum_b \sum_{i=1}^m |\langle \mathbf{X}, (\mathbf{D}_u[i] - \mathbf{D}_v[i]) \rangle_{N_b}|^2 \\ &\leq \|\alpha\|_F^2 \sum_b \sum_{i=1}^m \|\mathbf{X}\|_{F, N_b}^2 \cdot \|\mathbf{D}_u[i] - \mathbf{D}_v[i]\|_{F, N_b}^2 \\ &\leq (\|\alpha\|_F \lambda)^2 \sum_{b, i} \|\mathbf{D}_u[i] - \mathbf{D}_v[i]\|_{F, N_b}^2 \end{aligned} \quad (12)$$

497 and observe that

$$\sum_{b, i} \|\mathbf{D}_u[i] - \mathbf{D}_v[i]\|_{F, N_b}^2 = \sum_{b \in \mathcal{B}} \sum_{i=1}^m \|\mathbf{D}_u[i] - \mathbf{D}_v[i]\|_{F, N_b}^2 = |\mathcal{B}| \cdot \|(\mathbf{D}_u - \mathbf{D}_v)\|_F^2, \quad (13)$$

498 where $|\mathcal{B}|$ is the area of the domain of \mathbf{X} . Then Eq. 12 becomes

$$\sum_{b \in \mathcal{B}} |\mathbf{Z}_u(b) - \mathbf{Z}_v(b)|^2 \leq (\|\alpha\|_F \lambda)^2 |\mathcal{B}| \cdot \|(\mathbf{D}_u - \mathbf{D}_v)\|_F^2, \quad (14)$$

499 which proves that $\|\mathbf{Z}_u - \mathbf{Z}_v\|_F \leq (\|\alpha\|_F \lambda) \sqrt{|\mathcal{B}|} \cdot \|(\mathbf{D}_u - \mathbf{D}_v)\|_F$ as claimed.

500 □

501 **Proposition A.2.** *Assume filter atoms $\mathbf{D}_u, \mathbf{D}_v$ are orthogonal matrices, then $\mathcal{S}_{Gras} = \mathcal{S}_{Atom}$.*

502 *Proof.* Since $\mathbf{D}_u, \mathbf{D}_v \in \mathbb{R}^{k^2 \times m}$ are orthogonal matrices, i.e., $\mathbf{D}_u^T \mathbf{D}_u = \mathbf{D}_v^T \mathbf{D}_v = \mathbf{I}$, the Grassmann
503 similarity can be represented as,

$$S_{Gras}(\mathcal{F}_u, \mathcal{F}_v) = \frac{1}{m} \sum_i^m \cos \theta_i = \frac{1}{m} \sum_i^m \sigma_i, \quad (15)$$

504 where $\sigma_i = \Sigma_{ii}, U\Sigma V = \mathbf{D}_u^T \mathbf{D}_v$.

505 \mathcal{S}_{Atom} is defined as,

$$\mathcal{S}_{Atom}(\mathcal{F}_u, \mathcal{F}_v) = \cos(\mathbf{D}_u, \mathbf{D}_v) = \frac{\langle \text{vec}(\mathbf{D}_u), \text{vec}(\mathbf{D}_v) \rangle}{\|\text{vec}(\mathbf{D}_u)\|_F \cdot \|\text{vec}(\mathbf{D}_v)\|_F}. \quad (16)$$

506 Analyze each part separately, we have $\langle \text{vec}(\mathbf{D}_u), \text{vec}(\mathbf{D}_v) \rangle = \text{Tr}(\mathbf{D}_u^T \mathbf{D}_v) = \sum_i^m \sigma_i$,
507 $\|\text{vec}(\mathbf{D}_u)\|_F = \sqrt{\text{Tr}(\mathbf{D}_u^T \mathbf{D}_u)} = \sqrt{\text{Tr}(\mathbf{I})} = \sqrt{m}$, and also $\|\text{vec}(\mathbf{D}_v)\|_F = \sqrt{m}$. In total,
508 the filter subspace similarity becomes,

$$\mathcal{S}_{Atom}(\mathcal{F}_u, \mathcal{F}_v) = \cos(\mathbf{D}_u, \mathbf{D}_v) = \frac{\sum_i^m \sigma_i}{m}, \quad (17)$$

509 which equals \mathcal{S}_{Gras} . The claimed theorem is proved.

510 □

511 **Lemma A.3.** For two positive semidefinite matrices \mathbf{A}, \mathbf{B} ,

$$\text{Tr}(\mathbf{A}\mathbf{B}) \geq \sigma_{\min}(\mathbf{A})\text{Tr}(\mathbf{B}), \quad (18)$$

512 where σ_{\min} denotes the minimum eigenvalue of \mathbf{A} .

513 *Proof.* It is equivalent to prove that,

$$\text{Tr}((\mathbf{A} - \sigma_{\min}(\mathbf{A})\mathbf{I})\mathbf{B}) \geq 0. \quad (19)$$

514 Let \mathbf{C}, \mathbf{D} be matrices such that $\mathbf{A} - \sigma_{\min}(\mathbf{A})\mathbf{I} = \mathbf{C}^T \mathbf{C}$, $\mathbf{B} = \mathbf{D}^T \mathbf{D}$, then

$$\begin{aligned} \text{Tr}((\mathbf{A} - \sigma_{\min}(\mathbf{A})\mathbf{I})\mathbf{B}) &= \text{Tr}(\mathbf{C}^T \mathbf{C} \mathbf{D}^T \mathbf{D}) \\ &= \text{Tr}(\mathbf{C} \mathbf{D}^T \mathbf{D} \mathbf{C}^T) \\ &= \text{Tr}((\mathbf{D} \mathbf{C}^T)^T (\mathbf{D} \mathbf{C}^T)) \geq 0. \end{aligned} \quad (20)$$

515 □

516 **Theorem A.4.** Suppose the forward of decomposed convolution layer for the u -th model is $\mathbf{Z}_u =$
517 $\alpha \mathbf{X} \mathbf{D}_u$. $\mathbf{Z}_u, \mathbf{Z}_v$ nearly have zero-mean since \mathbf{X}_p is preprocessed to be normalized. CCA coefficient
518 is defined as $S(\mathbf{Z}_u, \mathbf{Z}_v) = \sqrt{\frac{1}{c} \sum_{i=1}^c \sigma_i^2}$, where σ_i^2 denotes the i -th eigenvalue of $\Lambda_{u,v} = \mathbf{Q}_u^T \mathbf{Q}_v$,
519 $\mathbf{Q}_u = \mathbf{Z}_u (\mathbf{Z}_u^T \mathbf{Z}_u)^{-\frac{1}{2}}$. Then $S(\mathbf{Z}_u, \mathbf{Z}_v)$ is upper bounded,

$$S(\mathbf{Z}_u, \mathbf{Z}_v) \leq \frac{c^{\frac{3}{2}} \mathcal{T}}{c} \cos(\mathbf{D}_u, \mathbf{D}_v), \quad (21)$$

520 where $\mathcal{T} = \text{Tr}(\mathbf{X}^T \alpha^T \alpha \mathbf{X})$, $\mathcal{C} = \sigma_{\min}(\mathbf{X}^T \alpha^T \alpha \mathbf{X})$.

