

## A. Appendix

### A.1. Symbols and Notations of C-Flat Turbo

- 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:  $\theta + \epsilon_1^* = \theta + \rho \cdot (\mathbf{g}_s - \mathbf{g}) / \|\mathbf{g}_s - \mathbf{g}\|$ ;
- 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  $\mathbf{g}$ ;
- the SAM loss:  $\mathcal{L}_{SAM}(\theta) = \mathcal{L}(\theta) + \mathcal{R}_\rho^0(\theta) = \max_{\|\epsilon_0\| \leq \rho} \mathcal{L}(\theta + \epsilon_0)$  with its gradient  $\mathbf{g}_s$ ;
- the C-Flat loss:  $\mathcal{L}_{CFlat}(\theta) = \mathcal{L}(\theta) + \mathcal{R}_\rho^0(\theta) + \lambda \cdot \mathcal{R}_\rho^1(\theta) = \max_{\|\epsilon_0\| \leq \rho} \mathcal{L}(\theta + \epsilon_0) + \lambda \cdot \rho \max_{\|\epsilon_1\| \leq \rho} \|\nabla \mathcal{L}(\theta + \epsilon_1)\|$ , with its gradient  $\mathbf{g}_s + \lambda \mathbf{g}_f$ ;
- the gradient of proxy model:  $\mathbf{g}_0 = \nabla \mathcal{L}(\theta + \epsilon_1^*)$ ;
- the gradient of proxy perturbed model:  $\mathbf{g}_1 = \nabla \mathcal{L}(\theta + \epsilon_1^* + \rho \cdot \nabla \mathcal{L}(\theta + \epsilon_1^*) / \|\nabla \mathcal{L}(\theta + \epsilon_1^*)\|)$ ;
- the empirical loss term:  $\mathbf{g} = \nabla \mathcal{L}(\theta)$ ;
- the zeroth-order sharpness term:  $\mathbf{g}_s - \mathbf{g} = \nabla \mathcal{R}_\rho^0(\theta)$ ;
- the first-order flatness term:  $\mathbf{g}_f = \nabla \mathcal{R}_\rho^1(\theta)$ ;
- the direction-invariant sharpness component:  $\mathbf{g}_{vs} = \mathbf{g}_s - \frac{\langle \mathbf{g}_s, \mathbf{g} \rangle}{\|\mathbf{g}\|^2} \mathbf{g}$ ;
- the direction-invariant flatness component:  $\mathbf{g}_{vf} = \mathbf{g}_f - \frac{\langle \mathbf{g}_f, \mathbf{g}_0 \rangle}{\|\mathbf{g}_0\|^2} \mathbf{g}_0$ .

### A.2. Derivation of Equation 5

Following [6, 61], the gradient of the first-order flatness loss  $\mathcal{R}_\rho^1$  is:

$$\begin{aligned} \nabla_\theta \mathcal{R}_\rho^1(\theta) &= \rho \cdot \nabla_\theta \max_{\epsilon \in B(0, \rho)} \|\nabla \mathcal{L}(\theta + \epsilon)\| \quad (10) \\ &= \rho \cdot \nabla_\theta \|\nabla \mathcal{L}(\theta + \epsilon_1^*)\| \\ &= \rho \cdot \left( \frac{\partial}{\partial \theta} \|\nabla \mathcal{L}(\theta + \epsilon_1^*)\| + \frac{\partial \epsilon_1^*}{\partial \theta} \cdot \nabla_\epsilon \|\nabla \mathcal{L}(\theta + \epsilon)\| \Big|_{\epsilon = \epsilon_1^*} \right) \\ &\approx \rho \cdot \nabla_\theta \|\nabla \mathcal{L}(\theta + \epsilon_1^*)\| \end{aligned}$$

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:

$$\begin{aligned} \epsilon_1^* &= \arg \max_{\epsilon \in B(0, \rho)} \|\nabla \mathcal{L}(\theta + \epsilon)\| \quad (11) \\ &\approx \arg \max_{\epsilon \in B(0, \rho)} \left( (\nabla_\theta \|\nabla \mathcal{L}(\theta)\|)^T \epsilon \right) \\ &= \rho \cdot \frac{\nabla_\theta \|\nabla \mathcal{L}(\theta)\|}{\|\nabla_\theta \|\nabla \mathcal{L}(\theta)\|\|} \\ &\approx \rho \cdot \frac{\nabla \mathcal{L}(\theta + \delta) - \nabla \mathcal{L}(\theta)}{\|\nabla \mathcal{L}(\theta + \delta) - \nabla \mathcal{L}(\theta)\|} \\ &= \rho \cdot \frac{\mathbf{g}_s - \mathbf{g}}{\|\mathbf{g}_s - \mathbf{g}\|}, \end{aligned}$$

