



Null-Space Filtering for Data-Free Continual Model Merging: Preserving Stability, Promoting Plasticity

Zihuan Qiu^{1,3}, Lei Wang¹, Yang Cao³, Runlong Zhang¹, Bing Su³, Yi Xu², Fanman Meng¹, Linfeng Xu¹, Qingbo Wu¹, Hongliang Li¹

¹ University of Electronic Science and Technology of China, ² Dalian University of Technology, ³ Jiigan Technology



Paper



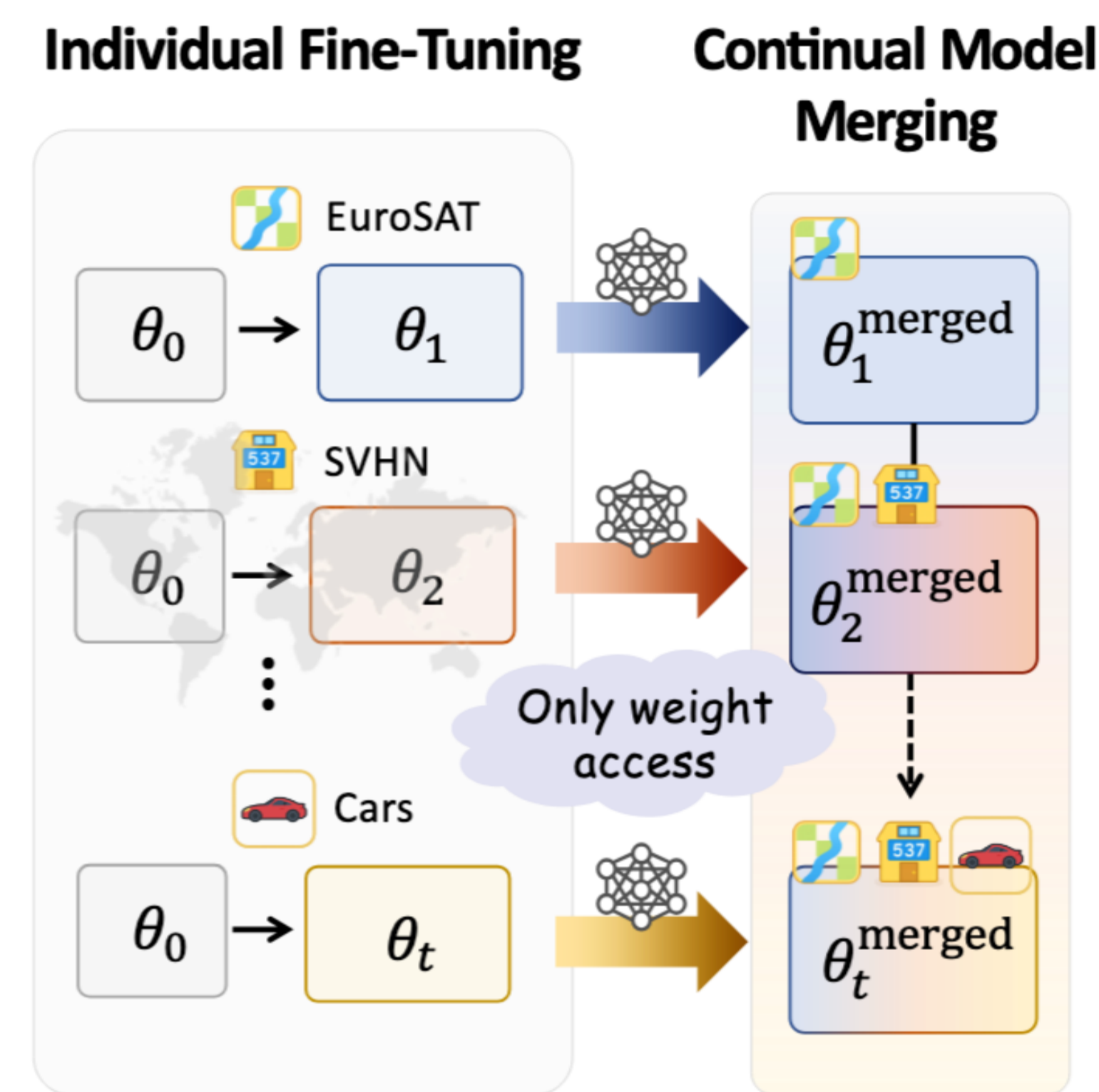
Code



Background and Motivation

Data-free continual model merging (DFCMM)

- A sequence of task-specific models are merged into a single backbone
- Tasks arrive sequentially and must be integrated without retraining
- Only the current model and merged model are accessible at each step



Desiderata
STABILITY: avoiding interference with earlier tasks;
PLASTICITY: adapting faithfully to each new task.