521 *Proof.* Consider $\mathcal{S}^2 = \frac{1}{c} \sum_{i=1}^c \sigma_i^2$.

$$\mathcal{S}^2 = \frac{1}{c} \sum_{i=1}^c \sigma_i^2 = \frac{1}{c} \text{Tr}(\Lambda_{u,v} \Lambda_{u,v}^T). \quad (22)$$

522 where

$$\text{Tr}(\Lambda_{u,v} \Lambda_{u,v}^T) = \text{Tr}(\mathbf{Q}_u^T \mathbf{Q}_v \mathbf{Q}_v^T \mathbf{Q}_u) = \text{Tr}(\mathbf{Q}_v \mathbf{Q}_v^T \mathbf{Q}_u \mathbf{Q}_u^T). \quad (23)$$

523 As defined above, we have

$$\begin{aligned} Q_u Q_u^\top &= \mathbf{Z}_u (\mathbf{Z}_u^\top \mathbf{Z}_u)^{-\frac{1}{2}} (\mathbf{Z}_u^\top \mathbf{Z}_u)^{-\frac{1}{2}} \mathbf{Z}_u^\top = \mathbf{Z}_u (\mathbf{Z}_u^\top \mathbf{Z}_u)^{-1} \mathbf{Z}_u^\top \\ Q_v Q_v^\top &= \mathbf{Z}_v (\mathbf{Z}_v^\top \mathbf{Z}_v)^{-\frac{1}{2}} (\mathbf{Z}_v^\top \mathbf{Z}_v)^{-\frac{1}{2}} \mathbf{Z}_v^\top = \mathbf{Z}_v (\mathbf{Z}_v^\top \mathbf{Z}_v)^{-1} \mathbf{Z}_v^\top. \end{aligned} \quad (24)$$

524 Then Equation 23 becomes,

$$\begin{aligned} \mathbf{Tr}(\Lambda_{u,v} \Lambda_{u,v}^\top) &= \mathbf{Tr}(\mathbf{Z}_u (\mathbf{Z}_u^\top \mathbf{Z}_u)^{-1} \mathbf{Z}_u^\top \mathbf{Z}_v (\mathbf{Z}_v^\top \mathbf{Z}_v)^{-1} \mathbf{Z}_v^\top) \\ &= \mathbf{Tr}((\mathbf{Z}_u^\top \mathbf{Z}_u)^{-1} \mathbf{Z}_u^\top \mathbf{Z}_v (\mathbf{Z}_v^\top \mathbf{Z}_v)^{-1} \mathbf{Z}_v^\top \mathbf{Z}_u). \end{aligned} \quad (25)$$

525 By Cauchy-Schwartz Inequality,

$$\mathbf{Tr}(\Lambda_{u,v} \Lambda_{u,v}^\top) \leq \mathbf{Tr}((\mathbf{Z}_u^\top \mathbf{Z}_u)^{-1}) \mathbf{Tr}((\mathbf{Z}_v^\top \mathbf{Z}_v)^{-1}) \mathbf{Tr}(\mathbf{Z}_u^\top \mathbf{Z}_v)^2. \quad (26)$$

526 Then we analyze these terms individually,

$$\begin{aligned} \mathbf{Tr}(\mathbf{Z}_u^\top \mathbf{Z}_v) &= \mathbf{Tr}(\mathbf{D}_u^\top \mathbf{X}^\top \alpha^\top \alpha \mathbf{X} \mathbf{D}_v) = \mathbf{Tr}(\mathbf{X}^\top \alpha^\top \alpha \mathbf{X} \mathbf{D}_v \mathbf{D}_u^\top) \\ &\leq \mathbf{Tr}(\mathbf{X}^\top \alpha^\top \alpha \mathbf{X}) \mathbf{Tr}(\mathbf{D}_u^\top \mathbf{D}_v) \leq \mathcal{T} \cdot \mathbf{Tr}(\mathbf{D}_u^\top \mathbf{D}_v) \end{aligned} \quad (27)$$

527 As for $\mathbf{Tr}((\mathbf{Z}_u^\top \mathbf{Z}_u)^{-1})$, let $\lambda_1, \lambda_2, \dots, \lambda_c$ be eigenvalues for $\mathbf{Z}_u^\top \mathbf{Z}_u$ listed in descending order ($\lambda_1 \geq$
528 $\lambda_2 \geq \dots \geq \lambda_c$), and assume the condition number of $\mathbf{Z}_u^\top \mathbf{Z}_u$ and $\mathbf{Z}_v^\top \mathbf{Z}_v$ satisfy $\lambda_{max}/\lambda_{min} \leq \gamma$, then,

$$\mathbf{Tr}((\mathbf{Z}_u^\top \mathbf{Z}_u)^{-1}) = \sum_{i=1}^c \frac{1}{\lambda_i} \leq c \cdot \frac{1}{\lambda_c} \leq \frac{\gamma c}{\lambda_1}, \quad (28)$$

529 where $\lambda_1 = \|\mathbf{Z}_u^\top \mathbf{Z}_u\|_2$, $\|\cdot\|_2$ denotes the operator norm induced by the vector L_2 -norm. With the
530 norm inequalities of any positive semidefinite matrix A ,

$$\|A\|_2 \geq \frac{1}{\sqrt{c}} \|A\|_F \geq \frac{1}{c} \|A\|_* \geq \frac{1}{c} \mathbf{Tr}(A), \quad (29)$$

531 where $\|\cdot\|_F, \|\cdot\|_*$ denote the Frobenius norm and the nuclear norm, respectively.

532 Equation (30) then becomes,

$$\mathbf{Tr}((\mathbf{Z}_u^\top \mathbf{Z}_u)^{-1}) \leq c \cdot \frac{1}{\|\mathbf{Z}_u^\top \mathbf{Z}_u\|_2} \leq \frac{\gamma c^2}{\mathbf{Tr}(\mathbf{Z}_u^\top \mathbf{Z}_u)}. \quad (30)$$

533 By Lemma A.3,

$$\begin{aligned} \mathbf{Tr}(\mathbf{Z}_u^\top \mathbf{Z}_u) &= \mathbf{Tr}(\mathbf{D}_u^\top \mathbf{X}^\top \alpha^\top \alpha \mathbf{X} \mathbf{D}_u) \\ &= \mathbf{Tr}(\mathbf{X}^\top \alpha^\top \alpha \mathbf{X} \mathbf{D}_u \mathbf{D}_u^\top) \\ &\geq \sigma_{min}(\mathbf{X}^\top \alpha^\top \alpha \mathbf{X}) \mathbf{Tr}(\mathbf{D}_u^\top \mathbf{D}_u) \\ &\geq \mathcal{C} \cdot \mathbf{Tr}(\mathbf{D}_u^\top \mathbf{D}_u) \\ &\geq \mathcal{C} \cdot \|\mathit{vec}(\mathbf{D}_u)\|_2^2, \end{aligned} \quad (31)$$

534 where $\mathit{vec}(\cdot)$ denotes vectorization of a matrix.

535 Then Equation 30 is further derived as,

$$\mathbf{Tr}((\mathbf{Z}_u^\top \mathbf{Z}_u)^{-1}) \leq \frac{\gamma c^2}{\mathcal{C} \cdot \|\mathit{vec}(\mathbf{D}_u)\|_2^2}. \quad (32)$$

536 Similarly, we have

$$\mathbf{Tr}((\mathbf{Z}_v^\top \mathbf{Z}_v)^{-1}) \leq \frac{\gamma c^2}{\mathcal{C} \cdot \|\mathit{vec}(\mathbf{D}_v)\|_2^2}. \quad (33)$$

537 Finally, with $\text{Tr}(\mathbf{D}_u^\top \mathbf{D}_v) = \langle \text{vec}(\mathbf{D}_u), \text{vec}(\mathbf{D}_v) \rangle$, we have

$$\begin{aligned} \text{Tr}(\Lambda_{u,v} \Lambda_{u,v}^\top) &\leq \frac{\gamma^2 \mathcal{T}^2 c^4 (\langle \text{vec}(\mathbf{D}_u), \text{vec}(\mathbf{D}_v) \rangle)^2}{\mathcal{C}^2 \|\text{vec}(\mathbf{D}_u)\|_2^2 \cdot \|\text{vec}(\mathbf{D}_v)\|_2^2} \\ &\leq \frac{\gamma^2 \mathcal{T}^2 c^4}{\mathcal{C}^2} \cdot \cos^2(\mathbf{D}_u, \mathbf{D}_v), \end{aligned} \quad (34)$$

538 and thus,

$$\begin{aligned} \mathcal{S}(\mathbf{Z}_u, \mathbf{Z}_v) &= \sqrt{\frac{1}{\mathcal{C}} \text{Tr}(\Lambda_{u,v} \Lambda_{u,v}^\top)} \\ &\leq \frac{\gamma \mathcal{T} c^{\frac{3}{2}}}{\mathcal{C}} \cdot \cos(\mathbf{D}_u, \mathbf{D}_v). \end{aligned} \quad (35)$$

539 Then the claimed theorem is proved. □

540

541 **Lemma A.5.** For two matrices \mathbf{A} , \mathbf{B} , their frobenius norm satisfies,

$$\|\mathbf{A}\mathbf{B}\|_F = \|\mathbf{A}\|_F \|\mathbf{B}\|_F \sqrt{1 - \frac{\Delta_1}{\|\mathbf{A}\|_F^2 \|\mathbf{B}\|_F^2}}, \quad (36)$$

542 where $\Delta_1 = \sum_{ij} (\sum_k A_{ik}^2) (\sum_k B_{kj}^2) \cdot \sin^2(\langle A_{i:}, B_{:j} \rangle)$.