where  $\mathbf{g} = \nabla \mathcal{L}(\theta)$ ,  $\mathbf{g}_s = \nabla \mathcal{L}(\theta + \delta)$ , and  $\delta = \rho' \cdot \frac{\mathbf{g}}{\|\mathbf{g}\|}$  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 Eq. 10 continues as:

$$\begin{aligned} \nabla_\theta \mathcal{R}_\rho^1(\theta) &\approx \rho \cdot \nabla_\theta \|\nabla \mathcal{L}(\theta_p)\| \quad (12) \\ &\approx \frac{\rho}{\rho'} \cdot \left[ \nabla \mathcal{L}(\theta_p + \rho' \cdot \frac{\nabla \mathcal{L}(\theta_p)}{\|\nabla \mathcal{L}(\theta_p)\|}) - \nabla \mathcal{L}(\theta_p) \right], \end{aligned}$$

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

In our symbol table, we define

$$\mathbf{g}_0 := \nabla \mathcal{L}(\theta_p) = \nabla \mathcal{L}(\theta + \epsilon_1^*), \quad (13)$$

$$\mathbf{g}_1 := \nabla \mathcal{L}(\theta_p + \rho' \cdot \frac{\nabla \mathcal{L}(\theta_p)}{\|\nabla \mathcal{L}(\theta_p)\|}), \quad (14)$$

Then Eq. 12 can be written as the finite-difference form

$$\mathbf{g}_f \approx \frac{\rho}{\rho'} (\mathbf{g}_1 - \mathbf{g}_0). \quad (15)$$

For simplicity of notation, we absorb the constant factor  $\rho/\rho'$  into  $\lambda$ , so that at the level of directions we can view

$$\mathbf{g}_f \propto \mathbf{g}_1 - \mathbf{g}_0.$$

Finally, we define the direction-invariant component of  $\mathbf{g}_f$  with respect to  $\mathbf{g}_0$ :

$$\mathbf{g}_{vf} = \mathbf{g}_f - \frac{\langle \mathbf{g}_f, \mathbf{g}_0 \rangle}{\|\mathbf{g}_0\|^2} \cdot \mathbf{g}_0, \quad (16)$$

which is orthogonal to  $\mathbf{g}_0$  and is the quantity tracked by the EMA variable in the implementation.

### A.3. Hyperparameter Settings

We report the main hyperparameter settings for methods trained on CIFAR100 in Table 4. For the other datasets, we follow the original settings from the open-source repository and keep  $k = 5$  and  $\beta = 0.8$  fixed to balance efficiency and performance.

We also study the sensitivity of the scheduler and trigger threshold in C-Flat Turbo on CIFAR100. For the scheduler  $k_t = k_0 + c \cdot \frac{t}{N}$ , larger  $k_0$  generally improves throughput with negligible performance changes, while  $c = 10$  provides a good trade-off between speed and accuracy, as shown in Table 5. This suggests that moderately increasing the turbo interval over the task sequence is sufficient to obtain most of the efficiency gains. For the trigger threshold  $\|\mathbf{g}_{0j}\|^2 > \mu_{f,j} + m \cdot \sigma_{f,j}$ , Table 6 shows that  $m$  between 0.2 and 1 performs best overall. Larger values improve efficiency but tend to degrade performance, whereas smaller values may lead to less stable trigger behavior.

Table 4. Hyperparameter settings for CIFAR100.

Methods	Epochs	LR	BS	Tasks	Exemplar	$\rho$	$\lambda$	$k$	$\beta$
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. Sensitivity analysis of the scheduler hyperparameters in C-Flat Turbo on CIFAR100.

Optimizer	C-Flat Turbo								
$k_0$	2	2	2	5	5	5	10	10	10
$c$	5	10	20	5	10	20	5	10	20
Avg	92.19	92.18	92.05	92.08	92.08	92.01	92.07	92.04	92.04
Last	87.67	87.59	87.45	87.62	87.54	87.42	87.51	87.59	87.46
Img/s	99.17	99.91	63.19	104.86	102.74	99.05	104.64	106.89	107.39

Table 6. Sensitivity analysis of the trigger threshold hyperparameter in C-Flat Turbo on CIFAR100.

$m$	0.1	0.2	0.5	1	2	5
Avg	91.98	92.07	92.05	92.06	91.86	91.94
Last	87.43	87.48	87.39	87.50	87.07	87.16
Img/s	75.56	105.43	97.69	102.74	154.52	135.18

#### A.4. 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 Table 7, although larger models require more cached gradients, the overall memory overhead remains almost negligible relative to the expansion typically caused by large architectures.