How to bridge data-level desiderata with parameter-space optimization?

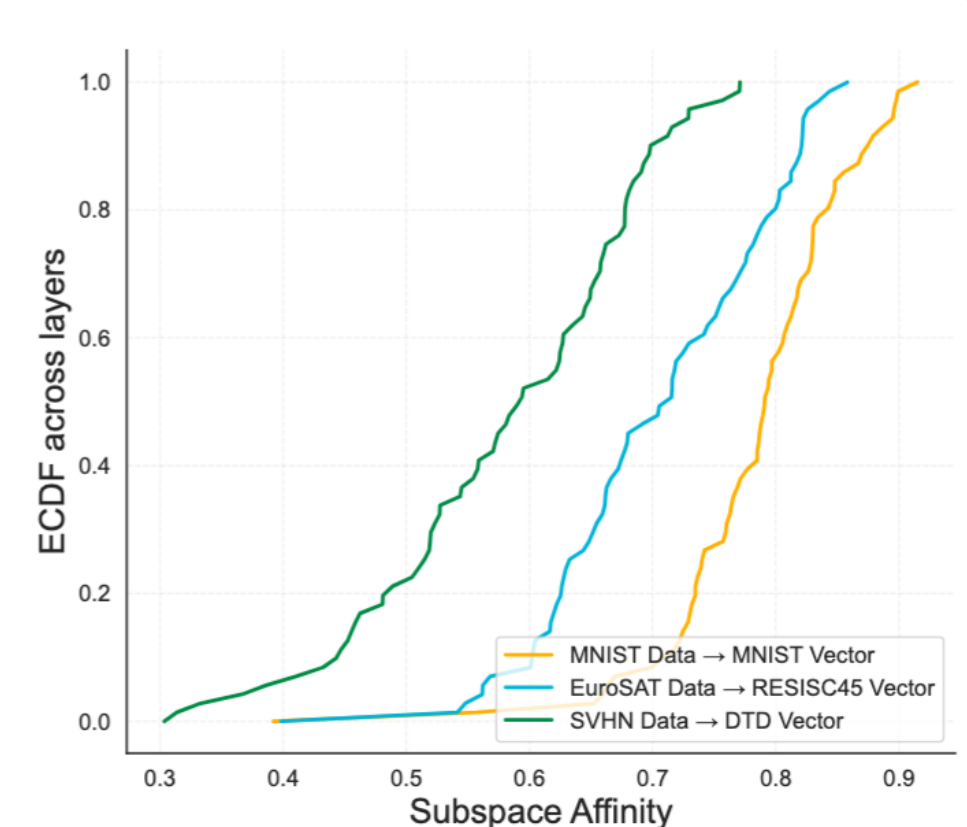
