

Sharpness-Aware Minimization Efficiently Selects Flatter Minima Late In Training

Zhanpeng Zhou*¹, Mingze Wang*², Yuchen Mao¹, Bingrui Li³, Junchi Yan¹

¹Shanghai Jiao Tong University, ²Peking University, ³Tsinghua University



SHANGHAI JIAO TONG
UNIVERSITY



北京大学
PEKING UNIVERSITY



清华大学

Background: Sharpness-Aware Minimization

Sharpness-Aware Minimization (SAM) [1]

The central idea is to minimize the worst-case loss within a neighborhood of current weights, i.e.,

$$\min_{\boldsymbol{\theta}} \mathcal{L}_S^{\text{SAM}}(\boldsymbol{\theta}; \rho), \quad \text{where } \mathcal{L}_S^{\text{SAM}}(\boldsymbol{\theta}; \rho) = \max_{\|\boldsymbol{\epsilon}\|_2 \leq \rho} \mathcal{L}_S(\boldsymbol{\theta} + \boldsymbol{\epsilon}).^*$$

* $\mathcal{L}_S(\boldsymbol{\theta})$ denotes the total loss over the training set.

However, finding $\boldsymbol{\epsilon}$ can be computationally intractable in practice. Thus Foret et al. [1] used first-order approximation, i.e.,

$$\boldsymbol{\epsilon} \approx \arg \max_{\|\boldsymbol{\epsilon}\|_2 \leq \rho} (\mathcal{L}_S(\boldsymbol{\theta}) + \boldsymbol{\epsilon}^\top \nabla \mathcal{L}_S(\boldsymbol{\theta})) = \rho \nabla \mathcal{L}_S(\boldsymbol{\theta}) / \|\nabla \mathcal{L}_S(\boldsymbol{\theta})\|_2.$$

Consequently, the update rule of SAM with stochastic gradient is,

$$\boldsymbol{\theta}^{t+1} = \boldsymbol{\theta}^t - \eta \nabla \mathcal{L}_{\xi_t} \left(\boldsymbol{\theta} + \rho \nabla \mathcal{L}_{\xi_t}(\boldsymbol{\theta}) / \|\nabla \mathcal{L}_{\xi_t}(\boldsymbol{\theta})\|_2 \right).^*$$

* $\mathcal{L}_{\xi_t}(\boldsymbol{\theta})$ denotes the loss over a randomly sampled mini-batch ξ_t at iteration t . η denotes the learning rate.

Background: Implicit Bias

*The effectiveness of gradient-based optimization methods can be attributed to their **implicit bias** toward solutions with favorable properties [2].*

Implicit Bias of SGD

- SGD and its variants tends to find flat minima, which often generalize well.

Implicit Bias of SAM

- SAM tends to find flatter minima over SGD, which represents a form of implicit bias*.

*Though SAM is inspired from sharpness regularization, its practical implementation, which minimizes a first-order approximation of the original objective, doesn't explicitly achieve this.

Understanding the mechanism behind the implicit bias of SAM towards flatter minima is crucial to explain its effectiveness.