543 *Proof.* According to the definition of frobenius norm $\|\mathbf{A}\|_F = \sqrt{\sum_{ij} |A_{ij}|^2}$ we have,

$$\|\mathbf{A}\mathbf{B}\|_F = \sqrt{\sum_{ij} (\sum_k A_{ik} B_{kj})^2}. \quad (37)$$

544 Note that $(\sum_i x_i y_i)^2 = (\sum_i x_i^2) (\sum_i y_i^2) \cdot \cos^2(\langle x, y \rangle) = (\sum_i x_i^2) (\sum_i y_i^2) - (\sum_i x_i^2) (\sum_i y_i^2) \cdot$
545 $\sin^2(\langle x, y \rangle)$, where $\langle x, y \rangle$ is the angle of two vectors x and y . We have,

$$\begin{aligned} &\sqrt{\sum_{ij} (\sum_k A_{ik} B_{kj})^2} \\ &= \sqrt{\sum_{ij} \left[(\sum_k A_{ik}^2) (\sum_k B_{kj}^2) - (\sum_k A_{ik}^2) (\sum_k B_{kj}^2) \cdot \sin^2(\langle A_{i:}, B_{:j} \rangle) \right]} \\ &= \sqrt{\sum_{ik} A_{ik}^2} \sqrt{\sum_{kj} B_{kj}^2} \sqrt{1 - \frac{\sum_{ij} (\sum_k A_{ik}^2) (\sum_k B_{kj}^2) \cdot \sin^2(\langle A_{i:}, B_{:j} \rangle)}{\sum_{ik} A_{ik}^2 \sum_{kj} B_{kj}^2}} \\ &= \|\mathbf{A}\|_F \|\mathbf{B}\|_F \sqrt{1 - \frac{\sum_{ij} (\sum_k A_{ik}^2) (\sum_k B_{kj}^2) \cdot \sin^2(\langle A_{i:}, B_{:j} \rangle)}{\|\mathbf{A}\|_F^2 \|\mathbf{B}\|_F^2}} \\ &= \|\mathbf{A}\|_F \|\mathbf{B}\|_F \sqrt{1 - \frac{\Delta_1}{\|\mathbf{A}\|_F^2 \|\mathbf{B}\|_F^2}}, \end{aligned} \quad (38)$$

546 where $A_{i:}$ is the i -th row of \mathbf{A} and $B_{:j}$ is the j -th column of \mathbf{B} , $\Delta_1 = \sum_{ij} (\sum_k A_{ik}^2) (\sum_k B_{kj}^2) \cdot$
547 $\sin^2(\langle A_{i:}, B_{:j} \rangle)$. As $A_{i:}$ and $B_{:j}$ are more correlated, $\langle A_{i:}, B_{:j} \rangle \rightarrow 0$, thus, $\Delta_1 \ll \|\mathbf{A}\|_F^2 \|\mathbf{B}\|_F^2$. □

548

Lemma A.6.

$$\|\mathbf{A}^{1/2}\|_F = \|\mathbf{A}\|_F^{1/2} \left(1 + \frac{\Delta_1 \mathbf{A}^{1/2}}{\|\mathbf{A}\|_F^2}\right)^{1/4}. \quad (39)$$

549 *Proof.* According to Lemma A.5, we have,

$$\|\mathbf{A}\|_F^2 = \|\mathbf{A}^{1/2}\|_F^4 - \Delta_1. \quad (40)$$

550 Thus,

$$\|\mathbf{A}^{1/2}\|_F = \|\mathbf{A}\|_F^{1/2} \left(1 + \frac{\Delta_{1A^{1/2}}}{\|\mathbf{A}\|_F^2}\right)^{1/4}, \quad (41)$$

551 where $\Delta_{1A^{1/2}} = \sum_{ij} (\sum_k (A^{1/2})_{ik}^2) (\sum_k (A^{1/2})_{kj}^2) \cdot \sin^2(\langle (A^{1/2})_{i:}, (A^{1/2})_{:j} \rangle)$. As $(A^{1/2})_{i:}$ and
552 $(A^{1/2})_{:j}$ are more correlated, $\langle (A^{1/2})_{i:}, (A^{1/2})_{:j} \rangle \rightarrow 0$, thus, $\Delta_{1A^{1/2}} \ll \|\mathbf{A}\|_F^2$.

553

□

554 **Lemma A.7.** For three matrices \mathbf{A} , \mathbf{B} , and \mathbf{C} , their frobenius norm satisfies,

$$\|\mathbf{A}\|_F = \|\mathbf{A}\|_F \|\mathbf{B}\|_F \|\mathbf{C}\|_F \sqrt{1 - \frac{\Delta_2 + \Delta_3}{\|\mathbf{A}\|_F^2 \|\mathbf{B}\|_F^2 \|\mathbf{C}\|_F^2}}, \quad (42)$$

555 where $\Delta_2 = \frac{1}{2} [\|\mathbf{A}\|_F^2 \sum_{kj} (\sum_l B_{kl}^2) (\sum_l C_{lj}^2) \cdot \sin^2(\langle B_{k:}, C_{:j} \rangle) + \|\mathbf{C}\|_F^2 \sum_{il} (\sum_k A_{ik}^2) (\sum_k B_{kl}^2) \cdot$
556 $\sin^2(\langle A_{i:}, B_{:l} \rangle)]$ and $\Delta_3 = \frac{1}{2} [\sum_{ij} (\sum_k A_{ik}^2) (\sum_k (BC)_{kj}^2) \cdot \sin^2(\langle A_{i:}, (BC)_{:j} \rangle) +$
557 $\sum_{ij} (\sum_l (AB)_{il}^2) (\sum_l C_{lj}^2) \cdot \sin^2(\langle (AB)_{i:}, C_{:j} \rangle)]$.