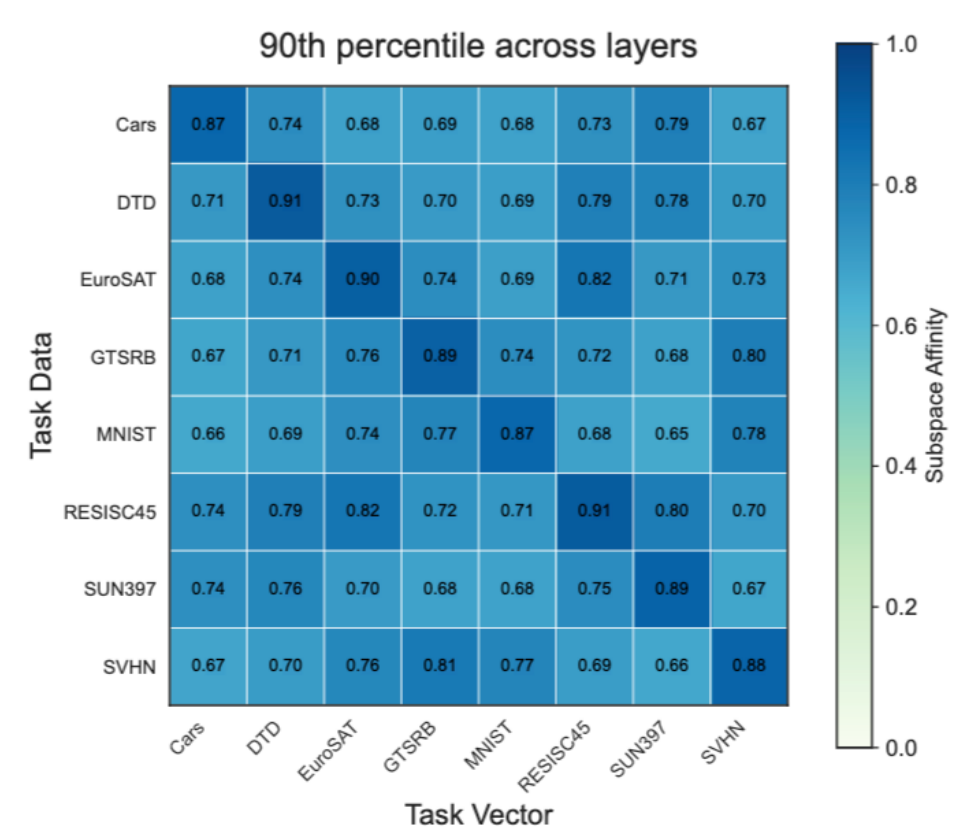
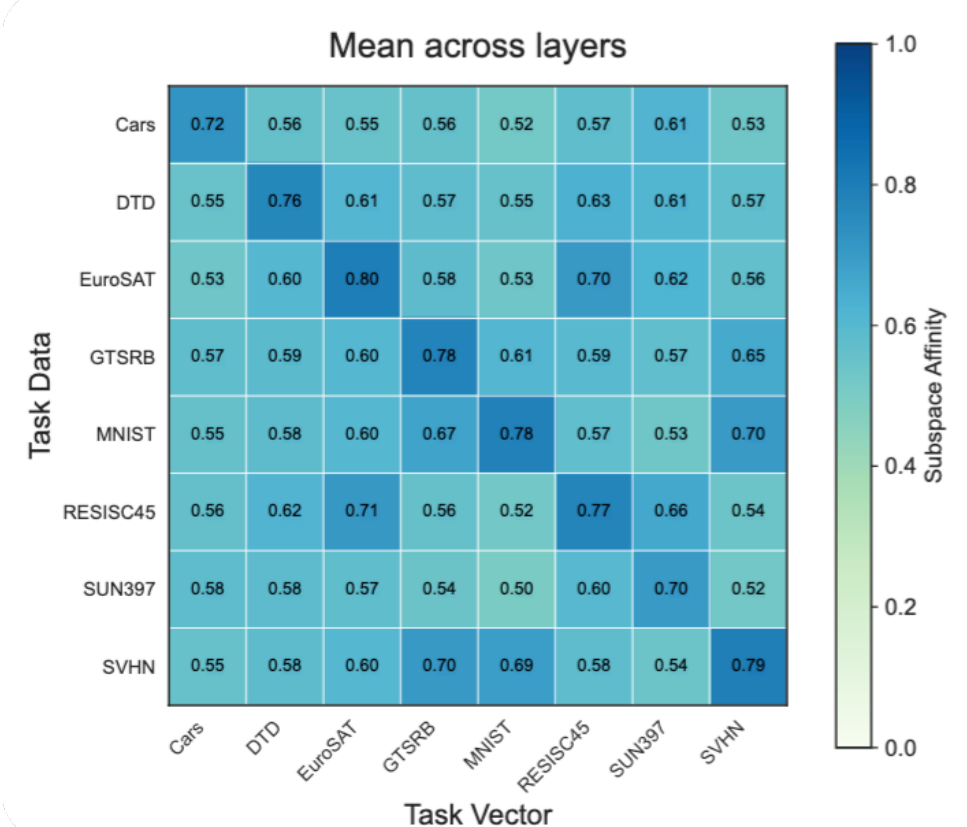


