

## A APPENDIX

### A.1 SYMBOLS AND NOTATIONS

- model parameter:  $\theta$ ;
- SAM perturbed model parameter:  $\theta + \epsilon_0^* = \theta + \rho \cdot \frac{\nabla \mathcal{L}(\theta)}{\|\nabla \mathcal{L}(\theta)\|}$ ;
- proxy model parameter:  $\theta + \epsilon_1^* = \theta + \rho \cdot (g_s - g) / \|g_s - g\|$ ;
- proxy perturbed model parameter:  $\theta + \epsilon_1^* + \rho \cdot \nabla \mathcal{L}(\theta + \epsilon_1^*) / \|\nabla \mathcal{L}(\theta + \epsilon_1^*)\|$ ;
- the empirical loss:  $\mathcal{L}(\theta)$ , with its gradient  $g$ ;
- the SAM loss:  $\mathcal{L}(\theta) + \mathcal{R}_\rho^0(\theta) = \max_{\|\epsilon_0\| \leq \rho} \mathcal{L}(\theta + \epsilon_0)$  with its gradient  $g_s$ ;
- the C-Flat loss:  $\mathcal{L}(\theta) + \mathcal{R}_\rho^0(\theta) + \lambda \cdot \mathcal{R}_\rho^1(\theta) = \max_{\|\epsilon_0\| \leq \rho} \mathcal{L}(\theta + \epsilon_0) + \rho \cdot \max_{\|\epsilon_1\| \leq \rho} \|\nabla \mathcal{L}(\theta + \epsilon_1)\|$ , with its gradient  $g_s + g_f$ ;
- the gradient of proxy model:  $g_0 = \nabla \mathcal{L}(\theta + \epsilon_1^*)$ ;
- the gradient of proxy perturbed model:  $g_1 = \nabla \mathcal{L}(\theta + \epsilon_1^* + \rho \cdot \nabla \mathcal{L}(\theta + \epsilon_1^*) / \|\nabla \mathcal{L}(\theta + \epsilon_1^*)\|)$ ;
- the empirical loss term:  $g = \nabla \mathcal{L}(\theta)$ ;
- the zeroth-order sharpness term:  $g_s - g = \nabla \mathcal{R}_\rho^0(\theta)$ ;
- the first-order flatness term:  $g_f = \nabla \mathcal{R}_\rho^1(\theta)$ ;

### A.2 DERIVATION OF EQUATION 5

Following Bian et al. (2024); Zhang et al. (2023b), the gradient of the first-order flatness loss  $\mathcal{R}_\rho^1$  is:

$$\begin{aligned} \nabla_\theta \mathcal{R}_\rho^1(\theta) &= \rho \cdot \nabla_\theta \max_{\epsilon \in B(0, \rho)} \|\nabla \mathcal{L}(\theta + \epsilon)\| \\ &= \rho \cdot \nabla_\theta \|\nabla \mathcal{L}(\theta + \epsilon_1^*)\| \\ &= \rho \cdot \left( \frac{\partial}{\partial \theta} \|\nabla \mathcal{L}(\theta + \epsilon_1^*)\| + \frac{\partial \epsilon_1^*}{\partial \theta} \cdot \nabla_\epsilon \|\nabla \mathcal{L}(\theta + \epsilon)\|_{\epsilon=\epsilon_1^*} \right) \\ &\approx \rho \cdot \nabla_\theta \|\nabla \mathcal{L}(\theta + \epsilon_1^*)\| \quad // \text{neglect higher-order term for tractability} \end{aligned} \quad (7)$$

