

# SHARPNESS-AWARE MINIMIZATION IN LOGIT SPACE EFFICIENTLY ENHANCES DIRECT PREFERENCE OPTIMIZATION

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Direct Preference Optimization (DPO) has emerged as a popular algorithm for aligning pretrained large language models with human preferences, owing to its simplicity and training stability. However, DPO suffers from the recently identified *squeezing effect* (also known as *likelihood displacement*), where the probability of preferred responses decreases unintentionally during training. To understand and mitigate this phenomenon, we develop a theoretical framework that models the coordinate-wise dynamics in logit space. Our analysis reveals that **negative-gradient updates** cause residuals to expand rapidly along high-curvature directions, which underlies the squeezing effect, whereas Sharpness-Aware Minimization (SAM) can suppress this behavior through its curvature-regularization effect. Building on this insight, we investigate *logits-SAM*, a computationally efficient variant that perturbs only the output layer with negligible overhead. Extensive experiments on Pythia-2.8B, Mistral-7B, and Gemma-2B-IT across multiple datasets and benchmarks demonstrate that logits-SAM consistently improves the effectiveness of DPO and integrates seamlessly with other DPO variants.

## 1 INTRODUCTION

Reinforcement learning from human feedback (RLHF) (Christiano et al., 2017; Stiennon et al., 2020; Ouyang et al., 2022) is a crucial technique for aligning pretrained large language models (LLMs) with human preferences to ensure helpfulness, harmlessness and safety (Bai et al., 2022; Dai et al., 2023). Its pipeline typically comprises three stages: supervised fine-tuning (SFT), reward modeling, and policy optimization. Classical policy optimization methods such as Proximal Policy Optimization (PPO) (Schulman et al., 2017), while widely used for their effectiveness, depend heavily on the quality of the learned reward model, rendering training complex and often unstable. Direct Preference Optimization (DPO) (Rafailov et al., 2024b) is a recently proposed and promising offline alternative that, by reparameterizing the implicit reward and optimizing a closed-form objective on preference data, trains the policy directly without explicitly fitting a reward model. DPO has gained traction due to its algorithmic simplicity and training stability.

Despite DPO and its many variants demonstrating state-of-the-art performance across a range of tasks, several potential issues remain. A particularly important one is the recently identified *squeezing effect* (Ren & Sutherland, 2024) (also known as *likelihood displacement* (Razin et al., 2024)), which describes an unintended decrease in the generation probability of preferred responses during DPO training, contrary to the intended goal of increasing it embodied in the DPO objective. This phenomenon can lead to performance degradation, reduced safety, and even alignment failure (Pal et al., 2024; Yuan et al., 2024; Rafailov et al., 2024a; Tajwar et al., 2024; Pang et al., 2024).

To understand the mechanism behind the squeezing effect and to identify an effective remedy, we develop a theoretical framework that elucidates the learning dynamics in both the parameter space and the logit space. Our analysis shows that gradient updates with a negative learning rate, **which are algorithmically equivalent to the negative-gradient updates induced by the negative objective associated with rejected answers in DPO**, cause the residual vector to expand rapidly along high-curvature directions, namely along the eigenvectors associated with large eigenvalues of the Hessian,

which is the source of the squeezing effect. This raises a natural question: *can curvature-aware training mitigate this unintended drift?*

We investigate *Sharpness-Aware Minimization* (SAM) (Foret et al., 2021), a bilevel optimization method widely used in supervised learning, and establish its dynamics in both the parameter and logit spaces. Our theory demonstrates that SAM effectively alleviates the squeezing effect through its intrinsic curvature regularization. Guided by these insights, we advocate using *logits-SAM* for DPO training, a computationally efficient variant of SAM that perturbs only the output-layer parameters. Although logits-SAM has been mentioned merely as a byproduct in prior work (Baek et al., 2024; Singh et al., 2025) and often overlooked, our study turns this neglected variant into a practically useful and effective technique by integrating it into DPO, where it efficiently mitigates the squeezing effect and consistently improves performance. To the best of our knowledge, this is the first work to analyze and apply SAM in the context of DPO.

**Contributions.** Our contributions are summarized as follows:

- We develop a theoretical framework that connects the parameter space and the logit space through geometric properties, enabling a unified analysis of learning dynamics in both domains. This framework yields unified dynamical equations for gradient descent (GD) and SAM that precisely track coordinate-wise evolution with controlled error terms.
- Our analysis identifies the root cause of the squeezing effect: under a negative learning rate, residuals expand rapidly along high-curvature directions. We rigorously show that SAM, through its intrinsic curvature regularization, effectively alleviates this phenomenon.
- Bridging theory and practice, we implement an efficient variant, *logits-SAM*, which perturbs only the output-layer parameters. Unlike vanilla SAM, it incurs virtually no additional overhead. Experiments on Pythia-2.8B and Mistral-7B across multiple datasets and benchmarks validate its effectiveness, demonstrating consistent performance gains for DPO and its variants.

## 2 PRELIMINARIES

### 2.1 PREFERENCE OPTIMIZATION

**SFT-RLHF pipeline.** Classical RLHF alignment proceeds in three phases: (i) *supervised fine-tuning* of a base policy on instruction-following data; (ii) *reward modeling* by fitting a scalar reward function on pairwise human preferences; and (iii) *policy optimization* to maximize the learned reward under a KL regularizer toward a reference policy.

