

DualOptim+: Bridging Shared and Decoupled Optimizer States for Better Machine Unlearning in Large Language Models

 Xuyang Zhong¹ Qizhang Li² Yiwen Guo² Chen Liu¹
¹City University of Hong Kong ²Independent Researcher

xuyang.zhong@my.cityu.edu.hk {liqizhang95, guoyiwen89}@gmail.com chen.liu@cityu.edu.hk

Comparison of Baselines and DualOptim+

Machine unlearning (MU) solves the optimization problem:

$$\min_{\theta} \mathcal{L}_f(\theta, \mathcal{D}_f) + \mathcal{L}_r(\theta, \mathcal{D}_r) \quad (1)$$

The update rules of baselines and DualOptim+ are presented as follows:

Joint:

$$\theta \leftarrow \theta - \mathcal{P}(\nabla_{\theta}(\mathcal{L}_f + \mathcal{L}_r)). \quad (2)$$

Alternate:

$$\begin{cases} \theta \leftarrow \theta - \mathcal{P}(\nabla_{\theta}\mathcal{L}_f) & \text{for forget data} \\ \theta \leftarrow \theta - \mathcal{P}(\nabla_{\theta}\mathcal{L}_r) & \text{for retain data} \end{cases} \quad (3)$$

DualOptim:

$$\begin{cases} \theta \leftarrow \theta - \mathcal{P}_f(\nabla_{\theta}\mathcal{L}_f) & \text{for forget data} \\ \theta \leftarrow \theta - \mathcal{P}_r(\nabla_{\theta}\mathcal{L}_r) & \text{for retain data} \end{cases} \quad (4)$$

DualOptim+ (Ours):

$$\begin{cases} B \leftarrow \beta B + (1 - \beta)\nabla_{\theta}\mathcal{L}_f & \text{for forget data} \\ B \leftarrow \beta B + (1 - \beta)\nabla_{\theta}\mathcal{L}_r & \text{for retain data} \end{cases} \quad (5)$$

$$\begin{cases} \Delta_f \leftarrow \beta\Delta_f + (1 - \beta)(\nabla_{\theta}\mathcal{L}_f - \hat{B}) & \text{for forget data} \\ \Delta_r \leftarrow \beta\Delta_r + (1 - \beta)(\nabla_{\theta}\mathcal{L}_r - \hat{B}) & \text{for retain data} \end{cases} \quad (6)$$