SAM Selects Flatter Minima Late In Training

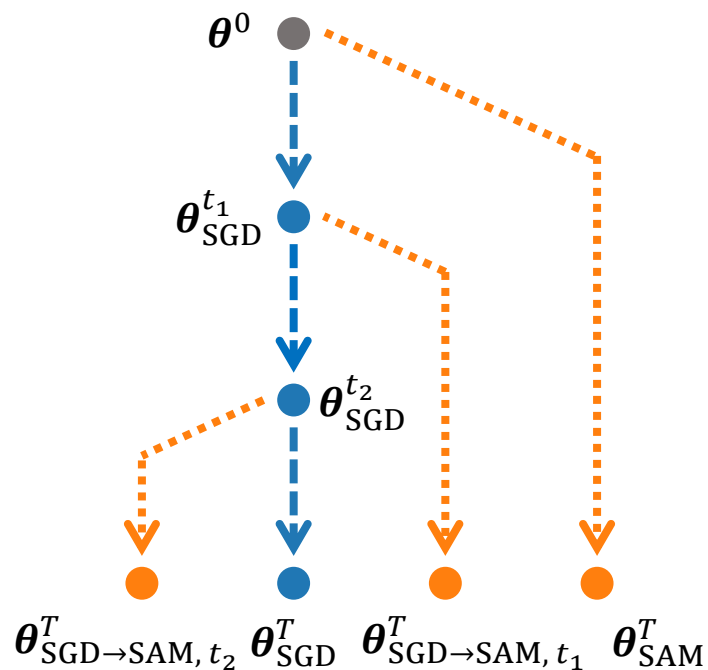


Fig. 1: **Illustration of the switching method.** Blue dashed lines represent SGD training, while orange dashed lines represent SAM training. t_1, t_2 denotes two switching points.

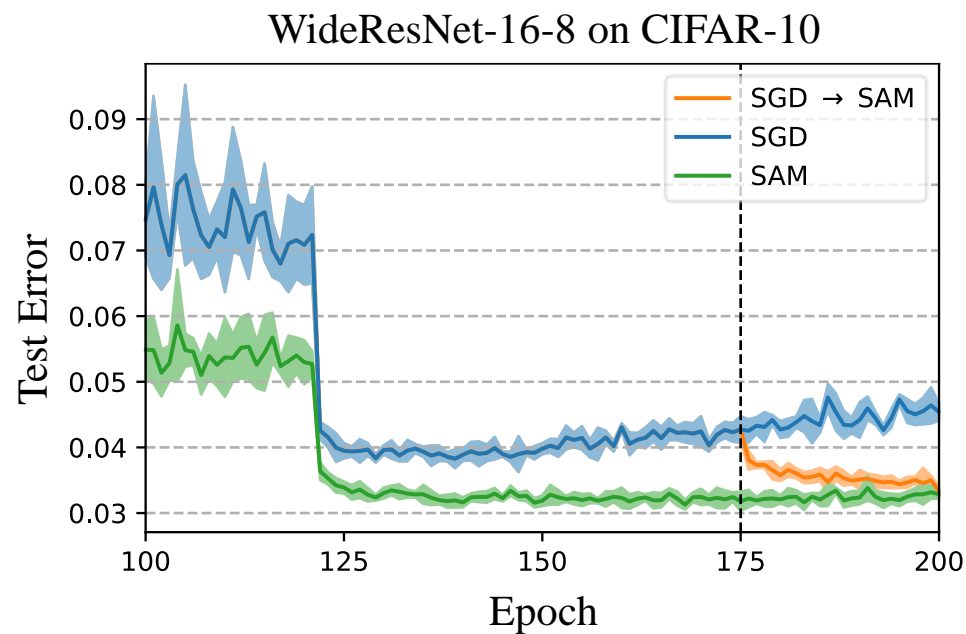
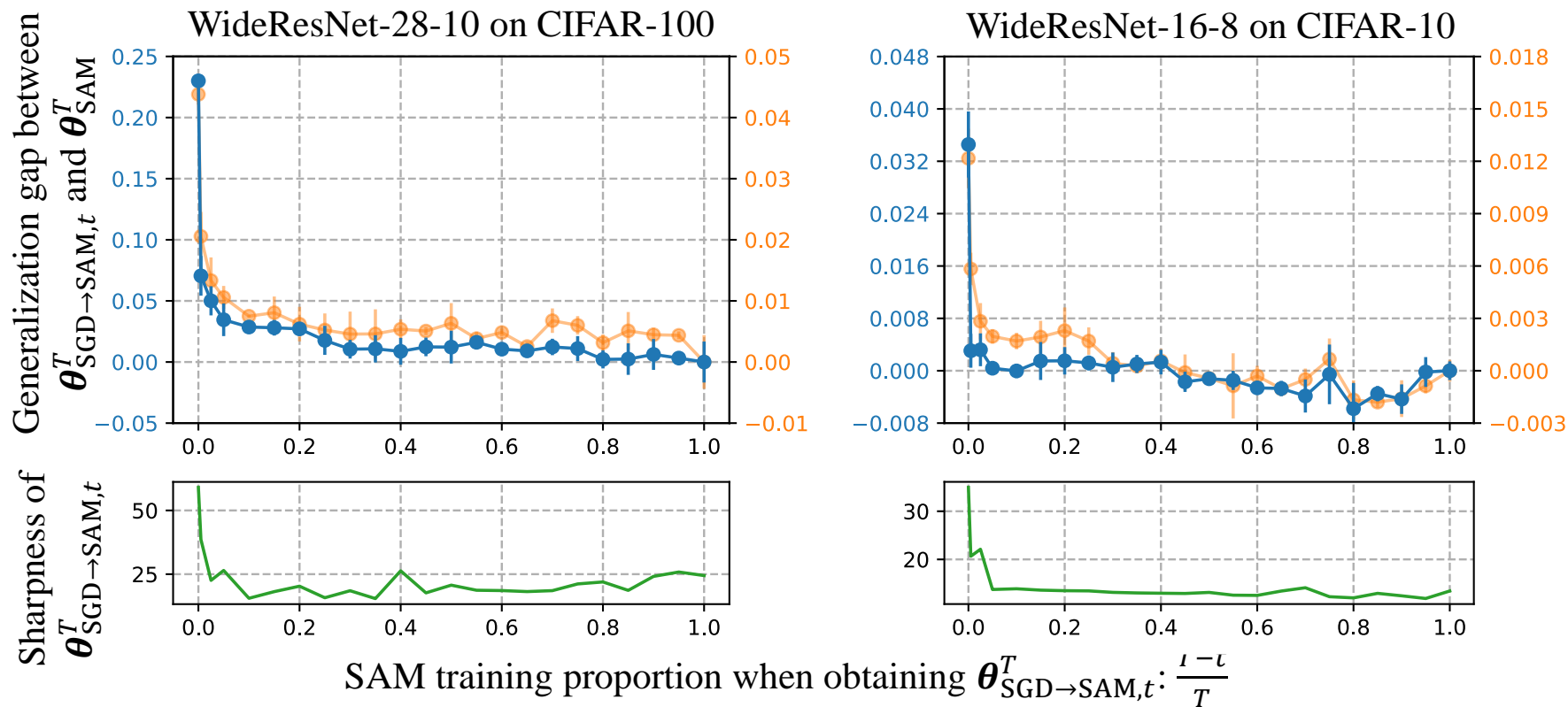


