

NegMerge: Consensual Weight Negation for Strong Machine Unlearning

Hyoseo Kim^{1,2*}, Dongyoon Han^{1†}, Junsuk Choe^{2†}

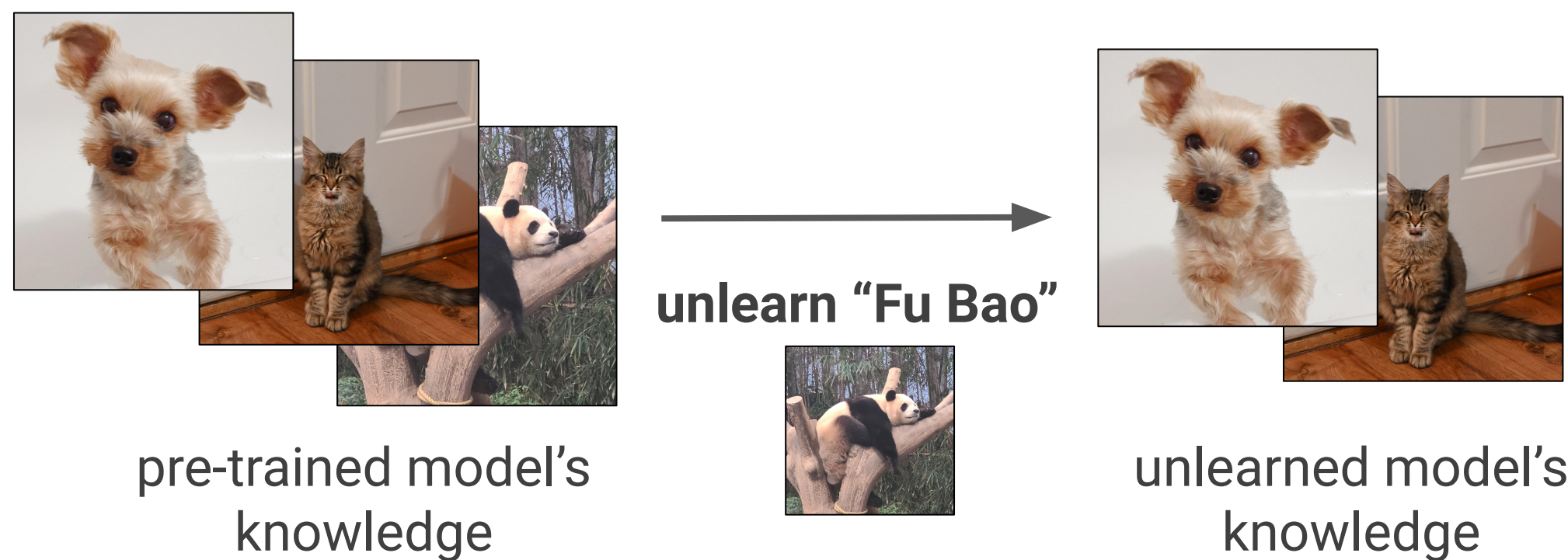
¹NAVER AI LAB, ²Sogang University

*Work done during an internship at NAVER AI Lab, †Corresponding author



What is Machine Unlearning?

- Removal of specific knowledge from a pre-trained model without impacting its remaining knowledge.
- Applications: protecting privacy and rights under regulations like GDPR, correcting inaccuracies, and removing harmful data.