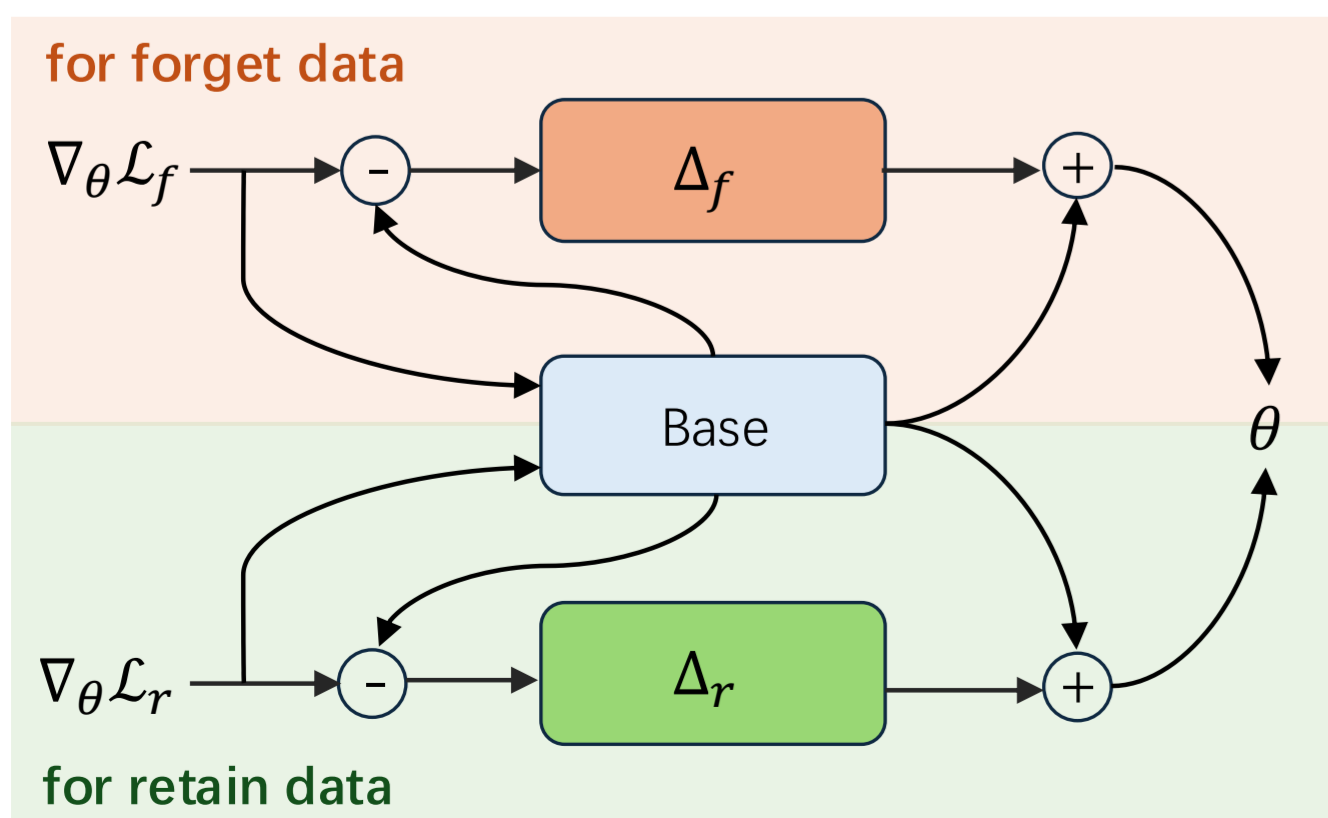
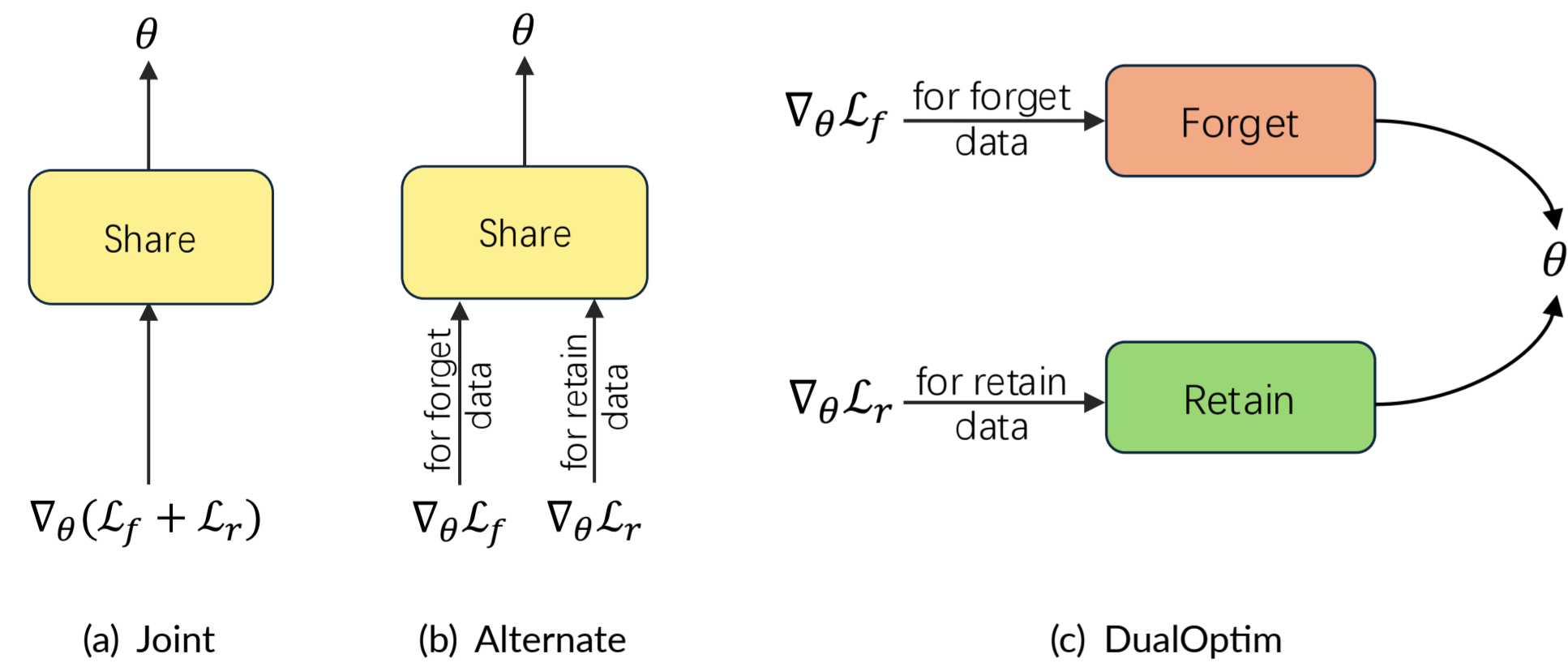


Figure 1. Comparison of baselines and DualOptim+.