Fig. 2: **SAM operates efficiently late in training.** Blue line represents θ_{SGD}^T , while orange lines represents θ_{SAM}^T . Orange line represents $\theta_{SGD \rightarrow SAM, t}^T$, where $t = 175$.

SAM Selects Flatter Minima Late In Training



—●— $\mathcal{L}_{D_{\text{test}}}(\theta_{\text{SGD} \rightarrow \text{SAM}, t}^T) - \mathcal{L}_{D_{\text{test}}}(\theta_{\text{SAM}}^T)$
—●— $\text{Err}_{D_{\text{test}}}(\theta_{\text{SGD} \rightarrow \text{SAM}, t}^T) - \text{Err}_{D_{\text{test}}}(\theta_{\text{SAM}, t}^T)$
— $\|H(\theta_{\text{SGD} \rightarrow \text{SAM}, t}^T)\|_2$

Fig. 3: **Few epochs of SAM substantially improves generalization/sharpness.** We vary t while keep T fixed to adjust the SAM training proportion of $\theta_{\text{SGD} \rightarrow \text{SAM}, t}^T$.

How SAM Selects Flatter Minima late In Training

A Two-Phase Picture

We identify a two-phase picture in training dynamics after switching to SAM in the late training phase. This two-phase picture is characterized into four key claims **(P1-4)**, as outlined in Tab. 1.

Phase I. (Escape)	(P1). <i>SAM rapidly escapes from the minimum found by SGD;</i> (P2). <i>However, the iterator remains within the current valley.</i>	Theorem 4.2 Proposition 4.1
Phase II. (Converge)	(P3). <i>SAM converges to a flatter minimum compared to SGD;</i> (P4). <i>The convergence rate of SAM is extremely fast.</i>	Theorem 4.1 Theorem 4.3

Tab. 1: **Overview of the two-phase picture and corresponding theoretical results.**

How SAM Selects Flatter Minima late In Training

To understand the two-phase picture, let us first use a toy but representative example.

Example 4.1.