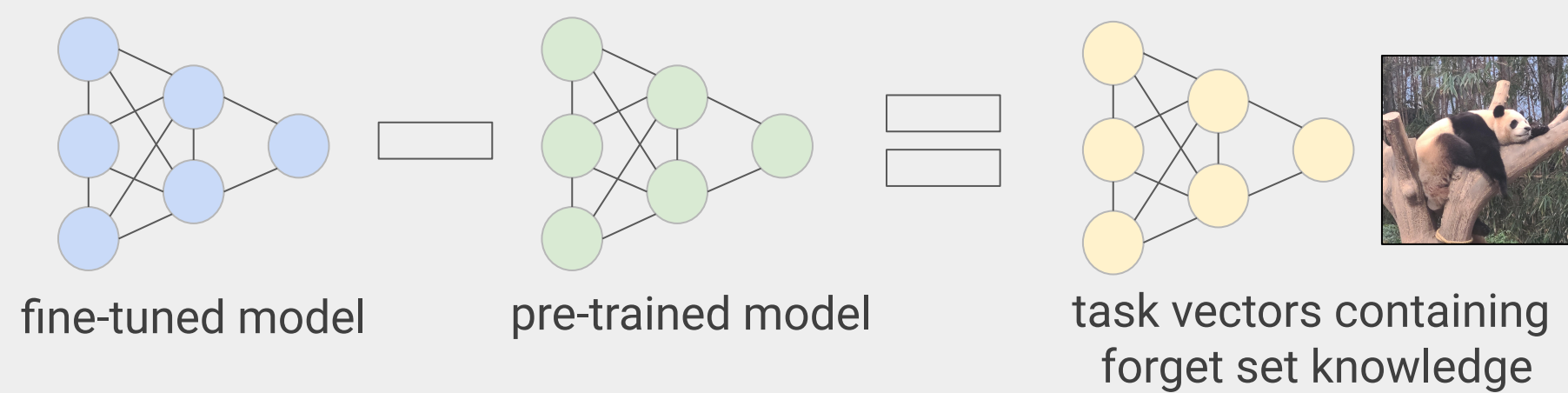


How to Unlearn?

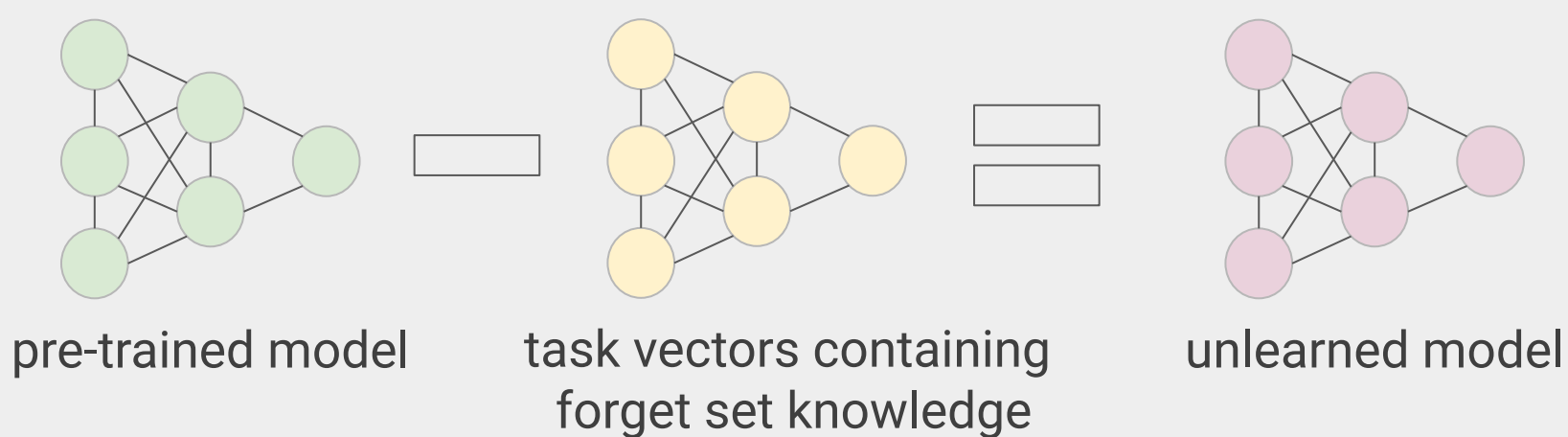
- Forgetting via negation^[1]
 - Adjust the model by subtracting the sum of the task vectors.

$$\theta_{unlearn} = \theta_{pre} - \lambda \underbrace{(\theta_{ft}^{forget} - \theta_{pre})}_{\text{sum of the task vectors}}$$