558 *Proof.* Based on Lemma A.5, we have,

$$\begin{aligned} & \|\mathbf{ABC}\|_F^2 \\ &= \|\mathbf{AB}\|_F^2 \|\mathbf{C}\|_F^2 - \sum_{ij} (\sum_l (AB)_{il}^2) (\sum_l C_{lj}^2) \cdot \sin^2(\langle (AB)_{i:}, C_{:j} \rangle) \\ &= \|\mathbf{A}\|_F^2 \|\mathbf{B}\|_F^2 \|\mathbf{C}\|_F^2 - \|\mathbf{C}\|_F^2 \sum_{il} (\sum_k A_{ik}^2) (\sum_k B_{kl}^2) \cdot \sin^2(\langle A_{i:}, B_{:l} \rangle) \\ & \quad - \sum_{ij} (\sum_l (AB)_{il}^2) (\sum_l C_{lj}^2) \cdot \sin^2(\langle (AB)_{i:}, C_{:j} \rangle) \end{aligned} \quad (43)$$

559 Symmetrically, we also have,

$$\begin{aligned} & \|\mathbf{ABC}\|_F^2 \\ &= \|\mathbf{A}\|_F^2 \|\mathbf{BC}\|_F^2 - \sum_{ij} (\sum_k A_{ik}^2) (\sum_k (BC)_{kj}^2) \cdot \sin^2(\langle A_{i:}, (BC)_{:j} \rangle) \\ &= \|\mathbf{A}\|_F^2 \|\mathbf{B}\|_F^2 \|\mathbf{C}\|_F^2 - \|\mathbf{A}\|_F^2 \sum_{kj} (\sum_l B_{kl}^2) (\sum_l C_{lj}^2) \cdot \sin^2(\langle B_{k:}, C_{:j} \rangle) \\ & \quad - \sum_{ij} (\sum_k A_{ik}^2) (\sum_k (BC)_{kj}^2) \cdot \sin^2(\langle A_{i:}, (BC)_{:j} \rangle) \end{aligned} \quad (44)$$

560 Thus,

$$\begin{aligned} & \|\mathbf{ABC}\|_F^2 \\ &= \frac{1}{2} [\|\mathbf{A}\|_F^2 \|\mathbf{B}\|_F^2 \|\mathbf{C}\|_F^2 - \|\mathbf{A}\|_F^2 \sum_{kj} (\sum_l B_{kl}^2) (\sum_l C_{lj}^2) \cdot \sin^2(\langle B_{k:}, C_{:j} \rangle) \\ & \quad - \sum_{ij} (\sum_k A_{ik}^2) (\sum_k (BC)_{kj}^2) \cdot \sin^2(\langle A_{i:}, (BC)_{:j} \rangle) \\ & \quad + \|\mathbf{A}\|_F^2 \|\mathbf{B}\|_F^2 \|\mathbf{C}\|_F^2 - \|\mathbf{C}\|_F^2 \sum_{il} (\sum_k A_{ik}^2) (\sum_k B_{kl}^2) \cdot \sin^2(\langle A_{i:}, B_{:l} \rangle) \\ & \quad - \sum_{ij} (\sum_l (AB)_{il}^2) (\sum_l C_{lj}^2) \cdot \sin^2(\langle (AB)_{i:}, C_{:j} \rangle)] \\ &= \|\mathbf{A}\|_F^2 \|\mathbf{B}\|_F^2 \|\mathbf{C}\|_F^2 - \Delta_2 - \Delta_3, \end{aligned} \quad (45)$$

561 where $\Delta_2 = \frac{1}{2}[\|A\|_F^2 \sum_{kj} (\sum_l B_{kl}^2)(\sum_l C_{lj}^2) \cdot \sin^2(\langle B_{k:}, C_{:j} \rangle) + \|C\|_F^2 \sum_{il} (\sum_k A_{ik}^2)(\sum_k B_{kl}^2) \cdot$
562 $\sin^2(\langle A_{i:}, B_{:l} \rangle)]$ and $\Delta_3 = \frac{1}{2}[\sum_{ij} (\sum_k A_{ik}^2)(\sum_k (BC)_{kj}^2) \cdot \sin^2(\langle A_{i:}, (BC)_{:j} \rangle) +$
563 $\sum_{ij} (\sum_l (AB)_{il}^2)(\sum_l C_{lj}^2) \cdot \sin^2(\langle (AB)_{i:}, C_{:j} \rangle)]$. Therefore,

$$\|\mathbf{ABC}\|_F = \|\mathbf{A}\|_F \|\mathbf{B}\|_F \|\mathbf{C}\|_F \sqrt{1 - \frac{\Delta_2 + \Delta_3}{\|\mathbf{A}\|_F^2 \|\mathbf{B}\|_F^2 \|\mathbf{C}\|_F^2}}. \quad (46)$$

564 As $A_{i:}$ and $B_{:l}$, $B_{k:}$ and $C_{:j}$ are more correlated, $\langle A_{i:}, B_{:l} \rangle$, $\langle B_{k:}, C_{:j} \rangle$, $\langle A_{i:}, (BC)_{:j} \rangle$, $\langle (AB)_{i:}, C_{:j} \rangle \rightarrow$
565 0, thus, $\Delta_2 \ll \|\mathbf{A}\|_F^2 \|\mathbf{B}\|_F^2 \|\mathbf{C}\|_F^2$ and $\Delta_3 \ll \|\mathbf{A}\|_F^2 \|\mathbf{B}\|_F^2 \|\mathbf{C}\|_F^2$.

566 \square

Lemma A.8.

$$\|\mathbf{A}^{-1/2} \mathbf{B} \mathbf{C}^{-1/2}\|_F = \kappa_F(\mathbf{A}^{1/2}) \kappa_F(\mathbf{C}^{1/2}) \frac{\|\mathbf{B}\|_F}{\|\mathbf{A}^{1/2}\|_F \|\mathbf{C}^{1/2}\|_F} \sqrt{1 - \frac{\Delta_2 + \Delta_3}{\|\mathbf{A}^{-1/2}\|_F^2 \|\mathbf{B}\|_F^2 \|\mathbf{C}^{-1/2}\|_F^2}}, \quad (47)$$

567 where $\kappa_F(\mathbf{A}^{1/2})$ and $\kappa_F(\mathbf{C}^{1/2})$ are the condition number of $\mathbf{A}^{1/2}$ and $\mathbf{C}^{1/2}$, $\kappa_F(\mathbf{A}^{1/2}) =$
568 $\sqrt{(\sum \sigma_i^2(\mathbf{A}^{1/2}))(\sum \frac{1}{\sigma_i^2(\mathbf{A}^{1/2}))}$ and $\kappa_F(\mathbf{C}^{1/2}) = \sqrt{(\sum \sigma_i^2(\mathbf{C}^{1/2}))(\sum \frac{1}{\sigma_i^2(\mathbf{C}^{1/2}))}$; $\sigma_i^2(\mathbf{A}^{1/2})$ are sin-
569 gular value of $\mathbf{A}^{1/2}$ and $\sigma_i^2(\mathbf{C}^{1/2})$ are singular value of $\mathbf{C}^{1/2}$.

570 *Proof.* Based on Lemma A.7, we have,

$$\|\mathbf{A}^{-1/2} \mathbf{B} \mathbf{C}^{-1/2}\|_F = \|\mathbf{A}^{-1/2}\|_F \|\mathbf{B}\|_F \|\mathbf{C}^{-1/2}\|_F \sqrt{1 - \frac{\Delta_2 + \Delta_3}{\|\mathbf{A}^{-1/2}\|_F^2 \|\mathbf{B}\|_F^2 \|\mathbf{C}^{-1/2}\|_F^2}}. \quad (48)$$

571 By the definition of condition number $\kappa_F(\mathbf{X}) = \|\mathbf{X}\|_F \|\mathbf{X}^{-1}\|_F = \sqrt{(\sum \sigma_i^2(\mathbf{X}))(\sum \frac{1}{\sigma_i^2(\mathbf{X}))}$,

$$\|\mathbf{A}^{-1/2} \mathbf{B} \mathbf{C}^{-1/2}\|_F = \kappa_F(\mathbf{A}^{1/2}) \kappa_F(\mathbf{C}^{1/2}) \frac{\|\mathbf{B}\|_F}{\|\mathbf{A}^{1/2}\|_F \|\mathbf{C}^{1/2}\|_F} \sqrt{1 - \frac{\Delta_2 + \Delta_3}{\|\mathbf{A}^{-1/2}\|_F^2 \|\mathbf{B}\|_F^2 \|\mathbf{C}^{-1/2}\|_F^2}}. \quad (49)$$