Table 7. Memory usage of different architectures. EASE uses a frozen backbone, while iCaRL updates the full backbone.

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

#### A.5. Comparison to Similar Works

We further compare C-Flat Turbo with SS-SAM and AE-SAM on L2P and EASE, as shown in Table 8. Overall, C-Flat Turbo consistently achieves the best performance on both benchmarks, while maintaining competitive efficiency.

Table 8. Comparison with SS-SAM and AE-SAM on L2P (left) and EASE (right).

Method	L2P			EASE		
	Avg	Last	Img/s	Avg	Last	Img/s
+SS-SAM	89.50	84.45	77.42 (70.20%)	90.90	<u>87.87</u>	91.65 (54.99%)
+AE-SAM	<u>89.70</u>	<u>84.48</u>	80.86 (73.32%)	<u>92.03</u>	87.37	112.58 (67.54%)
+C-Flat Turbo	<b>89.78</b>	<b>84.69</b>	65.50 (59.4%)	<b>92.36</b>	<b>87.96</b>	102.74 (61.6%)

#### A.6. Other CL Settings

We further evaluate C-Flat Turbo with DUCT on DIL and with DualPrompt on TIL and CIL. As shown in Tables 9 and 10, Turbo consistently improves performance over the corresponding baselines.

Table 9. Results on DomainNet under the DIL setting.

Method	Clipart	Infograph	Painting	Quickdraw	Real	Sketch	Avg
DUCT	74.21	60.52	67.73	<u>63.50</u>	69.51	68.63	67.35
+C-Flat	<b>74.93</b>	<u>60.88</u>	<b>68.10</b>	63.30	<u>70.22</u>	<u>69.41</u>	<u>67.81</u>
+C-Flat Turbo	<u>74.89</u>	<b>60.91</b>	<u>68.01</u>	<b>63.75</b>	<b>70.82</b>	<b>69.50</b>	<b>67.98</b>

Table 10. Results on CIFAR100 under the CIL and TIL settings.

Method	CIL		TIL	
	Avg	Last	Avg	Last
DualPrompt	85.77	<b>89.10</b>	97.79	<b>99.40</b>
+C-Flat	85.85	87.80	97.67	99.30
+C-Flat Turbo	<b>86.01</b>	87.00	<b>97.94</b>	<b>99.40</b>

#### A.7. Per-task Accuracy and Ablation Studies

Per-task accuracy provides a more detailed view of the continual learning process. As shown in Table 11, the reuse mechanism significantly reduces training time with minimal performance loss, while the linear step-size scheduler further improves speed, particularly for longer task sequences. 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 [63] 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 [24] and SAM-In-Late-Phase [72].

#### A.8. Detailed Evolution of Gradient Distances

Figure 7 shows the L2-norm distances between sharpness and flatness gradients and their reference gradients across

Table 11. Per-task accuracy and ablation study results for EASE trained on the 10-split CIFAR100 dataset.

Method	reuse	scheduler	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
+C-Flat 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>

