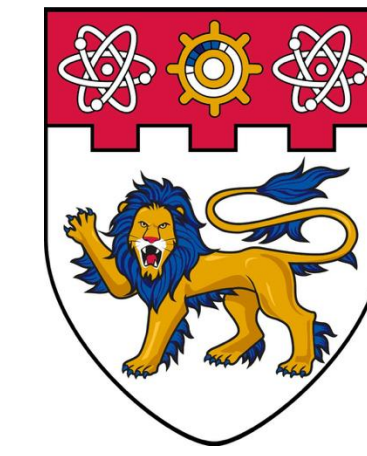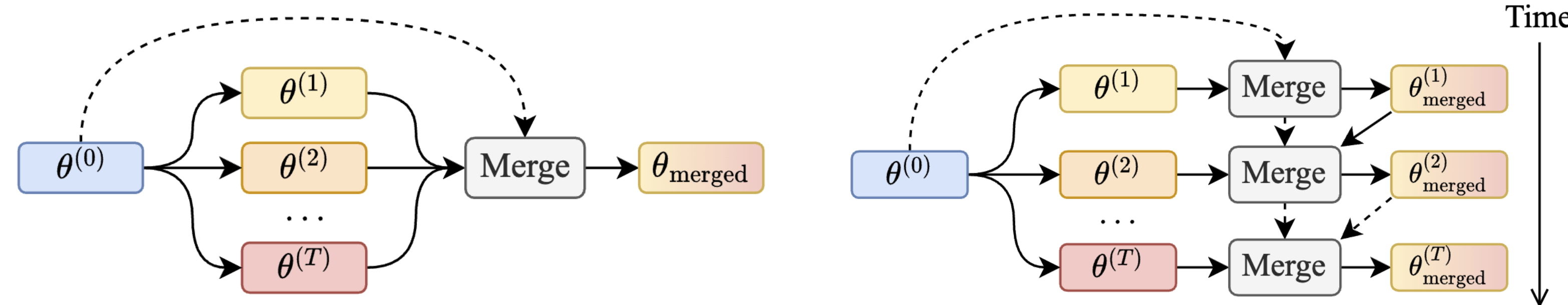# Merging on the Fly Without Retraining:
## A Sequential Approach to Scalable Continual Model Merging

Anke Tang, Enneng Yang, Li Shen, Yong Luo, Han Hu, Lefei Zhang, Bo Du, Dacheng Tao

## Two Types of Model Merging Paradigms



Conventional model merging
models are available **simultaneously**

Continual model merging
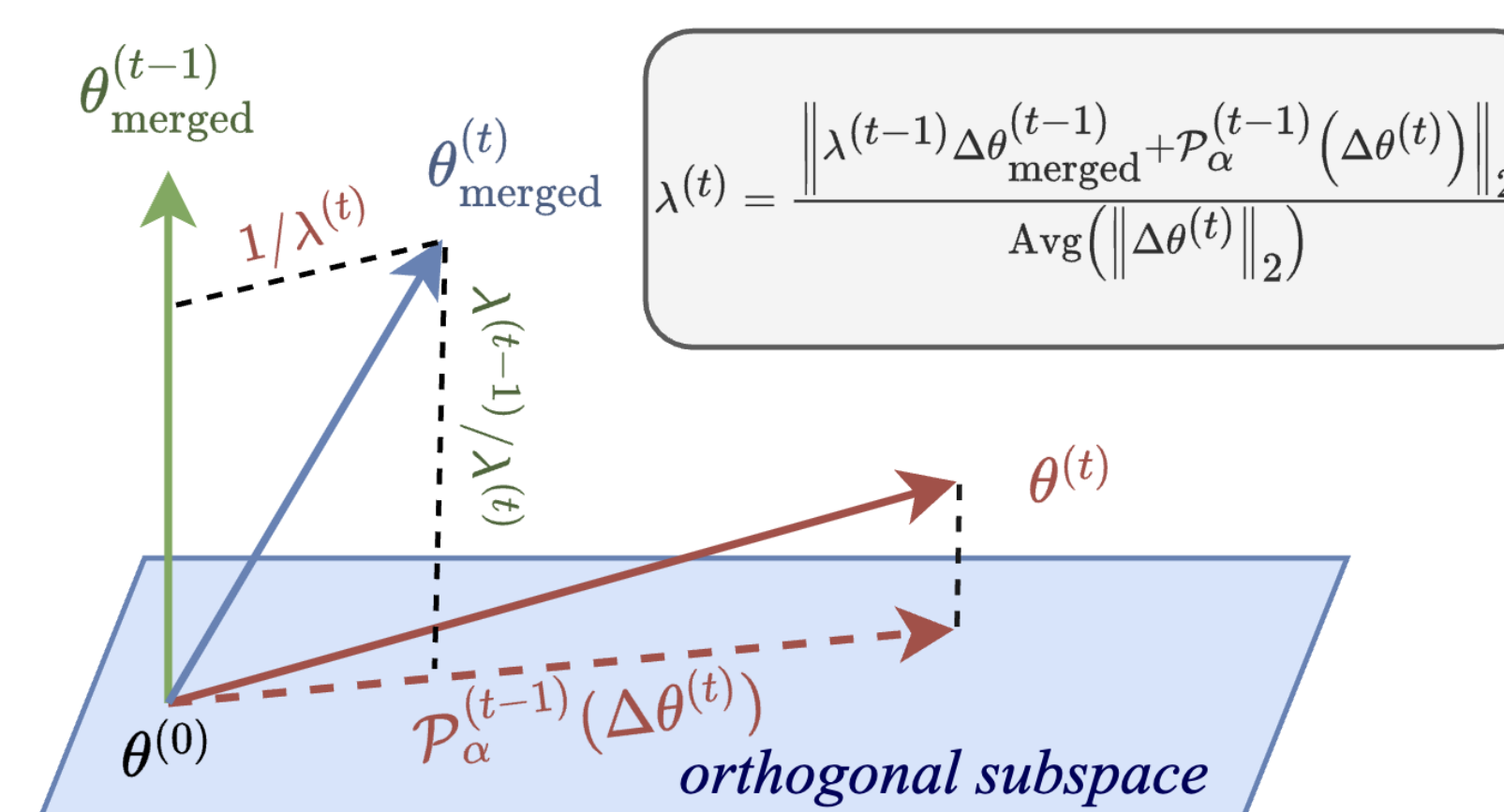models are available **sequentially**