Consider using the shallow neural network $f(u, v; x) = \tanh(v \tanh(ux))$ to fit a single data $(x = 1, y = 0)$ under the squared loss $\ell(y; y') = (y - y')^2/2$, then the loss landscape is $\mathcal{L}(u, v) = \frac{1}{2} \tanh^2(v \tanh(u))$.

Note: If dynamics occurs around the set of global minima $\mathcal{M} = \{(u, v) | v = 0\}$, then small u implies flatter minima*.

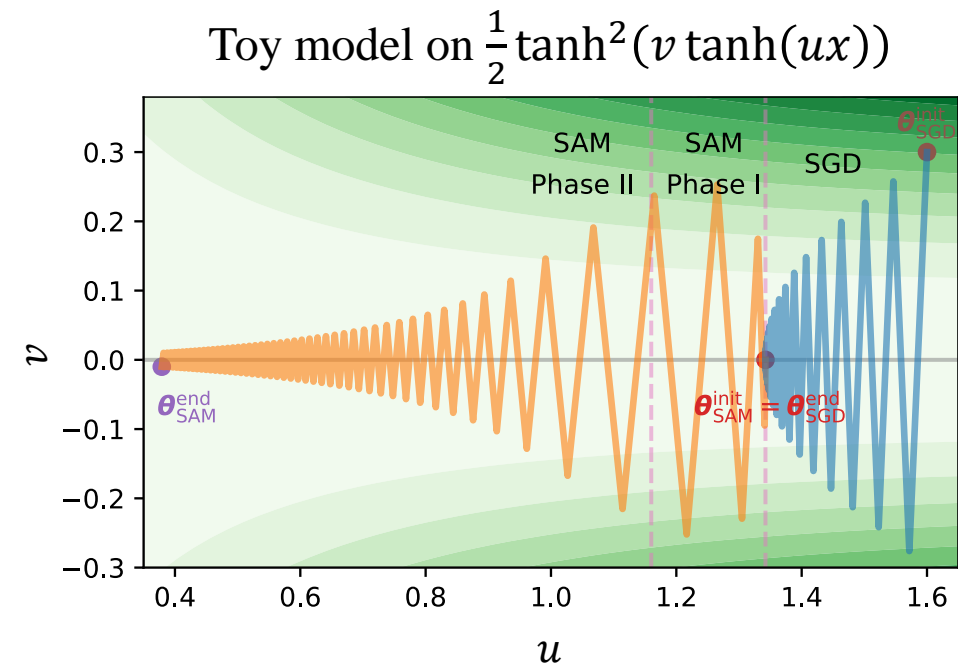


Fig. 4: **Visualization of the two-phase dynamics for Example 4.1.** The horizontal gray line represents \mathcal{M} . Blue lines trace SGD, while orange lines show SAM.

How SAM Selects Flatter Minima late In Training

Theoretical Support for P1 and P3: A Linear Stability Analysis [3].

Theorem 4.1 (P3)

Let θ^ be a global minimum that is linearly stable for SAM and suppose Assumption 4.1 (see main paper) holds, then we have $\|H(\theta^*)\|_F^2 \left(1 + \frac{\rho^2 \gamma}{B} \|H(\theta^*)\|_F^2\right) \leq \frac{B}{\eta^2 \gamma}$.*

* B denotes the mini-batch size, and $\gamma \geq 0$ is a constant defined in Assumption 4.1.

Theorem 4.1 characterizes the sharpness of the global minima selected by SAM. In Tab. 2, SAM probably selects flatter minima than SGD (P3).

	SAM (Theorem 4.1)	SGD [3]
Sharpness Bound	$\ H(\theta^*)\ _F^2 \left(1 + \frac{\rho^2 \gamma}{B} \ H(\theta^*)\ _F^2\right) \leq \frac{B}{\eta^2 \gamma}$	$\ H(\theta^*)\ _F^2 \leq \frac{B}{\eta^2 \gamma}$

Tab. 2: Comparison of the sharpness of global minima selected by SAM and SGD.