572 \square

573 **Theorem A.9.** Suppose the forward of decomposed convolution layer for the u -th model is $\mathbf{Z}_u =$
574 $\alpha \mathbf{X} \mathbf{D}_u$, CCA coefficient be $S(\mathbf{Z}_u, \mathbf{Z}_v) = \sqrt{\frac{1}{c} \sum_{i=1}^c \sigma_i^2}$, where σ_i^2 denotes the i -th eigenvalue of
575 $\Lambda_{u,v} = Q_u^\top Q_v$, $Q_u = \mathbf{Z}_u (\mathbf{Z}_u^\top \mathbf{Z}_u)^{-\frac{1}{2}}$. Then $S(\mathbf{Z}_u, \mathbf{Z}_v)$ is approximately linear to filter subspace
576 similarity,

$$S(\mathbf{Z}_u, \mathbf{Z}_v) = \frac{\gamma_1 \gamma_2 \gamma_3}{\sqrt{c}} \cos(\mathbf{D}_u, \mathbf{D}_v), \quad (50)$$

577 *Proof.* Based on $S(\mathbf{Z}_u, \mathbf{Z}_v) = \sqrt{\frac{1}{c} \sum_{i=1}^c \sigma_i^2}$ and $\|\Lambda_{u,v}\|_F = \sqrt{\sum_{i=1}^c \sigma_i^2}$, where σ_i are the singular
578 value of $\Lambda_{u,v}$,

$$S = \sqrt{\frac{1}{c} \sum_{i=1}^c \sigma_i^2} = \frac{1}{\sqrt{c}} \|\Lambda_{u,v}\|_F = \frac{1}{\sqrt{c}} \|(\mathbf{Z}_u^\top \mathbf{Z}_u)^{-\frac{1}{2}} \mathbf{Z}_u^\top \mathbf{Z}_v (\mathbf{Z}_v^\top \mathbf{Z}_v)^{-\frac{1}{2}}\|_F. \quad (51)$$

579 According to Lemma. A.8, we have

$$\frac{1}{\sqrt{c}} \|(\mathbf{Z}_u^\top \mathbf{Z}_u)^{-\frac{1}{2}} \mathbf{Z}_u^\top \mathbf{Z}_v (\mathbf{Z}_v^\top \mathbf{Z}_v)^{-\frac{1}{2}}\|_F = \frac{\gamma_1 \gamma_2}{\sqrt{c}} \frac{\|\mathbf{Z}_u^\top \mathbf{Z}_v\|_F}{\|(\mathbf{Z}_u^\top \mathbf{Z}_u)^{\frac{1}{2}}\|_F \|(\mathbf{Z}_v^\top \mathbf{Z}_v)^{\frac{1}{2}}\|_F}, \quad (52)$$

580 where $\gamma_1 = \kappa_F((\mathbf{Z}_u^\top \mathbf{Z}_u)^{\frac{1}{2}}) \cdot \kappa_F((\mathbf{Z}_v^\top \mathbf{Z}_v)^{\frac{1}{2}})$ and $\gamma_2 = \sqrt{1 - \frac{\Delta_2 + \Delta_3}{\|(\mathbf{Z}_u^\top \mathbf{Z}_u)^{-1/2}\|_F^2 \|\mathbf{Z}_u^\top \mathbf{Z}_v\|_F^2 \|(\mathbf{Z}_v^\top \mathbf{Z}_v)^{-1/2}\|_F^2}}$.

581 As $\mathbf{Z}_u = \alpha \mathbf{X} \mathbf{D}_u$ and $\mathbf{Z}_v = \alpha \mathbf{X} \mathbf{D}_v$, we have

$$\begin{aligned} & \frac{\gamma_1 \gamma_2}{\sqrt{c}} \frac{\|\mathbf{Z}_u^\top \mathbf{Z}_v\|_F}{\|(\mathbf{Z}_u^\top \mathbf{Z}_u)^{\frac{1}{2}}\|_F \|(\mathbf{Z}_v^\top \mathbf{Z}_v)^{\frac{1}{2}}\|_F} \\ &= \frac{\gamma_1 \gamma_2}{\sqrt{c}} \frac{\|\mathbf{D}_u^\top \mathbf{X}^\top \alpha^\top \alpha \mathbf{X} \mathbf{D}_v\|_F}{\|(\mathbf{D}_u^\top \mathbf{X}^\top \alpha^\top \alpha \mathbf{X} \mathbf{D}_u)^{\frac{1}{2}}\|_F \|(\mathbf{D}_v^\top \mathbf{X}^\top \alpha^\top \alpha \mathbf{X} \mathbf{D}_v)^{\frac{1}{2}}\|_F}. \end{aligned} \quad (53)$$

582 According to Lemma A.6,

$$\begin{aligned} & \frac{\gamma_1 \gamma_2}{\sqrt{c}} \frac{\|\mathbf{D}_u^\top \mathbf{X}^\top \alpha^\top \alpha \mathbf{X} \mathbf{D}_v\|_F}{\|(\mathbf{D}_u^\top \mathbf{X}^\top \alpha^\top \alpha \mathbf{X} \mathbf{D}_u)^{\frac{1}{2}}\|_F \|(\mathbf{D}_v^\top \mathbf{X}^\top \alpha^\top \alpha \mathbf{X} \mathbf{D}_v)^{\frac{1}{2}}\|_F} \\ &= \frac{\gamma_1 \gamma_2 \gamma_3}{\sqrt{c}} \frac{\|\mathbf{D}_u^\top \mathbf{X}^\top \alpha^\top \alpha \mathbf{X} \mathbf{D}_v\|_F}{\|(\mathbf{D}_u^\top \mathbf{X}^\top \alpha^\top \alpha \mathbf{X} \mathbf{D}_u)^{\frac{1}{2}}\|_F \|(\mathbf{D}_v^\top \mathbf{X}^\top \alpha^\top \alpha \mathbf{X} \mathbf{D}_v)^{\frac{1}{2}}\|_F}, \end{aligned} \quad (54)$$

583 where $\gamma_3 = (1 + \frac{\Delta_1}{\|(\mathbf{D}_u^\top \mathbf{X}^\top \alpha^\top \alpha \mathbf{X} \mathbf{D}_u)\|_F^2})^{-\frac{1}{4}} (1 + \frac{\Delta_1}{\|(\mathbf{D}_v^\top \mathbf{X}^\top \alpha^\top \alpha \mathbf{X} \mathbf{D}_v)\|_F^2})^{-\frac{1}{4}}$.

584 As Assumption 2.6 holds, it becomes

$$\begin{aligned} & \frac{\gamma_1 \gamma_2 \gamma_3}{\sqrt{c}} \frac{\|\mathbf{D}_u^\top \mathbf{X}^\top \alpha^\top \alpha \mathbf{X} \mathbf{D}_v\|_F}{\|(\mathbf{D}_u^\top \mathbf{X}^\top \alpha^\top \alpha \mathbf{X} \mathbf{D}_u)^{\frac{1}{2}}\|_F \|(\mathbf{D}_v^\top \mathbf{X}^\top \alpha^\top \alpha \mathbf{X} \mathbf{D}_v)^{\frac{1}{2}}\|_F} \\ &= \frac{\gamma_1 \gamma_2 \gamma_3}{\sqrt{c}} \frac{\|\mathbf{D}_u^\top \mathbf{D}_v\|_F \|\mathbf{X}^\top \alpha^\top \alpha \mathbf{X}\|_F}{\|\mathbf{D}_u\|_F^{\frac{1}{2}} \|\mathbf{X}^\top \alpha^\top \alpha \mathbf{X}\|_F^{\frac{1}{2}} \|\mathbf{D}_u\|_F^{\frac{1}{2}} \|\mathbf{D}_v\|_F^{\frac{1}{2}} \|\mathbf{X}^\top \alpha^\top \alpha \mathbf{X}\|_F^{\frac{1}{2}} \|\mathbf{D}_v\|_F^{\frac{1}{2}}} \\ &= \frac{\gamma_1 \gamma_2 \gamma_3}{\sqrt{c}} \frac{\|\mathbf{D}_u^\top \mathbf{D}_v\|_F}{\|\mathbf{D}_u\|_F \|\mathbf{D}_v\|_F} \\ &= \frac{\gamma_1 \gamma_2 \gamma_3}{\sqrt{c}} \cos(\mathbf{D}_u, \mathbf{D}_v). \end{aligned} \quad (55)$$