Key challenges

- Stability-plasticity trade-off: Preserve previous knowledge while adapting to new tasks
- Parameter interference: Task updates are entangled and may conflict in parameter space
- Fundamental gap: Data-level objectives vs parameter-space operations

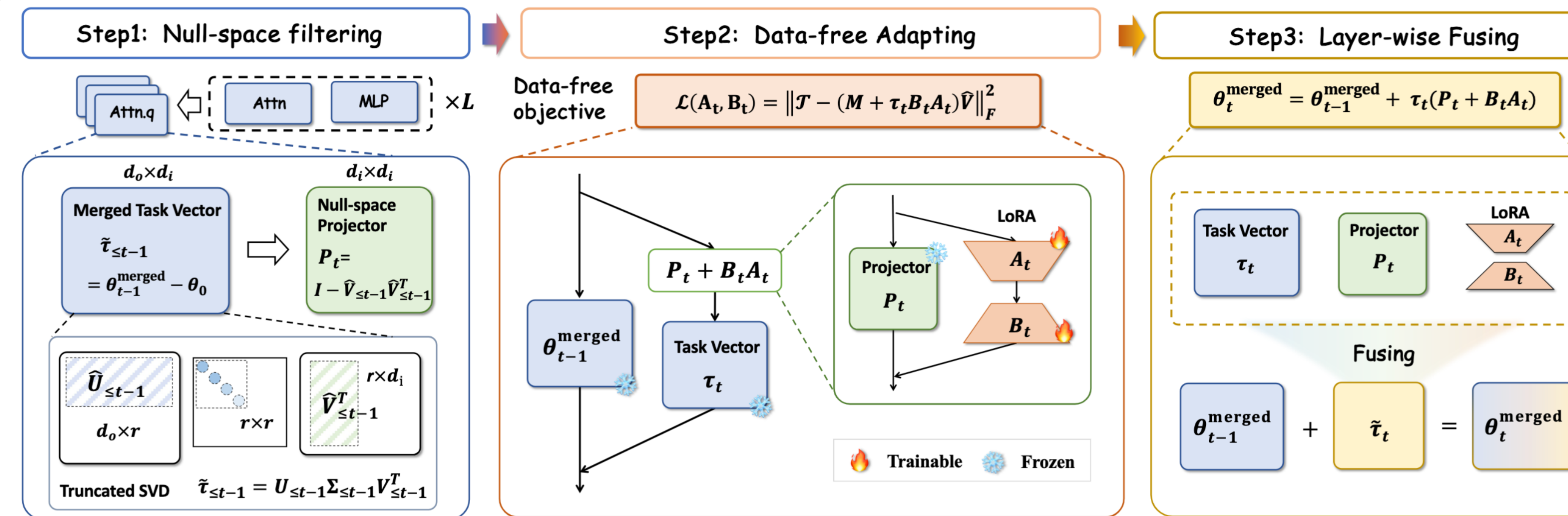
Key Insight: Task vectors align with representation subspaces

Implication: Parameter updates encode data-level structure



NUFILT: Null-Space Filtering

Data-free framework for continual model merging



Step 1: Filtering (Null-space projector preserves prior responses)

Project new task vector onto null-space of previous directions:

$$P_t^{(l)} = I - \hat{V}_{\leq t-1}^{(l)} \hat{V}_{\leq t-1}^{(l)\top} \quad \text{where } \hat{V}_{\leq t-1}^{(l)} \text{ is the top-}r_p \text{ right singular vectors of merged vector}$$

Property: for $x^{(l)} \in \text{span}(\hat{V}_{\leq t-1}^{(l)})$, $P_t^{(l)} x^{(l)} = 0$