How SAM Selects Flatter Minima late In Training

Theorem 4.2 (P1)

Let $\boldsymbol{\theta}^*$ be a global minimum that is linearly stable for SAM and suppose Assumption 4.1 (see main paper) holds. If $\|H(\boldsymbol{\theta}^*)\|_F^2 \left(1 + \frac{\rho^2 \gamma}{B} \|H(\boldsymbol{\theta}^*)\|_F^2\right) > \frac{B}{\eta^2 \gamma}$, then $\boldsymbol{\theta}^*$ is linearly non-stable for SAM and $\mathbb{E}[\mathcal{L}(\boldsymbol{\theta}^t)] \geq C^t \mathbb{E}[\mathcal{L}(\boldsymbol{\theta}^0)]$ holds for all $t > 0$ with $C > 1$.

Theorem 4.2 characterizes the necessary condition of a linearly stable minimum for SAM. As SGD minimum cannot meet the stability condition of SAM, SAM will escape from the minimum found by SGD exponentially fast (P1).

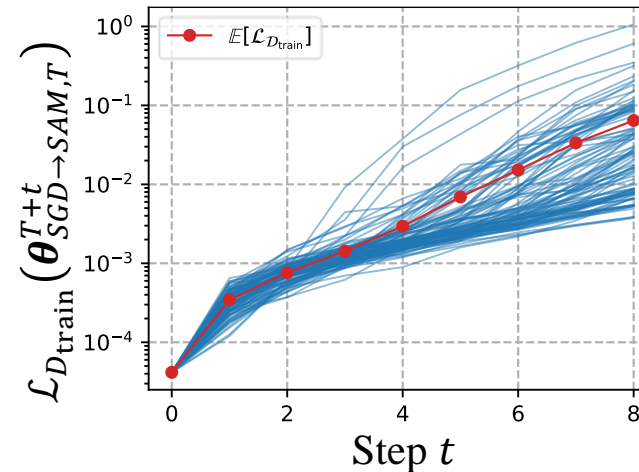


Fig. 5: **The exponentially fast escape from minima found by SGD.** Train loss $\mathcal{L}_{D_{\text{train}}}(\boldsymbol{\theta}_{\text{SGD} \rightarrow \text{SAM}, T}^{T+t})$ vs. step t .

How SAM Selects Flatter Minima late In Training

Theoretical Support for P2: Beyond Local Analysis [4].

Proposition 4.1 (P2)

Under Definition 4.2 (see main paper), assume the landscape is sub-quadratic in the valley $V = [-2b, 2b]$. Then, $\forall \eta, \rho$ s.t. $\eta < \min_{z \in V} b/|\mathcal{L}'(z)|$, $\rho \leq \min\{\frac{1}{a}, \eta \min_{0 < |z| < b} |\mathcal{L}'(2z)/\mathcal{L}'(z)|\}$, and $\theta_0 \in (-b, b)$, the full-batch SAM will remain within V , i.e., $\theta_t \in V, \forall t \in \mathbb{N}$.

Proposition 4.1 supports our key claim **P2** that SAM remains within the current valley during escape.

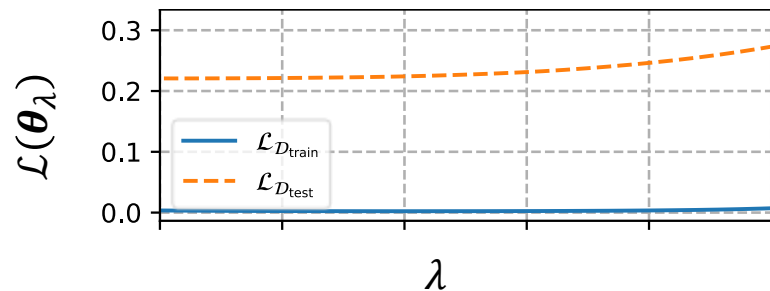


Fig. 6: **SAM converges to a flatter minimum within the same valley as the one found by SGD.** The loss of the interpolated model $\mathcal{L}_D(\theta_\lambda)$ vs. interpolation coefficient λ . Here, $\theta_\lambda = (1 - \lambda)\theta_{\text{SGD} \rightarrow \text{SAM}}^{\text{end}} + \lambda\theta_{\text{SGD}}^{\text{end}}$.

How SAM Selects Flatter Minima late In Training

Theoretical Support for P4: Convergence Analysis.