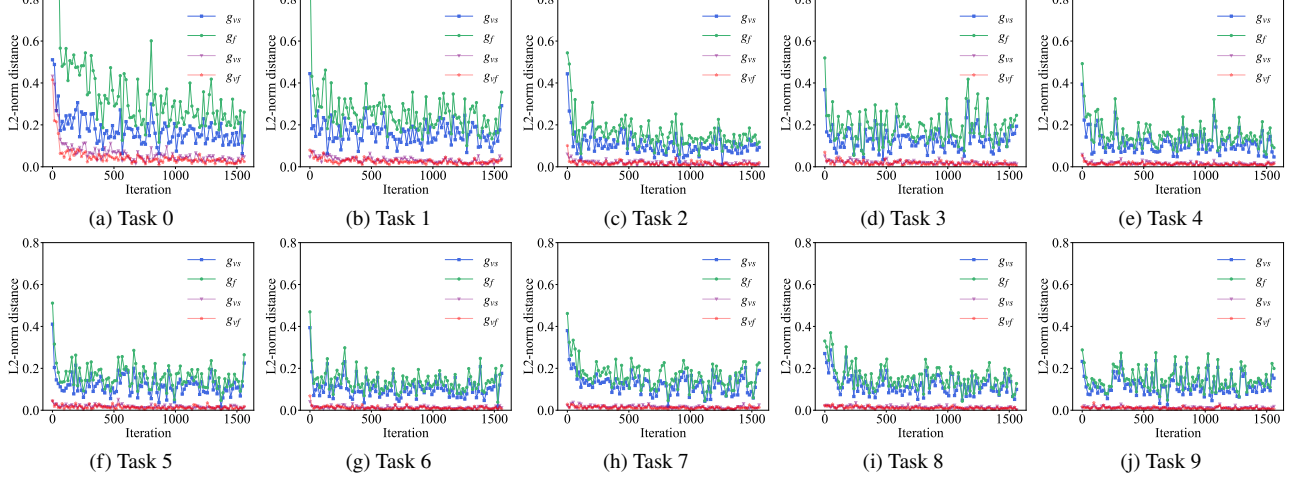


Figure 7. Visualization of L2-norm distances of the gradients every 5 steps across 10 tasks.

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.

### A.9. Convergence of Turbo in the Surrogate Steps

With the definitions in Appendix A.1, we denote C-Flat gradient at iteration  $t$  as

$$\mathbf{G}_t := \nabla \mathcal{L}_{\text{CFlat}}(\theta_t) = \mathbf{g}_s + \lambda \mathbf{g}_f, \quad (17)$$

then we can rewrite

$$\mathbf{G}_t = \underbrace{\mathbf{g}}_{\text{empirical gradient}} + \underbrace{(\mathbf{g}_s - \mathbf{g})}_{\text{sharpness increment}} + \lambda \underbrace{\mathbf{g}_f}_{\text{flatness term}}, \quad (18)$$

where  $\mathbf{g} = \nabla \mathcal{L}(\theta_t)$  is the empirical gradient.

In the exact C-Flat updates, the sharpness increment  $(\mathbf{g}_s - \mathbf{g})$  is obtained by subtracting the empirical gradient from the SAM gradient, and the flatness term  $\mathbf{g}_f$  is obtained via the finite-difference approximation  $\mathbf{g}_f \approx \frac{\rho'}{\rho'}(\mathbf{g}_1 - \mathbf{g}_0)$  as in Eq. 15.

In Turbo, we do not recompute  $\mathbf{g}_s$  and  $\mathbf{g}_f$  at every surrogate step. Instead, we maintain EMA states  $\mathbf{g}_{vs,t}$  and  $\mathbf{g}_{vf,t}$  to track

$$\mathbf{u}_t^{(s)} := \mathbf{g}_s - \mathbf{g}, \quad \mathbf{u}_t^{(f)} := \mathbf{g}_f,$$

and we construct deterministic surrogates

$$A_t(\mathbf{g}_{vs,t-1}) \approx \mathbf{u}_t^{(s)}, \quad (19)$$

$$B_t(\mathbf{g}_{vf,t-1}) \approx \mathbf{u}_t^{(f)}. \quad (20)$$

Therefore, at a Turbo surrogate step, the update direction is modeled as

$$\tilde{\mathbf{g}}_t = \underbrace{\mathbf{g}}_{\nabla \mathcal{L}(\theta_t)} + A_t(\mathbf{g}_{vs,t-1}) + \lambda B_t(\mathbf{g}_{vf,t-1}) + \xi_t, \quad (21)$$

where  $\mathbf{g} = \nabla \mathcal{L}(\theta_t)$  is the empirical gradient and  $\xi_t$  is a zero-mean stochastic noise term capturing minibatch sampling and finite-difference randomness.

To relate Eq. 21 to the true C-Flat gradient in Eq. 18, we decompose the deterministic surrogates as

$$A_t(\mathbf{g}_{vs,t-1}) = \mathbf{u}_t^{(s)} + \delta_t^{(s)}, \quad (22)$$

$$B_t(\mathbf{g}_{vf,t-1}) = \mathbf{u}_t^{(f)} + \delta_t^{(f)}, \quad (23)$$

where  $\mathbf{u}_t^{(s)} := \mathbf{g}_s - \mathbf{g}$  and  $\mathbf{u}_t^{(f)} := \mathbf{g}_f$ , and  $\delta_t^{(s)}, \delta_t^{(f)}$  are deterministic approximation errors (given  $\mathcal{F}_t$ ) coming from EMA tracking, orthogonal decomposition residuals, and finite differences.