## Why Continual Model Merging?

➢ **Memory Efficiency**: Conventional merging needs loading all expert models simultaneously, which scales linearly with the number of models. Continual merging keeps constant memory by only storing the current merged model + the new model.

➢ **Cheap**: The proposed continual model merging is training-free, so it avoids extra optimization and data access.

➢ **Better scalability**: Continual merging supports large task collections and frequent updates, unlike one-shot merging that becomes harder as the number of models grows.

## Abstract of this paper

Deep model merging represents an emerging research direction that combines multiple fine-tuned models to harness their specialized capabilities across different tasks and domains. Current model merging techniques focus on merging all available models simultaneously, with weight interpolation-based methods being the predominant approach. However, these conventional approaches are not well-suited for scenarios where models become available sequentially, and they often suffer from high memory requirements and potential interference between tasks. In this study, we propose a training-free projection-based continual merging method that processes models sequentially through orthogonal projections of weight matrices and adaptive scaling mechanisms. Our method operates by projecting new parameter updates onto subspaces orthogonal to existing merged parameter updates while using an adaptive scaling mechanism to maintain stable parameter distances, enabling efficient sequential integration of task-specific knowledge. Our approach maintains constant memory complexity to the number of models, minimizes interference between tasks through orthogonal projections, and retains the performance of previously merged models through adaptive task vector scaling. Extensive experiments on CLIP-ViT models demonstrate that our method achieves a 5-8% average accuracy improvement while maintaining robust performance in different task orderings.

## The Proposed Method

The **O**rthogonal **P**rojection-based **C**ontinual **M**erging (OPCM)



$$\lambda^{(t)} = \frac{\left\| \lambda^{(t-1)} \Delta\theta_{\text{merged}}^{(t-1)} + \mathcal{P}_\alpha^{(t-1)}\left(\Delta\theta^{(t)}\right) \right\|_2}{\text{Avg}\left(\left\|\Delta\theta^{(t)}\right\|_2\right)}$$

Mark this paper

Project page

FusionBench
(Benchmark on model merging)

---

**Algorithm 1** OPCM

1: Initialize $\theta_{\text{merged}}^{(1)} = \theta^{(1)}$, average norm of the task vectors $n = \|\Delta\theta^{(1)}\|_2$, scaling factor $\lambda^{(1)} = 1$.
2: **for** $t = 2$ to $T$ **do**
3:     **for** weight matrices $W \in \mathbb{R}^{m \times n}$ in linear layers **do**
4:         $\Delta W_{\text{merged}}^{(t-1)} \leftarrow W_{\text{merged}}^{(t-1)} - W^{(0)}$
5:         $\Delta W^{(t)} \leftarrow W^{(t)} - W^{(0)}$
6:         $\Delta W_{\text{proj}}^{(t)} \leftarrow \mathcal{P}_\alpha^{(t-1)}\left(\Delta W^{(t)}; \Delta W_{\text{merged}}^{(t-1)}\right)$
7:     **end for**
8:     **for** other parameters $p$ **do**
9:         $\Delta p^{(t)} \leftarrow p^{(t)} - p^{(0)}$
10:     **end for**
11:     $\Delta\theta_{\text{merged}}^{(t-1)} \leftarrow \theta_{\text{merged}}^{(t-1)} - \theta^{(0)}$
12:     Concatenate $\Delta W_{\text{proj}}^{(t)}$ and $\Delta p^{(t)}$ into $\Delta\theta_{\text{proj}}^{(t)}$
13:     $n \leftarrow \frac{t-1}{t} n + \frac{1}{t}\|\Delta\theta^{(t)}\|_2$
14:     $\lambda^{(t)} \leftarrow \left\| \lambda^{(t-1)} \Delta\theta_{\text{merged}}^{(t-1)} + \Delta\theta_{\text{proj}}^{(t)} \right\|_2 / n$
15:     $\theta_{\text{merged}}^{(t)} \leftarrow \theta^{(0)} + \frac{\lambda^{(t-1)} \Delta\theta_{\text{merged}}^{(t-1)} + \Delta\theta_{\text{proj}}^{(t)}}{\lambda^{(t)}}$
16: **end for**
17: **return** $\theta_{\text{merged}}^{(T)}$