Step 2: Adapting (LoRA adapter enables plasticity)

$$\mathcal{L}_{\text{stab}} = \|(\tilde{\tau}_{\leq t} - \tilde{\tau}_{\leq t-1}) X_o^\top\|_F^2 \leq 2\sigma_1(X_o)^2 (\|\tilde{\tau}_{\leq t} - \tilde{\tau}_{\leq t-1}\|_F^2 + r_o \zeta_o^2 \|\tilde{\tau}_{\leq t} - \tilde{\tau}_{\leq t-1}\|_2^2)$$

$$\mathcal{L}_{\text{plas}} = \|(\tilde{\tau}_{\leq t} - \tau_t) X_n^\top\|_F^2 \leq 2\sigma_1(X_n)^2 (\|\tilde{\tau}_{\leq t} - \tau_t\|_F^2 + r_n \zeta_n^2 \|\tilde{\tau}_{\leq t} - \tau_t\|_2^2)$$

Insert low-rank LoRA adapter:

$$P_t^{(l)} \rightarrow P_t^{(l)} + B_t^{(l)} A_t^{(l)}$$

Data-Free Objective:

$$\mathcal{L} = \|\mathcal{T} - (M + \tau_t^{(l)} B_t^{(l)} A_t^{(l)}) \hat{V}\|_F^2$$

\mathcal{T} : target behavior (preserve old tasks + match new task)

$\tau_t B_t A_t$: adaptive update (plasticity) \rightarrow learns new task

M : filtered update (stability) \rightarrow keeps previous knowledge

$$\mathcal{T} = \begin{bmatrix} \tilde{\tau}_{\leq t-1}^{(l)} \hat{V}_{\leq t-1}^{(l)} \\ \tau_t^{(l)} \hat{V}_t^{(l)} \end{bmatrix}, \quad \hat{V} = \begin{bmatrix} \hat{V}_{\leq t-1}^{(l)} \\ \hat{V}_t^{(l)} \end{bmatrix}, \quad M = \tilde{\tau}_{\leq t-1}^{(l)} + \tau_t^{(l)} P_t^{(l)}$$

Step 3: Fusing (Absorb extra params)

$$\theta_t^{\text{merged},(l)} = \theta_{t-1}^{\text{merged},(l)} + \tau_t^{(l)} (P_t^{(l)} + B_t^{(l)} A_t^{(l)})$$

Absorbed into weights, with no extra parameters or inference cost

Results

Effectiveness of each components

Method	Component	ACC(%) \uparrow			BWT(%) \uparrow		
		8 tasks	14 tasks	20 tasks	8 tasks	14 tasks	20 tasks
(1) Naive Merging ($\theta_t^{\text{merged}} = \theta_{t-1}^{\text{merged}} + \tau_t$)	\times / \times	62.1 \pm 0.0	46.5 \pm 0.0	34.3 \pm 0.0	-18.5 \pm 6.2	-25.8 \pm 2.2	-24.7 \pm 5.1
(2) Only Null-space ($\theta_t^{\text{merged}} = \theta_{t-1}^{\text{merged}} + \tau_t P_t$)	\checkmark / \times	80.0 \pm 1.1	76.7 \pm 1.0	67.0 \pm 0.7	-1.7 \pm 0.9	-4.2 \pm 0.9	-6.2 \pm 1.8
(3) Data-free Objective							
w/o null-space filter ($\theta_t^{\text{merged}} = \theta_{t-1}^{\text{merged}} + \tau_t B_t A_t$)	\times / \checkmark	75.8 \pm 1.9	63.7 \pm 1.6	51.7 \pm 1.0	-10.2 \pm 4.3	-17.1 \pm 2.7	-20.6 \pm 4.9
full method ($\theta_t^{\text{merged}} = \theta_{t-1}^{\text{merged}} + \tau_t (P_t + B_t A_t)$)	\checkmark / \checkmark	83.6 \pm 0.2	78.0 \pm 0.2	71.0 \pm 0.9	-2.7 \pm 0.7	-5.7 \pm 0.9	-8.9 \pm 2.3

State-of-the-art Results

Table 1: Comparative results of continual merging methods, reporting average accuracy and backward transfer over ten task orders (mean \pm std). EP and DA denote method assumptions: the need for extra parameters or data access. Best results are in **bold**, and the second best are underlined.