Substituting into Eq. 21 and comparing with Eq. 18, we obtain

$$\begin{aligned} \tilde{\mathbf{g}}_t &= \mathbf{g} + (\mathbf{u}_t^{(s)} + \delta_t^{(s)}) + \lambda(\mathbf{u}_t^{(f)} + \delta_t^{(f)}) + \xi_t \quad (24) \\ &= \underbrace{(\mathbf{g} + \mathbf{u}_t^{(s)} + \lambda\mathbf{u}_t^{(f)})}_{=: \mathbf{G}_t} + \underbrace{(\delta_t^{(s)} + \lambda\delta_t^{(f)})}_{=: \mathbf{b}_t} + \xi_t \\ &= \mathbf{G}_t + \mathbf{b}_t + \xi_t. \end{aligned}$$

**Assumption 1** (Generalized smoothness). *There exists  $L > 0$  such that for all  $\theta, \theta'$  and any generalized gradients  $\mathbf{g} \in \partial \mathcal{L}_{\text{CFlat}}(\theta), \mathbf{g}' \in \partial \mathcal{L}_{\text{CFlat}}(\theta')$ ,*

$$\|\mathbf{g} - \mathbf{g}'\| \leq L\|\theta - \theta'\|.$$

**Assumption 2** (Bounded gradients and noise). *There exist constants  $G_{\max}, \Sigma \geq 0$  such that for all  $t$  and any  $\mathbf{g} \in \partial \mathcal{L}_{\text{CFlat}}(\theta_t)$ ,*

$$\|\mathbf{g}\| \leq G_{\max}, \quad \mathbb{E}[\xi_t | \mathcal{F}_t] = \mathbf{0}, \quad \mathbb{E}[\|\xi_t\|^2 | \mathcal{F}_t] \leq \Sigma^2.$$

**Assumption 3** (Deterministic approximation bias). *The bias  $\mathbf{b}_t$  satisfies  $\|\mathbf{b}_t\| \leq \Delta_t$  almost surely, where  $\Delta_t \geq 0$  is  $\mathcal{F}_t$ -measurable. In Turbo,  $\Delta_t$  aggregates:*

1. EMA tracking error in estimating  $\mathbf{g}_s - \mathbf{g}$  and  $\mathbf{g}_f$  via  $\mathbf{g}_{v_s,t}$  and  $\mathbf{g}_{v_f,t}$ ;
2. orthogonal decomposition residuals in constructing the direction-invariant sharpness and flatness component  $\mathbf{g}_{v_s}$  and  $\mathbf{g}_{v_f}$ ;
3. finite-difference approximation errors in the proxy and proxy-perturbed gradients  $\mathbf{g}_0, \mathbf{g}_1$ .

**Assumption 4** (Stepsize and perturbation schedule). *We use the schedule in [61]:*

$$\eta_t = \frac{\eta_0}{\sqrt{t}}, \quad \rho_t = \frac{\rho_0}{\sqrt{t}},$$

with  $\eta_0$  chosen so that  $\eta_t \leq 1/(4L)$  for all  $t$ .

Finally, the EMA tracking of  $\mathbf{u}_t^{(s)}$  and  $\mathbf{u}_t^{(f)}$  yields a bias sequence  $\Delta_t$  satisfying

$$\sum_{t=1}^T \eta_t \Delta_t^2 \leq C_{\Delta,1} \sqrt{T} + C_{\Delta,2} \log T \quad (25)$$

for some constants  $C_{\Delta,1}, C_{\Delta,2} \geq 0$ .

**Justification of the deterministic bias assumption.** The weighted bound on  $\Delta_t$  in Eq. 25 is supported by the empirical directional stability of the cached components. As illustrated in Figure 3 and 7 in the main text, the direction invariant components  $\mathbf{g}_{v_s,t}$  and  $\mathbf{g}_{v_f,t}$  exhibit substantially higher

cosine similarity across consecutive iterations than the empirical gradient  $\mathbf{g}_t$  and the proxy gradient  $\mathbf{g}_{0,t}$ . This high stability suggests that the surrogate approximations  $A_t(\cdot)$  and  $B_t(\cdot)$  remain accurate within Turbo intervals (i.e.,  $\|\delta_t^{(s)}\|$  and  $\|\delta_t^{(f)}\|$  stay small), which in turn justifies modeling the resulting approximation bias  $\Delta_t$  controlled.

**Lemma 1** (Generalized descent lemma). *Under the generalized smoothness assumption, for any update  $\theta_{t+1} = \theta_t - \eta_t \tilde{\mathbf{g}}_t$  and any  $\mathbf{G}_t \in \partial \mathcal{L}_{\text{CFlat}}(\theta_t)$ ,*

$$\mathcal{L}_{\text{CFlat}}(\theta_{t+1}) \leq \mathcal{L}_{\text{CFlat}}(\theta_t) - \eta_t \langle \mathbf{G}_t, \tilde{\mathbf{g}}_t \rangle + \frac{L}{2} \eta_t^2 \|\tilde{\mathbf{g}}_t\|^2.$$