585 Thus, we have

$$\mathcal{S}(\mathbf{Z}_u, \mathbf{Z}_v) = \frac{\gamma_1 \gamma_2 \gamma_3}{\sqrt{c}} \cos(\mathbf{D}_u, \mathbf{D}_v). \quad (56)$$

586 Specifically, we have $\gamma_2 = \sqrt{1 - \frac{\Delta}{\gamma_1^2 \gamma_3^2 \cos^2(\mathbf{D}_u, \mathbf{D}_v)}}$, and since Δ are small, with Taylor expansion,
 587 $\gamma_2 \approx 1 - \frac{1}{2} \frac{\Delta}{\gamma_1^2 \gamma_3^2 \cos^2(\mathbf{D}_u, \mathbf{D}_v)}$. The term $\frac{1}{\cos^2(\mathbf{D}_u, \mathbf{D}_v)}$ causes non-linearity in the relation between
 588 CCA and filter subspace similarity. \square

589 A.2 Experiment Settings

590 **Model training of Federated Learning.** In each experiment we have 100 clients in total and
 591 sample a ratio $r = 0.1$ of all the clients on every round. All models are randomly initialized and
 592 trained for $T = 100$ communication rounds for the CIFAR datasets. At each round, the client
 593 executes 15 epochs of SGD with momentum to train the local model, the learning rate is 0.01 and
 594 momentum is 0.9. Accuracies are computed by taking the average local accuracies for all users at the
 595 final communication round. As shown in the Table 3, we have different settings for CIFAR-10 and
 596 CIFAR-100. For example, (100, 2) means 100 clients with 2 classes on each client. For each method,
 597 the training takes about 12 hours on Nvidia RTX A5000.

598 **Comparison with other FL approaches.** We compare our approach by evolving shared atom
 599 coefficients with various personalized federated learning methods and federated learning methods
 600 with local finetuning. Among these methods, FedPer [2] and FedRep[6] have the similar ideas by
 601 learning shared global representation and personalized local heads. Ditto [25] and FedProx [27]
 602 induce global regularization to improve the model performance. We also compare our method with
 603 FedAvg [32]. FedRep [6] approaches the common knowledge with shared representation. The codes

Table 3: Compare accuracy with different approaches

(# client, # classes per client)	CIFAR-100		CIFAR-10		
	(100, 5)	(100, 20)	(100, 2)	(100, 5)	(1000, 2)
FedAvg	82.39	62.92	86.37	70.63	86.12
FedProx	80.77	59.7	85.90	69.94	84.83
FedPer	81.46	62.52	81.74	68.24	81.74
FedRep	72.98	37.71	80.55	67.3	82.98
Local	81.21	49.25	90.24	72.05	97.80
Ours	81.03	52.13	83.37	65.63	82.54

604 are adapted from ¹. We evaluate the test accuracy on CIFAR-10 and CIFAR-100 with different FL
 605 setting. As shown in Table 3, our method achieves comparable performance among different methods.

606 **Fine-tuning models for ensemble.** We select 3 models with different similarity measures for
 607 ensemble. For feature-based similarity methods, we randomly select 1000 examples from CIFAR-100
 608 dataset. The fully-connected layer of each model is fine-tuned on the user’s local data with 100
 609 epochs. The fine-tuning takes about 12 hours on Nvidia RTX A5000. After fine-tuning, the accuracy
 610 is measured on local test data, with the predictions of current model and 3 selected models.

611 **A.3 Extra Experiments**

612 **Representation dependency on filter atoms.** We first validate the dependency of deep features on
 613 filter atoms in Proposition 2.1 with a simple experiment. The model \mathcal{F} here is a 2-layer CNN with
 614 coefficient α and atom \mathbf{D} generated from normal distribution $\mathcal{N}(0, 1)$. The input sample \mathbf{X} is also
 615 generated from normal distribution $\mathcal{N}(0, 1)$. Figure 8(a) shows the relation between $\|\mathbf{Z}_u - \mathbf{Z}_v\|_F$
 616 and $\|\mathbf{D}_u - \mathbf{D}_v\|_F$ by fixing coefficient α and input sample \mathbf{X} and randomly varying filter atoms \mathbf{D} .
 617 All the points are below the line which is the bound provided by Proposition 2.1, reflecting that the
 618 representation variations are dominated by filter atoms.

619 **Correlation between probing-based and filter subspace-based methods.** In addition, we em-
 620 pirically verify that CCA and filter subspace similarity have a strong correlation with AlexNet. In
 621 this experiment, 10 tasks are generated from CIFAR100 [21] with 10 classes in each task. Only the
 622 filter atoms of each task are trained while the atom coefficients are fixed. We calculate CCA and filter
 623 subspace similarity among 45 pairs of models. The correlation between CCA and filter subspace
 624 similarity is 0.8638 which is shown in Figure 9(b). Similarly, the correlation between CKA and filter
 625 subspace similarity is also reported in Figure 9 (Table). These results clearly show that the proposed
 626 filter subspace similarity has high linear relationship with popular probing-based similarities, which
 627 agrees with Theorem 2.5 and Theorem 2.7.

¹<https://github.com/lgcollins/FedRep>

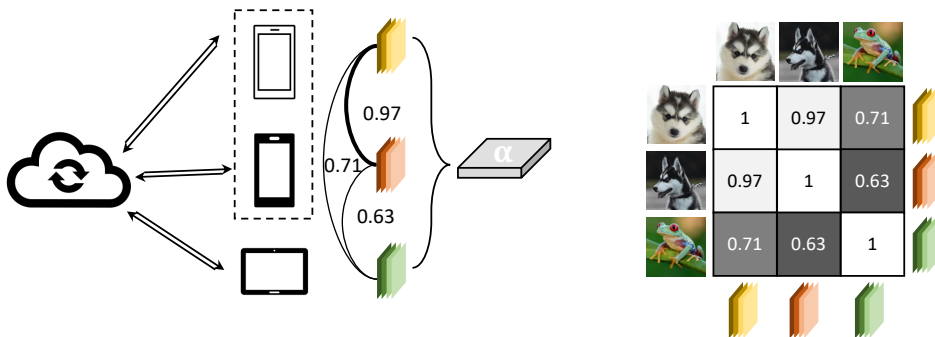


Figure 6: The shared coefficients and user-specific atoms represent common knowledge and personalized information. The filter subspace similarity is used to calculate the relations among users. Users with heterogeneous data result in lower similarity, as illustrated in a similarity matrix.

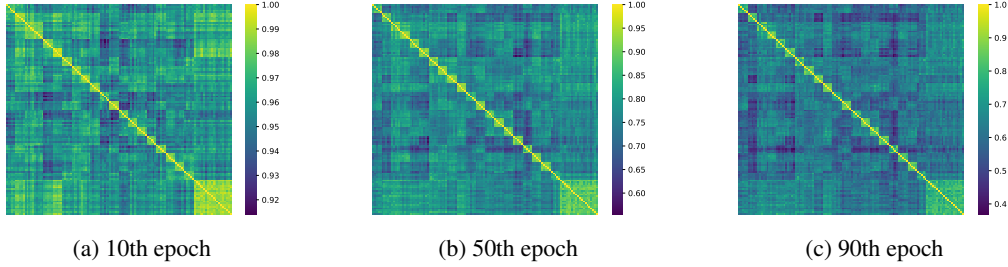


Figure 7: Similarity matrices that show relations among 120 users in FL with our filter subspace similarity through the training process.