(a) define task vector



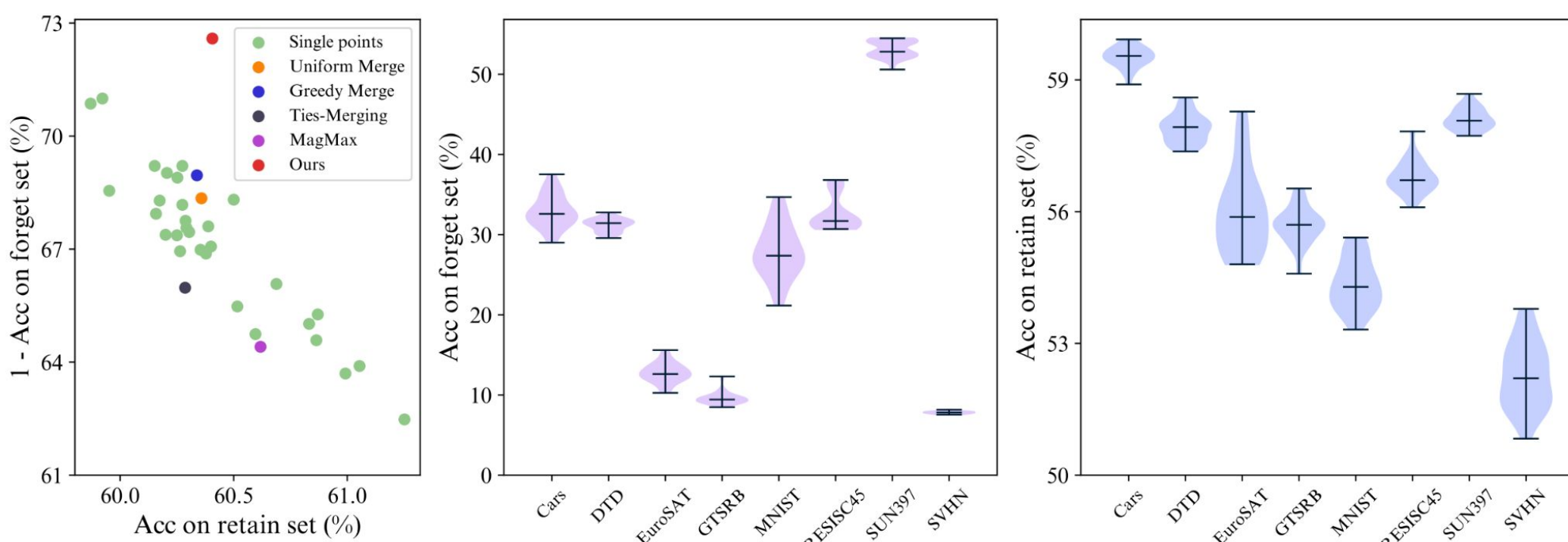
(a) subtract task vector



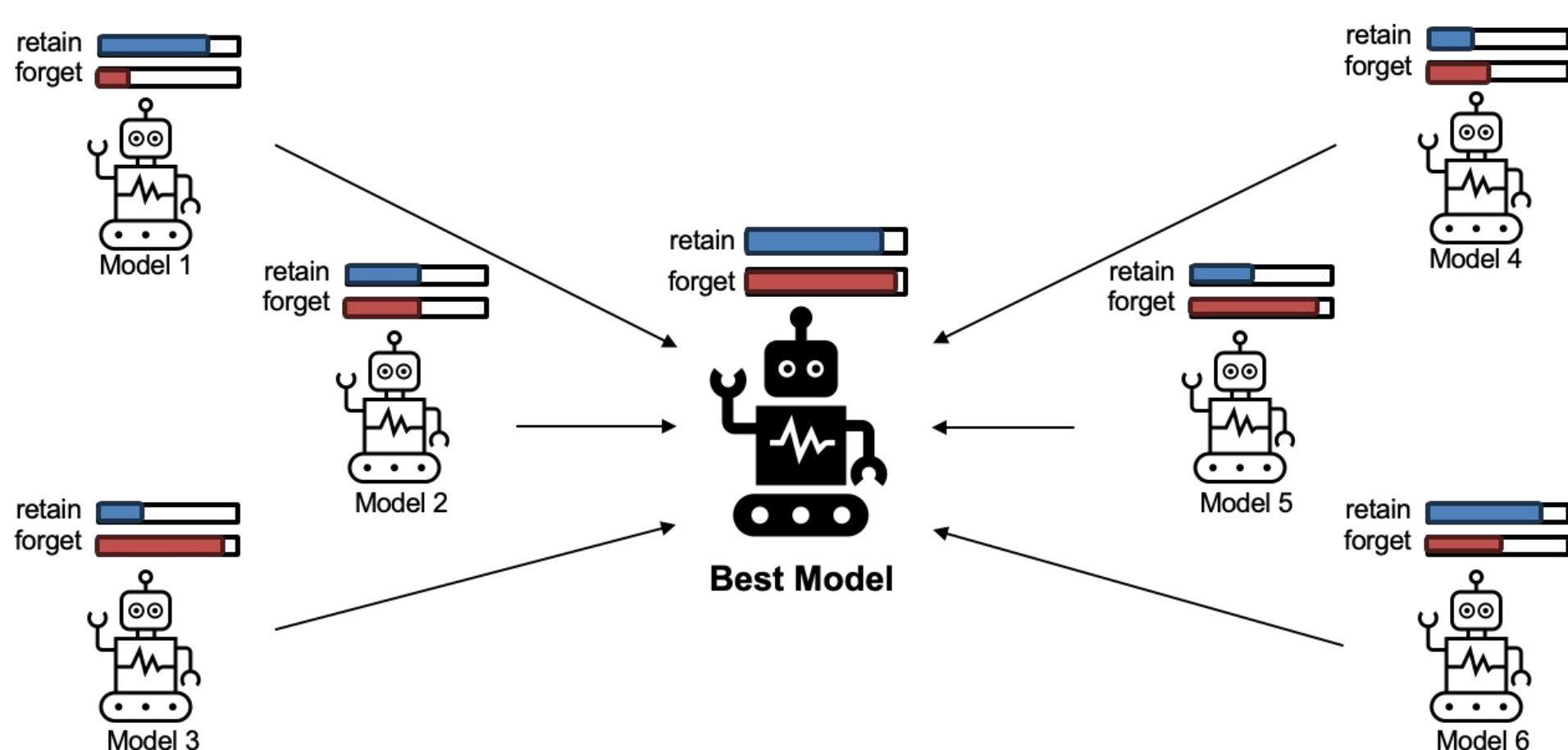
[1] Ilharco, Gabriel, et al. "Editing models with task arithmetic." The Eleventh International Conference on Learning Representations.