Pseudo Code

 The pseudo code of **DualOptim+** with Adam is shown below.

Algorithm 1 DualOptim+ with Adam

```

1: Input: parameter  $\theta$ , learning rate  $\eta$ , betas  $(\beta_1, \beta_2)$ , epsilon  $\epsilon$ , weight decay factor  $\lambda$ , forget objective  $\mathcal{L}_f$ , retain objective  $\mathcal{L}_r$ , total steps  $N$ , forget frequency  $F_f$ , retain frequency  $F_r$ 
2: Initialize:  $m_{\Delta_f} \leftarrow 0, m_{\Delta_r} \leftarrow 0, m_B \leftarrow 0, v_{\Delta_f} \leftarrow 0, v_{\Delta_r} \leftarrow 0, v_B \leftarrow 0, t_f \leftarrow 0, t_r \leftarrow 0$ 
3:  $\hat{m}_B, \hat{v}_B \leftarrow m_B, v_B$ 
4: for  $t = 1$  to  $N$  do
5:   if  $t \bmod (F_f + F_r) \leq F_f$  then
6:      $t_f \leftarrow t_f + 1$ 
7:      $g, m_{\Delta}, v_{\Delta}, t' \leftarrow \nabla_{\theta}\mathcal{L}_f(\theta), m_{\Delta_f}, v_{\Delta_f}, t_f$ 
8:   else
9:      $t_r \leftarrow t_r + 1$ 
10:     $g, m_{\Delta}, v_{\Delta}, t' \leftarrow \nabla_{\theta}\mathcal{L}_r(\theta), m_{\Delta_r}, v_{\Delta_r}, t_r$ 
11:   end if
12:    $\theta \leftarrow \theta - \eta\lambda\theta$ 
13:    $m_{\Delta} \leftarrow \beta_1 m_{\Delta} + (1 - \beta_1)(g - \hat{m}_B)$ 
14:    $v_{\Delta} \leftarrow \beta_2 v_{\Delta} + (1 - \beta_2)(g^2 - \hat{v}_B)$ 
15:    $\hat{m}_{\Delta}, \hat{v}_{\Delta} \leftarrow m_{\Delta}/(1 - \beta_1^{t'}), v_{\Delta}/(1 - \beta_2^{t'})$ 
16:    $\theta \leftarrow \theta - \eta(\hat{m}_B + \hat{m}_{\Delta})/(\sqrt{\hat{v}_B + \hat{v}_{\Delta}} + \epsilon)$ 
17:    $m_B \leftarrow \beta_1 m_B + (1 - \beta_1)g$ 
18:    $v_B \leftarrow \beta_2 v_B + (1 - \beta_2)g^2$ 
19:    $\hat{m}_B, \hat{v}_B \leftarrow m_B/(1 - \beta_1^{t'}), v_B/(1 - \beta_2^{t'})$ 
20: end for
21: Output: parameter  $\theta$ 

```

Theoretical Analysis on DualOptim+

Theorem 1 (Convergence of Base and Delta States). Assume the expectations of gradients $g_{f,t}$ and $g_{r,t}$ over time $\mathbb{E}_t[g_{f,t}] = mG$, $\mathbb{E}_t[g_{r,t}] = nG$ exist, where $m, n \in [-1, 1]$, and G is a non-negative constant. We have:

$$\begin{aligned} \lim_{T \rightarrow \infty} B_{(F_f+F_r)T} &= \frac{\beta^{F_f}(1 - \beta^{F_f})m + (1 - \beta^{F_f})n}{1 - \beta^{F_f+F_r}}G, \\ \lim_{T \rightarrow \infty} \Delta_{f,(F_f+F_r)T} &= \frac{F_f\beta^{F_f-1}(1 - \beta)(1 - \beta^{F_f})}{(1 - \beta^{F_f})(1 - \beta^{F_f+F_r})}(m - n)G, \\ \lim_{T \rightarrow \infty} \Delta_{r,(F_f+F_r)T} &= \frac{F_r\beta^{F_r-1}(1 - \beta)(1 - \beta^{F_r})}{(1 - \beta^{F_r})(1 - \beta^{F_f+F_r})}(n - m)G. \end{aligned} \quad (7)$$

- Positive correlation:** When $m = n$, B converges to mG , and Δ_f and Δ_r converge to 0. **DualOptim+ acts like Alternate.**
- Negative correlation:** When $m = -\frac{1 - \beta^{F_f}}{\beta^{F_f}(1 - \beta^{F_f})}n$, B converges to 0. **DualOptim+ acts like DualOptim.**

Experiments

 Table 1. Performance on **Fictitious Unlearning** task. The model is TOFU-tuned Phi-1.5.