Method	Requirement	ViT-B/32			ViT-B/16			ViT-L/14		
		EP / DA	8 tasks	14 tasks	20 tasks	8 tasks	14 tasks	20 tasks	8 tasks	14 tasks
PRE-TRAINED	- / -	48.1	56.9	55.6	55.4	62.0	59.8	64.9	69.1	65.6
INDIVIDUAL	- / -	90.4	89.3	89.8	92.4	91.3	91.6	94.3	93.4	93.5
C. FINE-TUNED	- / -	79.8	67.4	62.6	82.9	72.2	68.2	90.0	70.9	77.7
WEIGHT AVERAGING	\times / \times	66.3 \pm 0.0	65.4 \pm 0.0	61.1 \pm 0.0	72.3 \pm 0.0	69.7 \pm 0.0	64.8 \pm 0.0	80.0 \pm 0.0	77.5 \pm 0.0	71.1 \pm 0.0
TASK ARITHMETIC	\times / \times	67.5 \pm 0.0	66.5 \pm 0.0	60.0 \pm 0.0	77.1 \pm 0.0	70.9 \pm 0.6	64.2 \pm 0.0	82.1 \pm 0.0	77.9 \pm 0.0	70.3 \pm 0.0
TIES-MERGING	\times / \times	49.0 \pm 0.2	66.2 \pm 0.6	59.9 \pm 0.7	66.8 \pm 3.7	70.5 \pm 0.8	63.0 \pm 1.6	64.3 \pm 7.0	78.0 \pm 0.6	68.3 \pm 0.9
MAGMAX-IND	\times / \times	70.7 \pm 0.0	67.0 \pm 0.0	61.2 \pm 0.0	76.7 \pm 1.8	67.0 \pm 0.0	62.5 \pm 0.0	83.4 \pm 0.0	71.2 \pm 0.0	71.2 \pm 0.0
LW ADAMERGING	\times / \checkmark	53.4 \pm 3.2	59.8 \pm 1.6	59.7 \pm 7.4	59.9 \pm 2.3	64.3 \pm 1.2	61.5 \pm 1.1	68.8 \pm 2.9	73.1 \pm 5.7	66.9 \pm 1.1
LoRA-WEMOE	\checkmark / \checkmark	68.8 \pm 7.8	63.8 \pm 3.4	49.6 \pm 15.4	72.6 \pm 3.7	67.9 \pm 2.9	55.0 \pm 7.0	75.6 \pm 7.8	74.0 \pm 5.0	56.9 \pm 19.8
WUDI-MERGING	\times / \times	74.7 \pm 6.6	67.0 \pm 6.9	63.7 \pm 3.8	81.0 \pm 4.7	75.0 \pm 4.1	69.6 \pm 4.7	87.5 \pm 3.3	84.2 \pm 3.7	78.1 \pm 2.8
ISO-C	\times / \times	71.7 \pm 1.2	73.2 \pm 1.8	67.6 \pm 0.8	78.5 \pm 1.2	79.7 \pm 1.3	73.0 \pm 1.1	86.9 \pm 0.5	86.9 \pm 1.8	80.9 \pm 0.8
KNOTS-TIES	\times / \times	54.4 \pm 6.9	67.8 \pm 0.4	60.5 \pm 1.4	57.9 \pm 8.4	71.6 \pm 0.3	62.7 \pm 1.1	68.3 \pm 5.7	78.8 \pm 0.3	69.7 \pm 0.8
TSPV-M	\times / \times	68.2 \pm 4.8	63.3 \pm 4.8	58.8 \pm 3.3	75.4 \pm 4.0	69.2 \pm 3.1	63.1 \pm 1.3	82.2 \pm 3.6	78.1 \pm 3.6	70.5 \pm 1.2
OPCM	\times / \times	75.5 \pm 0.5	71.9 \pm 0.3	65.7 \pm 0.2	81.8 \pm 0.3	77.1 \pm 0.5	70.3 \pm 0.2	87.0 \pm 0.4	83.5 \pm 0.2	76.0 \pm 0.2
NUFILT (Ours)	\times / \times	83.6\pm0.2	78.0\pm0.2	71.0\pm0.9	87.3\pm0.1	83.1\pm0.3	78.1\pm0.9	91.6\pm0.1	89.2\pm0.1	84.7\pm0.8