Here,  $\epsilon_1^*$  denotes the optimal perturbation that maximizes the gradient norm within the  $\ell_2$ -ball  $B(0, \rho)$ . To make the computation tractable, we approximate it using first-order Taylor expansion and finite differences:

$$\begin{aligned} \epsilon_1^* &= \arg \max_{\epsilon \in B(0, \rho)} \|\nabla \mathcal{L}(\theta + \epsilon)\| \\ &\approx \arg \max_{\epsilon \in B(0, \rho)} \left( (\nabla_\theta \|\nabla \mathcal{L}(\theta)\|)^T \epsilon \right) \quad // \text{first-order Taylor expansion} \\ &= \rho \cdot \frac{\nabla_\theta \|\nabla \mathcal{L}(\theta)\|}{\|\nabla_\theta \|\nabla \mathcal{L}(\theta)\|\|} \quad // \text{direction of steepest ascent} \\ &\approx \rho \cdot \frac{\nabla \mathcal{L}(\theta + \delta) - \nabla \mathcal{L}(\theta)}{\|\nabla \mathcal{L}(\theta + \delta) - \nabla \mathcal{L}(\theta)\|} \quad // \text{finite-difference approximation} \\ &= \rho \cdot \frac{g_s - g}{\|g_s - g\|}, \end{aligned} \quad (8)$$

where  $g = \nabla \mathcal{L}(\theta)$ ,  $g_s = \nabla \mathcal{L}(\theta + \delta)$ , and  $\delta = \rho' \cdot \frac{g}{\|g\|}$  is a small perturbation in the direction of the gradient, with  $\rho' \ll \rho$ .

Let  $\theta_p = \theta + \epsilon_1^*$  be the perturbed model after maximizing the gradient norm. Then Equation 7 continues as:

$$\begin{aligned} \nabla_\theta \mathcal{R}_\rho^1(\theta) &\approx \rho \cdot \nabla_\theta \|\nabla \mathcal{L}(\theta_p)\| \\ &\approx \frac{\rho}{\rho'} \cdot \left[ \nabla \mathcal{L}\left(\theta_p + \rho' \cdot \frac{\nabla \mathcal{L}(\theta_p)}{\|\nabla \mathcal{L}(\theta_p)\|}\right) - \nabla \mathcal{L}(\theta_p) \right], \end{aligned} \quad (9)$$

where  $\rho' \ll \rho$  is a small step size for finite-difference approximation.

Then, the direction-invariant component of  $\mathbf{g}_f$  with respect to  $\mathbf{g}_0$  is defined as:

$$\mathbf{g}_{vf} = \mathbf{g}_f - \frac{\langle \mathbf{g}_f, \mathbf{g}_0 \rangle}{\|\mathbf{g}_0\|^2} \cdot \mathbf{g}_0, \quad (10)$$

which represents the orthogonal projection of  $\mathbf{g}_f$  onto the subspace perpendicular to  $\mathbf{g}_0$  (i.e.,  $\mathbf{g}_{vf} \perp \mathbf{g}_0$ ), capturing a direction-invariant update toward flatness.

### A.3 C-FLAT TURBO ALGORITHM

---

#### Algorithm 1 C-Flat Turbo

---

**Input:** Training phase  $T$ , training data  $S^T$ , model parameter  $\theta$ , total iterations  $J$ , oracle loss function  $\mathcal{L}$ , learning rate  $\eta$ , C-Flat coefficient  $\lambda$ , Turbo step  $k$ ,  $\mu_{s,0} = \mu_{f,0} = 0$ ,  $\sigma_{s,0} = \sigma_{f,0} = 1e-8$ .