Challenges of Machine Unlearning

- Highly sensitive to the hyperparameters used for fine-tuning.
- Trade-off in unlearning performance and retaining performance.



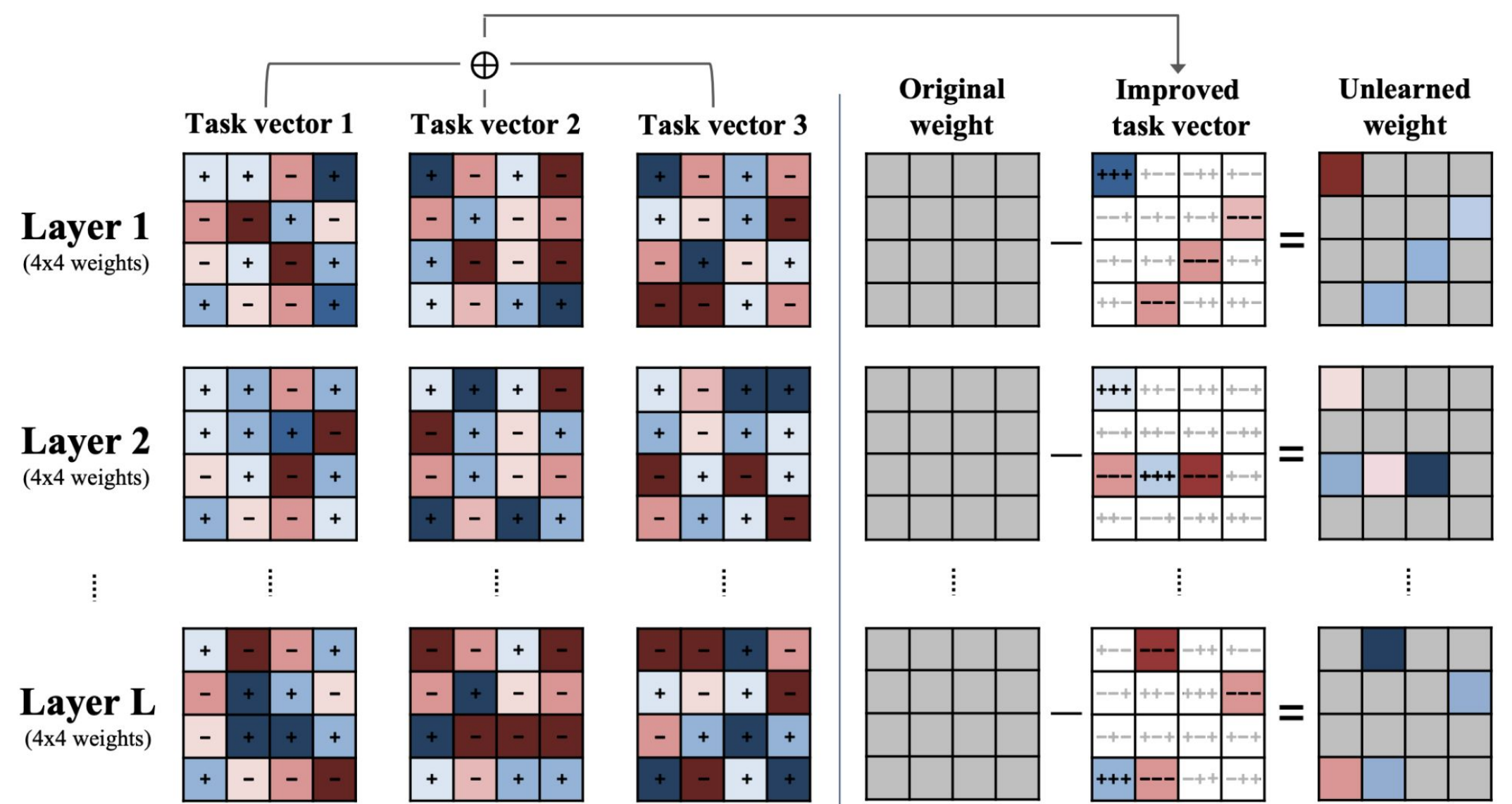
- Effective unlearning requires the fine-tuned model to maintain high performance on the forget set without degrading performance on the retain set.
- To achieve this, merge them all instead of selecting just one.