Loss	Method	Phi 1.5											
		forget 1% data				forget 5% data				forget 10% data			
		UFE \uparrow	TFE \uparrow	MU \uparrow	OVR \uparrow	UFE \uparrow	TFE \uparrow	MU \uparrow	OVR \uparrow	UFE \uparrow	TFE \uparrow	MU \uparrow	OVR \uparrow
IDK+GD	Joint	78.11	45.45	18.61	40.19	72.55	58.32	36.26	50.85	71.65	64.39	33.92	50.97
	Alternate	73.35	62.49	48.14	58.03	67.73	64.30	47.81	56.91	65.82	64.46	49.54	57.34
	DO	74.75	63.51	46.46	57.80	68.49	64.34	49.50	57.96	65.87	66.84	50.25	58.30
	DO 8bit	75.58	61.23	46.73	57.57	68.34	64.33	48.41	57.37	65.81	65.60	50.38	58.05
	DO+ DO+ 8bit	75.51	67.85	47.69	59.69	67.63	67.60	51.52	59.57	65.42	66.50	51.32	58.64
		73.69	68.36	47.53	59.28	67.56	65.94	50.36	58.55	65.26	65.59	51.30	58.36

 Table 2. Performance on **Real-world Unlearning** task. The model is Llama-3-8B-Instruct.

Loss	Method	Llama 3									
		Unlearning Task				Downstream Tasks					
		UFE \uparrow	TFE \uparrow	MU \uparrow	OVR \uparrow	ARC-c \uparrow	MMLU \uparrow	TruthfulQA \uparrow	TriviaQA \uparrow	GSM8K \uparrow	AVG \uparrow
	Initial	30.55	-	61.45	46.00	55.38	64.59	37.33	50.93	76.12	56.87
IDK+GD	Joint	85.54	72.96	27.38	53.32	46.79	62.85	33.41	7.58	74.32	44.99
	Alternate	85.49	69.95	29.19	53.45	49.77	63.31	35.62	12.71	74.15	47.11
	DO	85.25	69.60	28.06	52.73	48.47	63.20	35.29	10.33	72.35	45.93
	DO 8bit	85.28	69.60	27.68	52.56	48.75	63.08	35.05	11.56	72.35	46.16
	DO+ DO+ 8bit	85.72	69.94	27.96	52.90	50.85	64.43	36.35	11.17	76.02	47.77
		85.47	69.59	33.36	55.45	52.56	64.51	36.80	17.86	75.21	49.39

 Table 3. Performance on **Safety Alignment** task. The model is Alpaca-tuned Llama 3-8B-Instruct.

Method	Alpaca-Llama 3												
	Safety					Utility							
	I-Mali \uparrow	I-CoNa \uparrow	I-Cont \uparrow	Q-Harm \uparrow	AVG \uparrow	ARC-c \uparrow	MMLU \uparrow	TruthfulQA \uparrow	TriviaQA \uparrow	GSM8K \uparrow	AVG \uparrow	OVR \uparrow	XSTest \downarrow
Initial	28.00	38.76	55.00	64.00	46.44	45.56	52.53	29.74	12.11	13.12	30.61	33.56	0.40
Joint	94.67	96.63	97.50	97.00	96.45	47.04	51.63	33.74	12.18	14.10	31.74	54.84	28.00
Alternate	97.00	97.38	97.50	99.67	97.89	46.81	50.83	34.60	13.83	12.94	31.80	55.67	29.20
DO	95.67	97.94	97.50	99.30	97.61	47.36	50.13	33.58	14.02	12.99	31.62	55.76	30.27
DO 8bit	96.00	98.31	95.50	99.00	97.20	47.10	50.78	33.58	13.82	12.96	31.65	54.66	28.27
DO+ DO+ 8bit	96.00	97.56	97.50	98.67	97.43	47.27	51.89	32.25	15.39	14.23	32.81	56.45	28.13
			97.50	99.33	97.79	46.93	51.66	34.47	13.71	14.73	32.30	55.61	28.27

Takeaway Messages

- We propose **DualOptim+**, which introduces a **shared base state** to capture common representations and **decoupled delta states** to preserve task-specific residuals.
- DualOptim+ is a **plug-and-play** framework applicable to any multi-objective optimization and optimizers with stored states.
- Experiments confirm that DualOptim+ achieves a **superior trade-off between** forgetting efficacy and model utility in MU. Our method can also be extended to more general alignment tasks.

Codes on Github