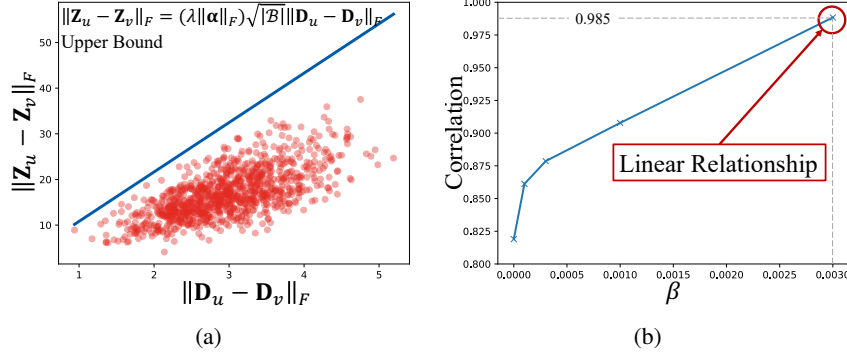


Figure 8: (a) The change of features $\|\mathbf{Z}_u - \mathbf{Z}_v\|_F$ is bounded by the change of atoms $\|\mathbf{D}_u - \mathbf{D}_v\|_F$. (b) The channel decorrelation leads to a higher correlation between CCA and filter subspace similarity. And the correlation can reach 0.985 with $\beta = 3 \times 10^{-3}$, which means a near linear relation between CCA and filter subspace similarity.

628 **Effect of channel decorrelation.** We further design a regularization term $\beta \sum_{i \neq j} (\mathbf{Z}_u^\top \mathbf{Z}_u)_{ij}^2$ to
 629 approach $(\mathbf{Z}_u^\top \mathbf{Z}_u)_{ii} \gg (\mathbf{Z}_u^\top \mathbf{Z}_u)_{ij}$ in Assumption. 2.6. As shown in Figure 8(b), the correlation
 630 between CCA and filter subspace similarity keeps increasing as β increases. The correlation reaches
 631 0.985 when $\beta = 3 \times 10^{-3}$, indicating a near-linear relationship, which is aligned with Theorem. 2.7.

632 **Similar representations across datasets.** Similar to [19], we can use filter subspace similarity to
 633 compare networks trained on different datasets. In Figure 10(a), we show that pairs of models that
 634 are both trained on CIFAR-10 and CIFAR-100 have high atom-based similarities. Models learned
 635 on two datasets respectively still show high similarity. In contrast, similarities between trained and
 636 untrained models are significantly lower.

637 **Limitation of probing-based methods.** As shown in Figure 10(b), to illustrate sensitivity of
 638 probing-based similarities to probing data, we perform a simple regression task with data, $\{(x_i =$
 639 $0, y_i, z_i)\}_{i=1}^n$, where $z_i = f(x_i, y_i) + \epsilon_i$ and $y_i, \epsilon_i \sim \mathcal{N}(0.5, 0.1)$. Two NN models \mathcal{F}_1 and \mathcal{F}_2

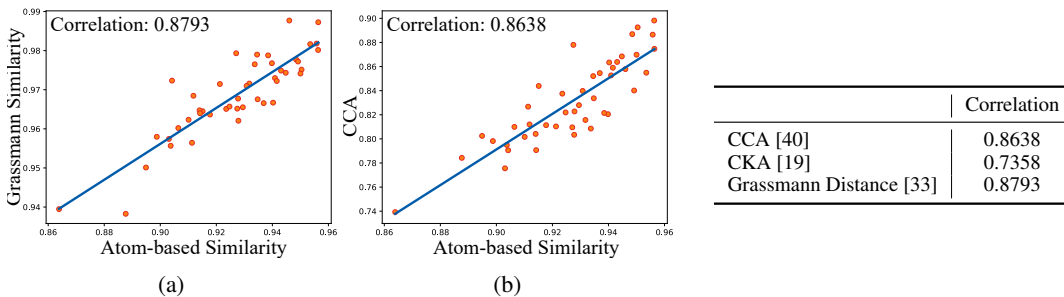


Figure 9: (a) Correlation between Grassmann similarity and filter subspace similarity; (b) Correlation between CCA and filter subspace similarity. (Table) Correlation between filter subspace similarity and other approaches.

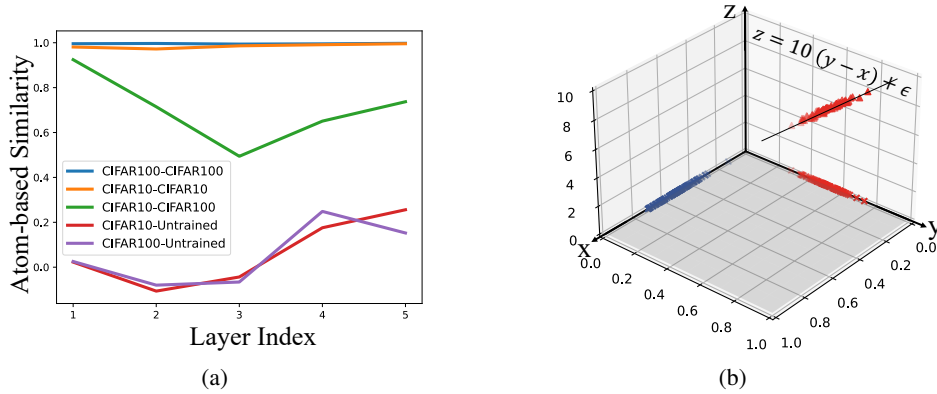


Figure 10: (a) Using filter subspace similarity, models trained on different datasets (CIFAR-10 and CIFAR-100) are similar among themselves, but they differ from untrained models. (b) Illustration of limitations of probing-based similarities. Input data from “red” ($\{(x_i = 0, y_i)\}$) and “blue” ($\{(x'_i = y_i, y'_i = 0)\}$) are orthogonal. Since two models are learned on “red” data, their similarity should be 1, which can be faithfully indicated by our atom similarity. However, probing-based similarities will become 0 with the “blue” probing data.

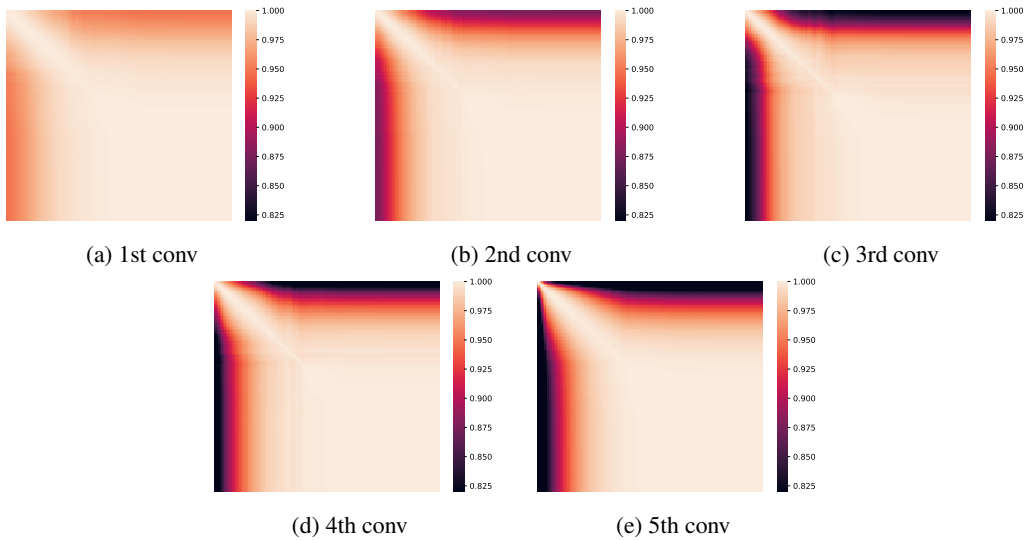


Figure 11: Similarity of AlexNet with atoms from different time point during the training.

640 with the same initialization and atom coefficients are trained for their different atoms to learn
 641 $\mathcal{F} : (X, Y) \rightarrow Z$. It is can be simply found that the filter subspace similarity of \mathcal{F}_1 and \mathcal{F}_2 is 1 and
 642 the probing-based similarity is also 1 with the same $\{(x_i = 0, y_i)\}$ as the probing data. However,
 643 if we choose $\{(x'_i = y_i, y'_i = 0)\}$ as the probing data, then the probing-based similarities directly
 644 become $\mathbf{0}$ as the data are now orthogonal to model parameters.