Overview of NegMerge

- Hyperparameter tuning generates multiple fine-tuned models.
- Merge all models based on the sign consensus of task vectors.

$$\tau_{merged} = \left(\frac{1}{n} \sum_{k=1}^n \tau_k \right) \odot \mathbf{1}_{\text{signs are equal}}$$



- Elements that consistently show the same sign across task vectors are attributed to the forget set.
- Components that exhibit differing signs are considered less related to the forget set.

Experiment Results

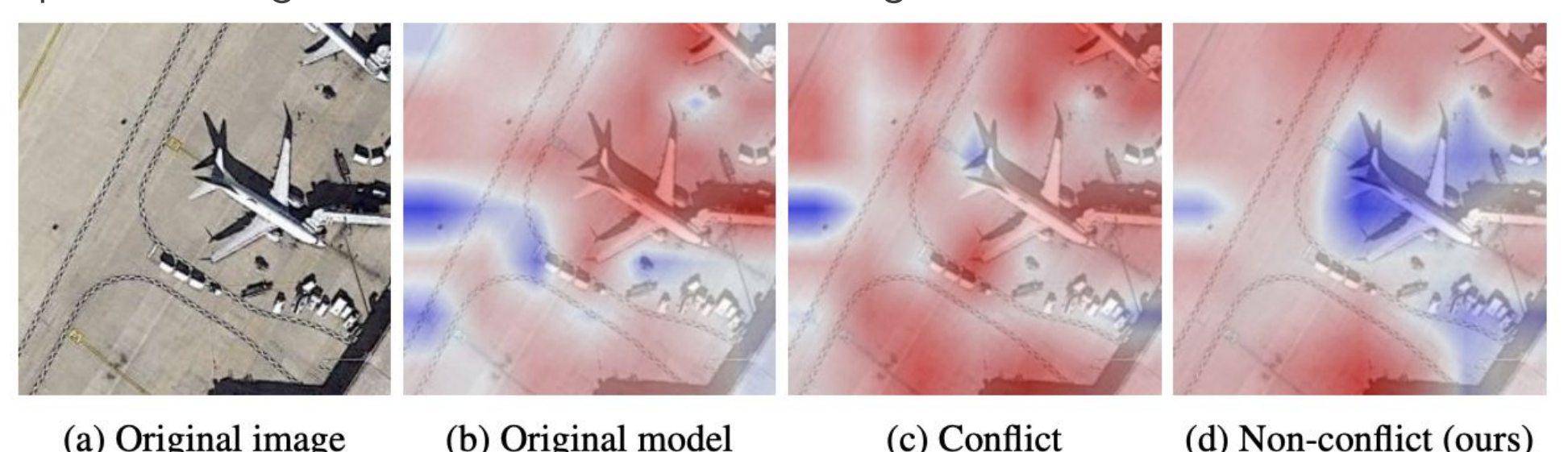
- Unlearning on CLIP ViT Models.
 - achieves the best reduction in accuracy on the forget set.