**Output:** Model trained at the current time  $T$  with C-Flat.

```

1: for  $j = 1$  to  $J$ , sample batch  $B_j^T$  from dataset  $S^T$  do
2:   Compute gradient:  $\mathbf{g} = \nabla \mathcal{L}(\theta)$ 
3:   Initialize update direction:  $\bar{\mathbf{g}} = \mathbf{g}$ 
4:   Update EMA statistics:  $\mu_{s,j}, \sigma_{s,j}, \mu_{f,j}, \sigma_{f,j}$  by eq. (6)
5:   if  $\|\mathbf{g}\|^2 \geq \mu_{s,j} + \sigma_{s,j}$  then
6:     if  $j \bmod k = 0$  then
7:       Compute  $\mathbf{g}_s = \nabla \mathcal{L}(\theta + \epsilon_0^*) - \nabla \mathcal{L}(\theta)$  by eq. (4)
8:       Cache  $\mathbf{g}_{vs} = \mathbf{g}_s \sin(\phi_s)$ , where  $\phi_s$  denotes the angle between  $\mathbf{g}$  and  $\mathbf{g}_s$ 
9:     else
10:      Simulate  $\mathbf{g}_s = \mathbf{g} + \beta \frac{\|\mathbf{g}\|}{\|\mathbf{g}_{vs}\|} \mathbf{g}_{vs}$ 
11:    end if
12:    Update direction:  $\bar{\mathbf{g}} = \mathbf{g}_s$ 
13:  end if
14:  if  $\|\mathbf{g}_0\|^2 \geq \mu_{f,j} + \sigma_{f,j}$  then
15:    if  $j \bmod k = 0$  then
16:      Compute  $\mathbf{g}_f = \rho \cdot \nabla \|\nabla \mathcal{L}(\theta + \epsilon_1^*)\|$  using  $\mathbf{g}_0$  and  $\mathbf{g}_1$  by eq. (5)
17:      Cache  $\mathbf{g}_{vf} = \mathbf{g}_f \sin(\phi_f)$ , where  $\phi_f$  denotes the angle between  $\mathbf{g}_0$  and  $\mathbf{g}_1$ 
18:    else
19:      Simulate  $\mathbf{g}_f = \mathbf{g}_0 + \beta \frac{\|\mathbf{g}_0\|}{\|\mathbf{g}_{vf}\|} \mathbf{g}_{vf}$ 
20:    end if
21:    Update direction:  $\bar{\mathbf{g}} = \bar{\mathbf{g}} + \lambda \cdot \mathbf{g}_f$ 
22:  end if
23:  Update turbo step  $k$  according to Section 4.6
24:  Update model parameter:  $\theta^T = \theta^T - \eta \cdot \bar{\mathbf{g}}$ 
25: end for
```

---

### A.4 HYPERPARAMETER CONFIGURATIONS

Here, we primarily provide the hyperparameter configurations for methods trained on CIFAR100. For other datasets, we follow the original settings from the open repository and keep  $k = 5$  and  $\beta = 0.8$  fixed.

### A.5 MEMORY USAGE

The cached gradients are used to substitute partial sharpness-aware gradient computations, so their memory usage heavily depends on the number of trainable parameters in the model. As shown in the table 5, although larger models require more cached gradients, the overall memory overhead remains almost negligible relative to the expansion typically caused by large architectures.

Table 4: Hyperparameter settings for CIFAR100.

	Epochs	LR	BS	Tasks	Exemplar/class	$\rho$	$\lambda$	$k$	$\beta$
iCaRL	20	1e-3	32	10	20	0.1	0.2	5	0.8
MEMO	20	1e-3	32	10	20	0.1	0.2	5	0.8
L2P	5	2e-3	16	10	-	0.02	0.2	5	0.8
Ranpac	5	1e-2	16	10	-	0.05	0.2	5	0.8
EASE	5	2.5e-3	16	10	-	0.05	0.2	5	0.8

Table 5: Memory usage of different architecture.

Method	Optimizer	Backbone	Training backbone	Trainable / total params	Memory
EASE	C-Flat	ViT-Base-16	×	1.19M / 86.99M	2.14GB
	Turbo				2.15GB
	C-Flat	ViT-Large-16	×	3.17M / 306.47M	5.34GB
	Turbo				5.37GB
iCaRL	C-Flat	ResNet-18	✓	11.17M / 11.17M	1.55GB
	Turbo				1.66GB
	C-Flat	ResNet-34	✓	21.28M / 21.28M	2.32GB
	Turbo				2.51GB

#### A.6 PER-TASK ACCURACY AND ABLATION STUDIES

Per-task accuracy provides a more detailed view of the continual learning process. As shown in Table 6, the reuse mechanism significantly reduces training speed with minimal performance loss, while the linear scheduler for step size further enhances speed, particularly for longer tasks. The adaptive trigger additionally accelerates training, as it allows basic single propagation gradient descent in certain stages. Regarding performance gains, prior works have shown that selectively applying SAM updates can outperform applying SAM throughout training. For instance, SS-SAM explicitly demonstrates that with appropriate scheduling, models can achieve comparable or even superior performance at substantially lower computational cost compared to training exclusively with SAM. Similar observations also have been reported for AE-SAM and SAM-In-Later-Phase.