*Proof.* By generalized smoothness and the descent inequality for  $L$ -smooth functions, we have

$$\mathcal{L}_{\text{CFlat}}(\theta_{t+1}) \leq \mathcal{L}_{\text{CFlat}}(\theta_t) + \langle \mathbf{G}_t, \theta_{t+1} - \theta_t \rangle + \frac{L}{2} \|\theta_{t+1} - \theta_t\|^2.$$

Substitute  $\theta_{t+1} - \theta_t = -\eta_t \tilde{\mathbf{g}}_t$  and rearrange.  $\square$

For EMA tracking, the bias sequence induced by  $\mathbf{g}_{v_s,t}$  and  $\mathbf{g}_{v_f,t}$  is absorbed into  $\Delta_t$ , and we simply assume that it satisfies the bound in Eq. 25, which is standard for exponential moving averages under mild regularity conditions.

**Single-step expected decrease.** Starting from the descent lemma and using  $\tilde{\mathbf{g}}_t = \mathbf{G}_t + \mathbf{b}_t + \xi_t$ , we obtain

$$\begin{aligned} \mathcal{L}_{\text{CFlat}}(\theta_{t+1}) &\leq \mathcal{L}_{\text{CFlat}}(\theta_t) - \eta_t \langle \mathbf{G}_t, \mathbf{G}_t + \mathbf{b}_t + \xi_t \rangle \quad (26) \\ &\quad + \frac{L}{2} \eta_t^2 \|\mathbf{G}_t + \mathbf{b}_t + \xi_t\|^2. \end{aligned}$$

Taking full expectation, and using  $\mathbb{E}[\langle \mathbf{G}_t, \xi_t \rangle | \mathcal{F}_t] = 0$ , we get

$$\begin{aligned} \mathbb{E}[\mathcal{L}_{\text{CFlat}}(\theta_{t+1})] &\leq \mathbb{E}[\mathcal{L}_{\text{CFlat}}(\theta_t)] - \eta_t \mathbb{E}[\|\mathbf{G}_t\|^2] \quad (27) \\ &\quad - \eta_t \mathbb{E}[\langle \mathbf{G}_t, \mathbf{b}_t \rangle] + \frac{L}{2} \eta_t^2 \mathbb{E}[\|\mathbf{G}_t + \mathbf{b}_t + \xi_t\|^2]. \end{aligned}$$

**a) Bias inner product term:** By Cauchy–Schwarz and  $\|\mathbf{b}_t\| \leq \Delta_t$ , we have

$$|\langle \mathbf{G}_t, \mathbf{b}_t \rangle| \leq \|\mathbf{G}_t\| \|\mathbf{b}_t\| \leq \frac{1}{2} \|\mathbf{G}_t\|^2 + \frac{1}{2} \Delta_t^2.$$

Using  $|\mathbb{E}[Z]| \leq \mathbb{E}[|Z|]$ , this implies

$$-\eta_t \mathbb{E}[\langle \mathbf{G}_t, \mathbf{b}_t \rangle] \leq \frac{\eta_t}{2} \mathbb{E}[\|\mathbf{G}_t\|^2] + \frac{\eta_t}{2} \Delta_t^2. \quad (28)$$

**b) Second moment term:** For any three vectors  $x, y, z$ , Jensen's inequality gives

$$\|x + y + z\|^2 \leq 3(\|x\|^2 + \|y\|^2 + \|z\|^2).$$

Applying this with  $x = \mathbf{G}_t, y = \mathbf{b}_t$  and  $z = \xi_t$ , we obtain

$$\|\mathbf{G}_t + \mathbf{b}_t + \xi_t\|^2 \leq 3\|\mathbf{G}_t\|^2 + 3\|\mathbf{b}_t\|^2 + 3\|\xi_t\|^2,$$

hence, using  $\|\mathbf{b}_t\| \leq \Delta_t$  and  $\mathbb{E}[\|\xi_t\|^2 | \mathcal{F}_t] \leq \Sigma^2$ ,

$$\mathbb{E}[\|\mathbf{G}_t + \mathbf{b}_t + \xi_t\|^2] \leq 3\mathbb{E}\|\mathbf{G}_t\|^2 + 3\Delta_t^2 + 3\Sigma^2. \quad (29)$$