Method	ViT-B/32		ViT-B/16		ViT-L/14		Time (sec)
	Acc $D_f(\downarrow)$	Acc $D_r(\uparrow)$	Acc $D_f(\downarrow)$	Acc $D_r(\uparrow)$	Acc $D_f(\downarrow)$	Acc $D_r(\uparrow)$	
Pre-trained	48.13	63.33	55.49	68.32	65.19	75.54	-
Task Arithmetic							
<i>Paper number*</i>	24.00	60.90	21.30	65.40	19.00	72.90	-
Single Best Model [†]	23.63	60.60	20.64	64.04	19.17	72.09	-
Uniform Merge	22.50	60.55	21.51	64.60	18.10	71.91	12±0.1
Greedy Merge [‡]	23.31	60.75	21.34	64.54	17.71	71.99	607±2.6
TIES-Merging	26.21	61.08	23.78	64.72	22.70	72.41	128±10.1
MagMax	25.24	60.95	24.45	64.78	21.71	72.55	24±1.8
NegMerge (Ours)	20.76	60.36	19.24	64.54	17.32	72.08	37±1.2
Linear Task Arithmetic							
<i>Paper number*</i>	10.90	60.80	11.30	64.80	-	-	-
Single Best Model [†]	8.88	60.16	6.92	64.62	-	-	-
Uniform Merge	9.12	60.47	6.84	65.26	-	-	19±2.3
Greedy Merge [‡]	8.73	60.27	6.80	64.72	-	-	1696±35.3
TIES-Merging	10.66	60.38	8.44	65.12	-	-	378±8.0
MagMax	11.33	60.67	8.65	65.17	-	-	164±2.4
NegMerge (Ours)	8.03	60.58	6.60	65.40	-	-	194±1.6

- 10% Random Data Forgetting on CIFAR-10 using ResNet-18.
 - achieves the smallest average gap of 1.07.

Methods	Used Splits	Acc $D_r(\simeq)$	Acc $D_f(\simeq)$	Acc $D_{test}(\simeq)$	MIA(\simeq)	Avg. Gap(\downarrow)
Retrain *	retain	100.00±0.00	94.76±0.69	94.26±0.02	12.88±0.09	0.00
Random Labeling *	all	99.67±0.14	92.39±0.31	92.83±0.38	37.36±0.06	7.15
Influence *		99.20±0.22	98.93±0.28	93.20±1.03	2.67±0.01	4.06
SalUn *		99.62±0.12	97.15±0.43	93.93±0.29	14.39±0.82	1.15
Finetune *	retain	99.88±0.08	99.37±0.55	94.06±0.27	2.70±0.01	3.78
ℓ_1 -sparse *		97.74±0.33	95.81±0.62	91.59±0.57	9.84±0.00	2.26
Gradient Ascent *		99.50±0.38	99.31±0.54	94.01±0.47	1.70±0.01	4.12
Boundary Shrink *		98.29±2.50	98.22±2.52	92.69±2.99	8.96±0.13	2.67
Boundary Expanding *	forget	99.42±0.33	99.41±0.30	93.85±1.02	7.47±1.15	2.76
Random Labeling		99.99±0.00	99.98±0.02	95.04±0.11	2.15±1.94	4.19
SalUn		99.88±0.04	99.89±0.04	94.42±0.05	9.51±2.07	2.20
Task Arithmetic						
Single Best Model [†]		98.36±0.51	94.85±0.16	91.49±0.80	10.91±0.72	1.62
Uniform Merge	forget	98.70±0.91	95.83±2.17	92.36±1.16	10.14±2.93	1.75
TIES-Merging		98.38±0.17	95.45±0.32	92.23±0.14	9.36±0.31	1.96
MagMax		98.38±0.12	97.97±0.77	91.53±0.00	8.45±2.60	3.00
NegMerge (Ours)		99.15±0.24	96.63±0.59	92.71±0.39	12.87±1.29	1.07

- Impact of Sign Conflicts on Unlearning: Grad-CAM



Conclusion

- Propose a novel machine unlearning technique, NegMerge, based on task arithmetic and model merging.
- Tested on the CLIP ViT models and the standard ResNet18 classifier, achieving SOTA across nine datasets.