645 A.4 Training dynamics.

646 We investigate the training dynamics of AlexNet [22] and VGG [47] separately on CIFAR-100 [21]
 647 and ImageNet [44]. The details of training dynamics of models with atoms from different time point
 648 during the training are shown in Figure 11 and Figure 12. Moreover, we examine the similarity
 649 between the two participated models shared the same initialization trained only with atoms on two
 650 different tasks. The results is shown in Figure 13 and Figure 14. The difference is less on the first few
 651 layers, but more on the middle layers. It reflects the middle layer is more critical than other layers,
 652 which is aligned with previous work [36].

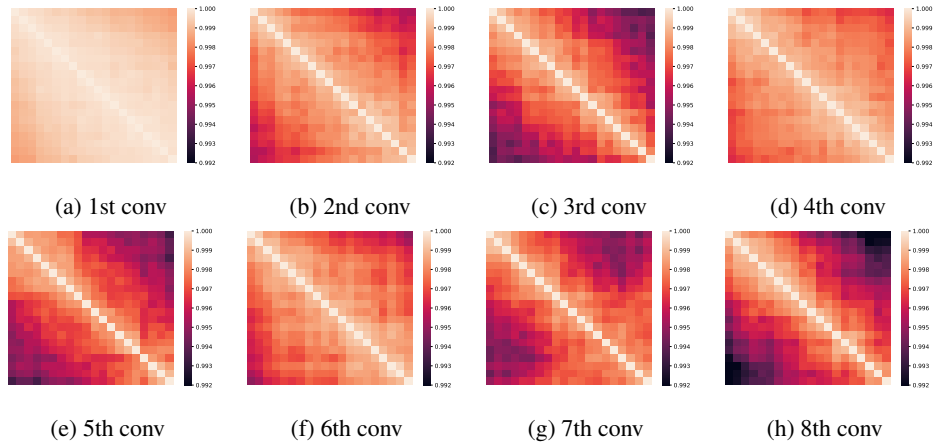


Figure 12: Similarity of VGG with atoms from different time point during the training.

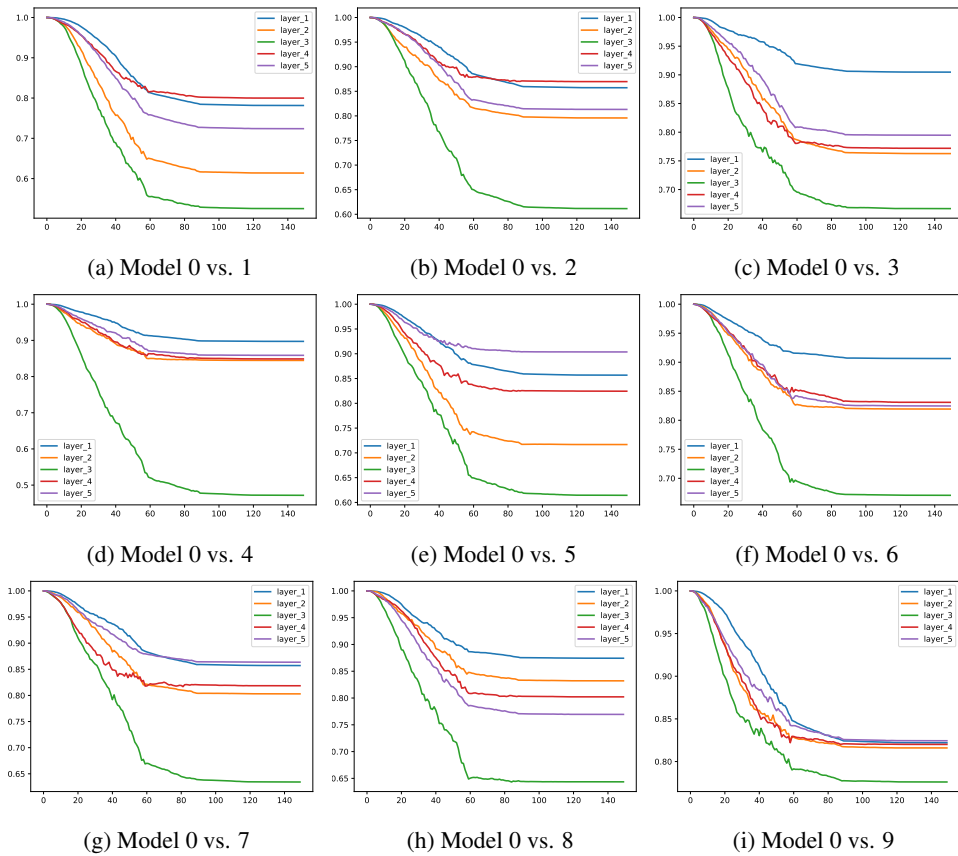


Figure 13: Similarity of AlexNet trained on different tasks during the training.

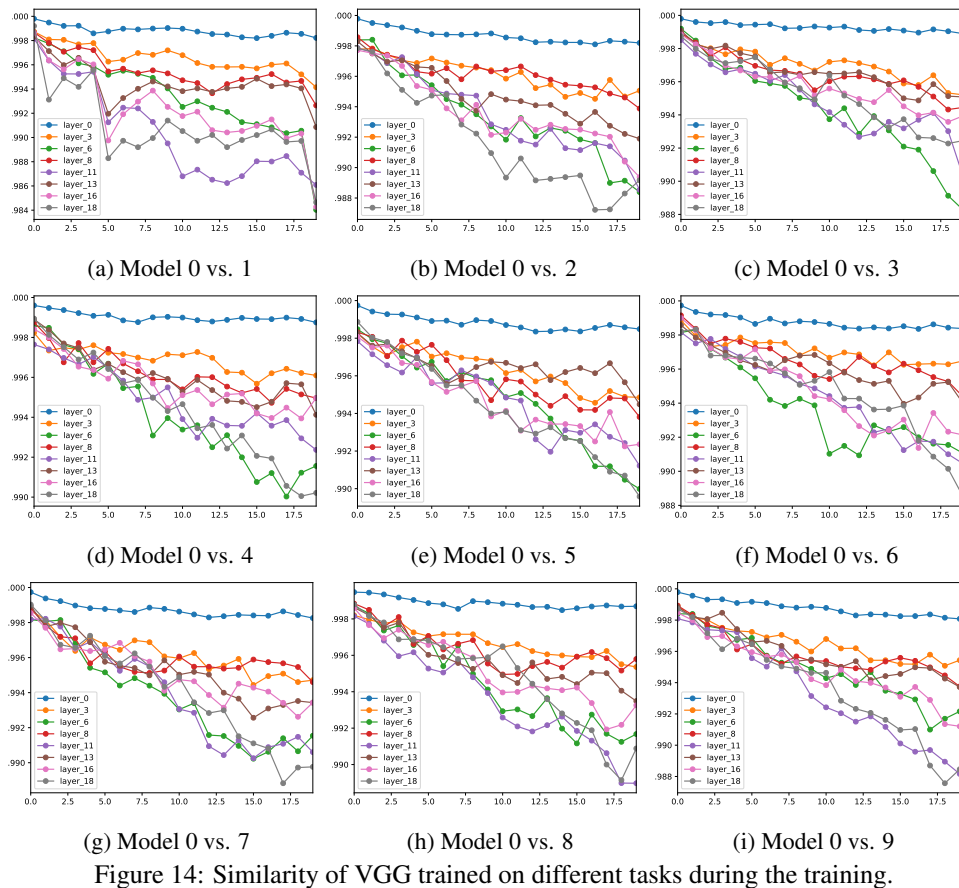


Figure 14: Similarity of VGG trained on different tasks during the training.