Substituting Eq. 28 and Eq. 29 into Eq. 27, we obtain

$$\begin{aligned} \mathbb{E}[\mathcal{L}_{\text{CFlat}}(\boldsymbol{\theta}_{t+1})] &\leq \mathbb{E}[\mathcal{L}_{\text{CFlat}}(\boldsymbol{\theta}_t)] - \eta_t \mathbb{E}\|\mathbf{G}_t\|^2 + \frac{\eta_t}{2} \mathbb{E}\|\mathbf{G}_t\|^2 \\ &\quad + \frac{\eta_t}{2} \Delta_t^2 + \frac{L}{2} \eta_t^2 (3\mathbb{E}\|\mathbf{G}_t\|^2 + 3\Delta_t^2 + 3\Sigma^2) \\ &= \mathbb{E}[\mathcal{L}_{\text{CFlat}}(\boldsymbol{\theta}_t)] - \left(\frac{\eta_t}{2} - \frac{3L}{2}\eta_t^2\right) \mathbb{E}\|\mathbf{G}_t\|^2 \\ &\quad + \left(\frac{\eta_t}{2} + \frac{3L}{2}\eta_t^2\right) \Delta_t^2 + \frac{3L}{2}\eta_t^2 \Sigma^2. \end{aligned} \quad (30)$$

Under the stepsize condition  $\eta_t \leq 1/(4L)$ , we have

$$\frac{\eta_t}{2} - \frac{3L}{2}\eta_t^2 \geq \frac{\eta_t}{8},$$

and hence

$$\frac{\eta_t}{8} \mathbb{E}\|\mathbf{G}_t\|^2 \leq \mathbb{E}[\mathcal{L}_{\text{CFlat}}(\boldsymbol{\theta}_t)] - \mathbb{E}[\mathcal{L}_{\text{CFlat}}(\boldsymbol{\theta}_{t+1})] + \Psi_t,$$

where

$$\Psi_t := \left(\frac{\eta_t}{2} + \frac{3L}{2}\eta_t^2\right) \Delta_t^2 + \frac{3L}{2}\eta_t^2 \Sigma^2.$$

**Summation.** Summing from  $t = 1$  to  $T$ , we obtain

$$\sum_{t=1}^T \frac{\eta_t}{8} \mathbb{E}\|\mathbf{G}_t\|^2 \leq \mathcal{L}_{\text{CFlat}}(\boldsymbol{\theta}_1) - L^* + \sum_{t=1}^T \Psi_t,$$

where  $L^* = \inf_{\boldsymbol{\theta}} \mathcal{L}_{\text{CFlat}}(\boldsymbol{\theta})$ . Hence

$$\sum_{t=1}^T \eta_t \mathbb{E}\|\mathbf{G}_t\|^2 \leq 8(\mathcal{L}_{\text{CFlat}}(\boldsymbol{\theta}_1) - L^*) + 8 \sum_{t=1}^T \Psi_t. \quad (31)$$

Using  $\eta_t = \eta_0/\sqrt{t}$  and the assumption  $\sum_{t=1}^T \eta_t \Delta_t^2 \leq C_{\Delta,1}\sqrt{T} + C_{\Delta,2} \log T$ , together with  $\sum_{t=1}^T \eta_t^2 \leq \eta_0^2(1 + \log T)$ , we obtain

$$\sum_{t=1}^T \Psi_t \leq C_1 \sqrt{T} + C_2 \log T + C_3$$

for some constants  $C_1, C_2, C_3$ . Thus there exist  $C_A, C_B, C_C$  such that

$$\sum_{t=1}^T \eta_t \mathbb{E}\|\mathbf{G}_t\|^2 \leq C_A + C_B \sqrt{T} + C_C \log T.$$

Since  $\eta_T = \eta_0/\sqrt{T}$  and  $\eta_t$  is nonincreasing, we have

$$\eta_T \sum_{t=1}^T \mathbb{E}\|\mathbf{G}_t\|^2 \leq \sum_{t=1}^T \eta_t \mathbb{E}\|\mathbf{G}_t\|^2,$$

so

$$\frac{\eta_0}{\sqrt{T}} \sum_{t=1}^T \mathbb{E}\|\mathbf{G}_t\|^2 \leq C_A + C_B \sqrt{T} + C_C \log T.$$

Dividing by  $T$  yields

$$\frac{1}{T} \sum_{t=1}^T \mathbb{E}\|\mathbf{G}_t\|^2 \leq \frac{C_A}{\eta_0 \sqrt{T}} + \frac{C_B}{\eta_0} + \frac{C_C \log T}{\eta_0 \sqrt{T}}.$$

**Theorem 1** (Turbo convergence in surrogate steps). *Under the above assumptions and the EMA tracking condition Eq. 25, there exist constants  $\tilde{C}_1, \tilde{C}_2, \tilde{C}_3 \geq 0$  such that*

$$\frac{1}{T} \sum_{t=1}^T \mathbb{E}[\|\nabla \mathcal{L}_{\text{CFlat}}(\boldsymbol{\theta}_t)\|^2] \leq \frac{\tilde{C}_1 + \tilde{C}_2 \log T}{\sqrt{T}} + \tilde{C}_3.$$