WEIGHT AVERAGING	\times / \times	-11.5 \pm 2.2	-8.0 \pm 1.3	-7.1 \pm 2.1	-9.7 \pm 1.5	-7.1 \pm 1.4	-7.3 \pm 1.7	-7.3 \pm 1.4	-5.8 \pm 1.0	-6.4 \pm 1.5
TASK ARITHMETIC	\times / \times	-9.6 \pm 1.5	-1.3 \pm 1.6	-3.4 \pm 1.0	-4.2 \pm 1.0	-1.3 \pm 0.4	-3.6 \pm 0.4	-7.1 \pm 0.8	-1.8 \pm 0.3	-3.3 \pm 0.3
TIES-MERGING	\times / \times	-15.3 \pm 8.0	1.9\pm0.6	-1.5 \pm 0.7	-5.5 \pm 0.4	1.4\pm0.7	-1.5 \pm 1.2	-13.0 \pm 5.7	-1.1 \pm 0.4	-2.9 \pm 1.0
MAGMAX-IND	\times / \times	-8.3 \pm 1.3	-7.4 \pm 1.4	-7.2 \pm 1.6	-6.1 \pm 1.3	-7.4 \pm 2.0	-8.0 \pm 2.2	-5.0 \pm 0.8	-6.0 \pm 2.1	-6.5 \pm 2.1
LW ADAMERGING	\times / \checkmark	-32.5 \pm 3.6	-24.1 \pm 1.7	-22.7 \pm 4.3	-27.8 \pm 2.7	-22.1 \pm 1.4	-21.4 \pm 1.2	-24.3 \pm 3.3	-19.6 \pm 1.7	-21.7 \pm 1.1
LoRA-WEMOE	\checkmark / \checkmark	-20.4 \pm 9.0	-20.2 \pm 3.9	-24.5 \pm 10.0	-18.0 \pm 6.2	-18.8 \pm 3.4	-25.8 \pm 7.9	-17.8 \pm 5.9	-16.8 \pm 5.3	-27.9 \pm 17.2
WUDI-MERGING	\times / \times	-17.0 \pm 7.5	-22.8 \pm 7.3	-26.0 \pm 4.1	-12.6 \pm 5.4	-16.9 \pm 4.4	-18.5 \pm 14.2	-7.3 \pm 3.7	-9.4 \pm 4.0	-15.8 \pm 2.9
ISO-C	\times / \times	-10.2 \pm 1.2	-10.4 \pm 1.9	-10.3 \pm 1.4	-6.7 \pm 0.5	-7.1 \pm 1.1	-9.9 \pm 1.7	-3.7 \pm 0.6	-3.7 \pm 1.7	-5.5 \pm 1.3
KNOTS-TIES	\times / \times	-12.6 \pm 3.9	1.9\pm0.5	-1.3 \pm 0.7	-13.5 \pm 5.5	1.0\pm0.1	-2.3 \pm 0.6	-11.6 \pm 3.6	0.3\pm0.3	-2.3 \pm 0.7
TSPV-M	\times / \times	-24.0 \pm 5.6	-26.7 \pm 4.9	-31.7 \pm 3.4	-18.5 \pm 4.6	-22.7 \pm 3.1	-28.2 \pm 1.8	-13.0 \pm 4.0	-15.6 \pm 3.9	-23.3 \pm 1.3
OPCM	\times / \times	-6.3 \pm 1.1	-6.0 \pm 1.0	-7.8 \pm 1.5	-4.8 \pm 0.7	-5.1 \pm 1.4	-6.3 \pm 2.2	-2.6 \pm 1.0	-4.3 \pm 0.7	-6.5 \pm 1.8
NUFILT (Ours)	\times / \times	-2.7\pm0.7	-5.7\pm0.9	-8.9\pm2.3	-1.6\pm0.5	-3.5\pm0.6	-7.1\pm1.9	-1.1\pm0.3	-2.0\pm0.3	-4.6\pm0.7

Conclusion:

- Established subspace alignment property of task vectors
- NUFILT = null-space filtering + projection-aware adaptation
- SOTA performance, minimal forgetting, no extra data/params