Table 6: Per-task accuracy and ablation study results for EASE trained on the 10-split CIFAR100 dataset.

Method	reuse	sche.	trigger	T1	T2	T3	T4	T5	T6	T7	T8	T9	T10	Avg	Img/s
EASE	×	×	×	98.40	96.25	94.63	93.88	91.80	90.92	90.47	88.09	87.58	87.17	91.92	166.67
+C-Flat	×	×	×	98.50	96.45	94.87	94.08	91.94	91.05	90.64	88.44	87.93	87.58	92.15	44.25
	✓	×	×	98.50	96.37	94.77	93.94	91.92	91.05	90.67	88.28	87.81	87.45	92.08	67.20
+Turbo	✓	✓	×	98.40	96.31	94.74	93.89	91.90	91.00	90.70	88.25	87.75	87.40	92.03	74.63
	✓	✓	✓	98.50	<b>96.60</b>	<b>95.07</b>	<b>94.15</b>	<b>92.08</b>	<b>91.27</b>	<b>90.73</b>	<b>88.52</b>	<b>88.00</b>	<b>87.57</b>	<b>92.25</b>	<b>102.74</b>

#### A.7 DETAIL DISTANCE EVOLUTION OF GRADIENTS

Figure 7 shows the L2-norm distances between sharpness and flatness gradients and their reference gradients across tasks. While  $g$  and  $g_0$  exhibit significant fluctuations during training, the gradients  $g_{vs}$  and  $g_{vf}$ , core to zeroth-order sharpness and first-order flatness regularization, change much more slowly. This stability suggests their potential as shortcut directions for flat region exploration, bypassing the need for model ascent and backpropagation.

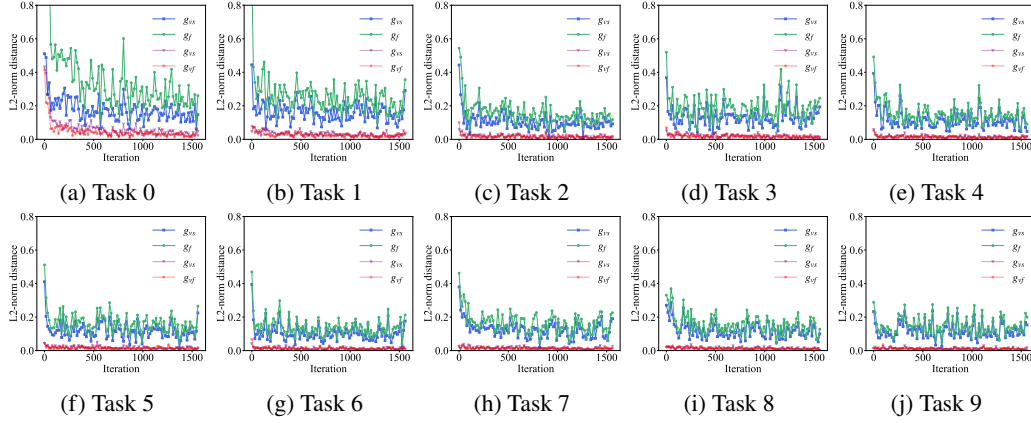


Figure 7: Visualization of L2-norm distances of the gradients every 5 steps across 10 tasks.

#### A.8 LIMITATIONS AND FUTURE WORK

This work is based on the stabilization evolution of sharpness and flatness during C-Flat optimization, where we approximate regularization terms using progressively memorized branches. Although efficient, there is still room for improvement in comparison to the speed of vanilla optimizers, particularly in minimizing the computational overhead of  $g_f$  in flatness. Then, the current framework has focused on validation using PTM-based and typical CIL tasks, but its application to other CL tasks, such as Vision-Language Models (VLMs), remains unexplored. Extending our methodology to these diverse CL settings could reveal its broader applicability and provide valuable insights into its generalization capabilities. Future work could also explore the integration of C-Flat Turbo with advanced learning paradigms, such as few-shot learning and lifelong learning, to evaluate its potential in real-world, dynamic environments.