Moreover, if the bias sequence satisfies  $\Delta_t = O(1/\sqrt{t})$  so that  $\sum_{t=1}^T \eta_t \Delta_t^2 = O(1 + \log T)$ , then  $\tilde{C}_3 = 0$  and Turbo recovers a GAM-style rate [6]:

$$\frac{1}{T} \sum_{t=1}^T \mathbb{E}[\|\nabla \mathcal{L}_{\text{CFlat}}(\boldsymbol{\theta}_t)\|^2] = O\left(\frac{1 + \log T}{\sqrt{T}}\right).$$

## A.10. C-Flat Turbo Algorithm

---

### Algorithm 1 C-Flat Turbo

---

**Input:** Training phase  $T$ , training data  $S^T$ , model parameter  $\boldsymbol{\theta}$ , total iterations  $J$ , oracle loss function  $\mathcal{L}$ , learning rate  $\eta$ , C-Flat coefficient  $\lambda$ , Turbo step  $k$ ,  $\mu_{s,0} = \mu_{f,0} = 0$ ,  $\sigma_{s,0} = \sigma_{f,0} = 10^{-8}$ .

**Output:** Model trained at the current time  $T$  with Turbo.

- 1: **for**  $j = 1$  to  $J$ , sample batch  $B_j^T$  from dataset  $S^T$  **do**
  - 2:   Compute empirical gradient:  $\mathbf{g} = \nabla \mathcal{L}(\boldsymbol{\theta})$
  - 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. 8
  - 5:   **if**  $\|\mathbf{g}\|^2 \geq \mu_{s,j} + \sigma_{s,j}$  **then**
  - 6:     **if**  $j \bmod k = 0$  **then**
  - 7:       Compute SAM gradient  $\mathbf{g}_s = \nabla \mathcal{L}(\boldsymbol{\theta} + \boldsymbol{\epsilon}_0^*)$  by Eq. 4
  - 8:       Cache direction-invariant sharpness component  $\mathbf{g}_{vs}$  by Eq. 6
  - 9:     **else**
  - 10:       Simulate sharpness increment using cached  $\mathbf{g}_{vs}$ , that means  $\mathbf{g}_s = \mathbf{g} + \beta \frac{\|\mathbf{g}\|}{\|\mathbf{g}_{vs}\|} \mathbf{g}_{vs}$
  - 11:     **end if**
  - 12:     Update direction:  $\bar{\mathbf{g}} = \mathbf{g}_s$
  - 13:   **end if**
  - 14:   Compute proxy model parameter  $\boldsymbol{\theta}_p = \boldsymbol{\theta} + \boldsymbol{\epsilon}_1^*$  by Eq. 5
  - 15:   Compute proxy gradient:  $\mathbf{g}_0 = \nabla \mathcal{L}(\boldsymbol{\theta}_p)$
  - 16:   **if**  $\|\mathbf{g}_0\|^2 \geq \mu_{f,j} + \sigma_{f,j}$  **then**
  - 17:     **if**  $j \bmod k = 0$  **then**
  - 18:       Compute proxy-perturbed gradient  $\mathbf{g}_1 = \nabla \mathcal{L}\left(\boldsymbol{\theta}_p + \rho' \frac{\mathbf{g}_0}{\|\mathbf{g}_0\|}\right)$
  - 19:       Compute first-order flatness gradient  $\mathbf{g}_f$  from  $\mathbf{g}_0, \mathbf{g}_1$  by Eq. 5
  - 20:       Cache direction-invariant flatness component  $\mathbf{g}_{vf}$  by Eq. 7
  - 21:     **else**
  - 22:       Simulate flatness direction using cached  $\mathbf{g}_{vf}$ , that means  $\mathbf{g}_f = \mathbf{g}_0 + \beta \frac{\|\mathbf{g}_0\|}{\|\mathbf{g}_{vf}\|} \mathbf{g}_{vf}$
  - 23:     **end if**
  - 24:     Update direction:  $\bar{\mathbf{g}} = \bar{\mathbf{g}} + \lambda \cdot \mathbf{g}_f$
  - 25:   **end if**
  - 26:   Optionally update the Turbo step  $k$  according to Section 4.6
  - 27:   Update model parameter:  $\boldsymbol{\theta}^T = \boldsymbol{\theta}^T - \eta \cdot \bar{\mathbf{g}}$
  - 28: **end for**
-