Theorem 4.3 (P4)

Under Assumption 4.2 and 4.3 (main paper), let $\{\boldsymbol{\theta}^t\}_t$ be the weights found by SAM. If $\eta \leq \min\{\frac{1}{L}, \frac{\mu B}{2L\sigma^2}\}$ and $\rho \leq \min\{\frac{1}{L}, \frac{\mu B}{4L\sigma^2}, \frac{\eta\mu^2}{24L^2}\}$, then we have $\mathbb{E}[\mathcal{L}(\boldsymbol{\theta}^t)] \leq \left(1 - \frac{\eta\mu}{2}\right)^t \mathcal{L}(\boldsymbol{\theta}^0), \forall t \in \mathbb{N}$.

Theorem 4.3 supports our key claim **P4** that the convergence rate of stochastic SAM is significantly fast, and notably faster than the previous result on SAM's convergence rate [5].

Is SAM Still Necessary in Early Phase?

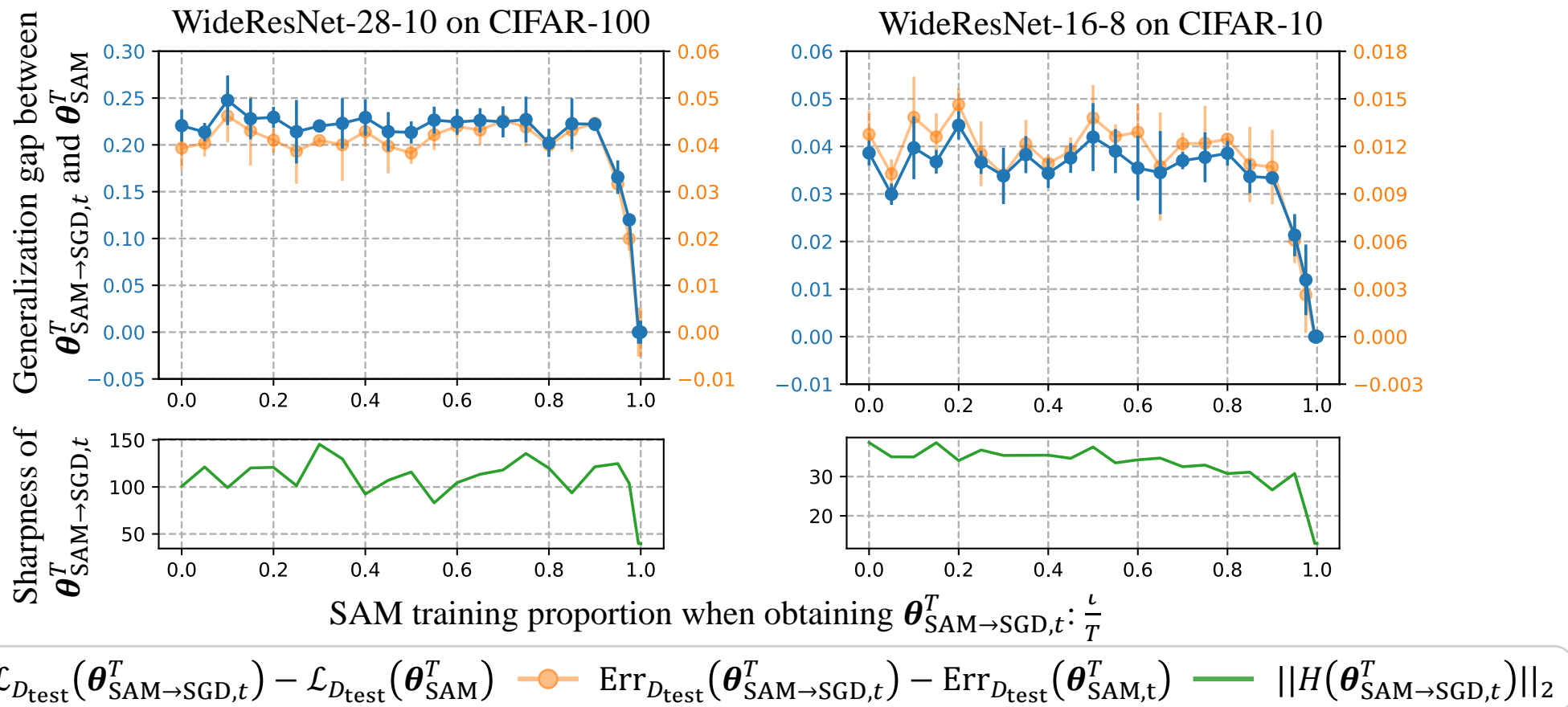


Fig. 7: **Early-phase SAM marginally improves generalization/sharpness.** We vary t while keep T fixed to adjust the SAM training proportion of $\theta_{\text{SAM} \rightarrow \text{SGD}, t}^T$.

Extend Findings to Adversarial Training

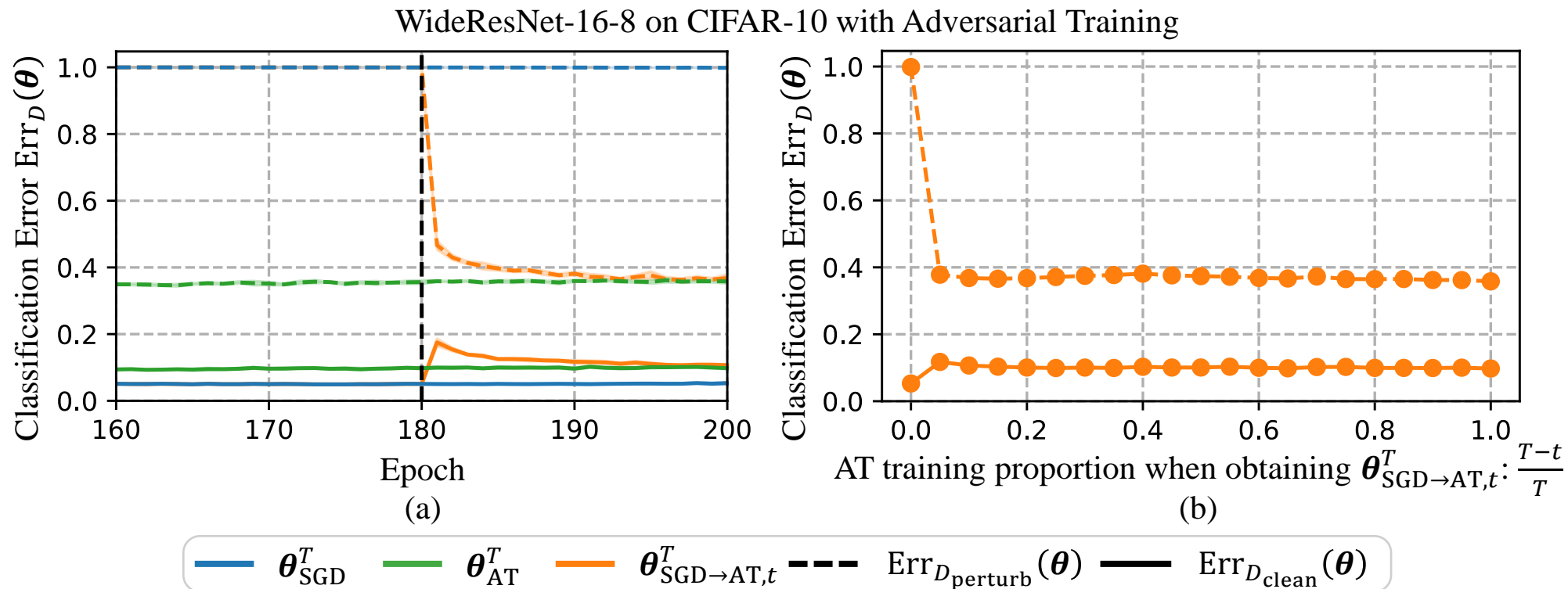


Fig. 8: **AT improves robustness efficiently even when applied only during the final few epochs of training.** (a) Robust/natural error vs. training epochs for model trained with different strategies. (b) Robust/natural error of $\theta_{\text{SGD} \rightarrow \text{AT}, t}^T$ vs. the proportion of AT epochs $\frac{T-t}{T}$.

Thank you!

Q&A