**DPO reparameterization.** DPO (Rafailov et al., 2024b) bypasses training an explicit reward model by expressing an *implicit* reward for a policy  $\pi_\theta$  as a log-likelihood ratio to a fixed reference policy  $\pi_{\text{ref}}$  (typically the SFT model):

$$r_\theta(\mathbf{x}, \mathbf{y}) = \beta \log \frac{\pi_\theta(\mathbf{y} | \mathbf{x})}{\pi_{\text{ref}}(\mathbf{y} | \mathbf{x})} + \beta \log Z(\mathbf{x}), \quad (1)$$

where  $\beta > 0$  is a temperature and  $Z(\mathbf{x})$  is a partition term independent of  $\theta$ . Combining equation 1 with the Bradley–Terry preference model (Bradley & Terry, 1952)  $p(\mathbf{y}^+ \succ \mathbf{y}^- | \mathbf{x}) = \sigma(r_\theta(\mathbf{x}, \mathbf{y}^+) - r_\theta(\mathbf{x}, \mathbf{y}^-))$  yields the standard DPO objective, optimized over a dataset  $\mathcal{D} = \{(\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-)\}$  of preferred/dispreferred pairs:

$$\mathcal{L}_{\text{DPO}}(\pi_\theta; \pi_{\text{ref}}) = -\mathbb{E}_{(\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-) \sim \mathcal{D}} \left[ \log \sigma \left( \beta \log \frac{\pi_\theta(\mathbf{y}^+ | \mathbf{x})}{\pi_{\text{ref}}(\mathbf{y}^+ | \mathbf{x})} - \beta \log \frac{\pi_\theta(\mathbf{y}^- | \mathbf{x})}{\pi_{\text{ref}}(\mathbf{y}^- | \mathbf{x})} \right) \right], \quad (2)$$

where  $\sigma(\cdot)$  is the logistic function.

### 2.2 SHARPNESS-AWARE MINIMIZATION

SAM regularizes training by explicitly penalizing *parameter-space sharpness*: it chooses parameters that minimize the worst-case loss within an  $\ell_2$  ball of radius  $\rho$  around  $\theta$ . Concretely, for supervised learning with examples  $(\mathbf{x}, \mathbf{y}) \sim \mathcal{D}$  and per-example loss  $f(\theta; \mathbf{x}, \mathbf{y})$ , the SAM objective is

$$\min_{\theta} \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \mathcal{D}} \left[ \max_{\|\epsilon\|_2 \leq \rho} f(\theta + \epsilon; \mathbf{x}, \mathbf{y}) \right]. \quad (3)$$

This formulation can be interpreted as a form of *curvature regularization*: by seeking minimizers whose neighborhoods exhibit consistently low loss, SAM favors flatter minima that often correlate with improved generalization. In practice, the inner maximization is approximated to first order by the perturbation  $\epsilon^*(\theta) = \rho \nabla_{\theta} f(\theta; \mathbf{x}, \mathbf{y}) / \|\nabla_{\theta} f(\theta; \mathbf{x}, \mathbf{y})\|$ , and one takes a descent step using the gradient at the perturbed point,  $\nabla_{\theta} f(\theta + \epsilon^*; \mathbf{x}, \mathbf{y})$ .

### 3 LEARNING DYNAMICS IN LOGIT SPACE

#### 3.1 SETTING

We adopt the same theoretical setting as in Ren & Sutherland (2024), namely multiclass logistic classification, where the features of the samples are fixed (also referred to as the kernel regime (Malladi et al., 2023)), and the learning rate can be either positive or negative, corresponding respectively to the objectives of  $\mathbf{y}^+$  and  $\mathbf{y}^-$  in DPO (for analytical convenience, we allow negative learning rates rather than explicitly modeling negative objectives). Prior work (Ren & Sutherland, 2024) has shown that the negative-gradient dynamics in DPO, including the characteristic squeezing effect, can be faithfully reproduced within this simplified setting. Several phenomena observed in the multi-class logistic regression abstraction also emerge empirically during real LLM fine-tuning. These findings suggest that analyzing DPO through this framework offers a theoretically tractable and practically relevant perspective on the learning behavior of LLMs.

Let  $\mathbf{x}$  be a training example with one-hot label  $\mathbf{y} \in \{0, 1\}^V$ ,  $\mathbf{1}^\top \mathbf{y} = 1$ . In the fixed-feature (kernel) regime,  $\phi(\mathbf{x}) \in \mathbb{R}^d$  are fixed and

$$\mathbf{z}^t = \mathbf{W}^t \phi(\mathbf{x}) \in \mathbb{R}^V, \quad \mathbf{p}^t = \text{softmax}(\mathbf{z}^t), \quad f(\mathbf{z}^t, \mathbf{y}) = - \sum_{k=1}^V \mathbf{y}_k \log p_k^t,$$

where  $\mathbf{W}^t \in \mathbb{R}^{V \times d}$  are trainable parameters,  $\mathbf{z}^t$  are the logits. For notational convenience, we write  $\phi(\mathbf{x})$  as  $\phi$ . We use  $\|\cdot\|$  to denote the  $\ell_2$  norm for vectors and the Frobenius norm for matrices. We use  $\otimes$  to denote the Kronecker product.

We denote the parameter Hessian by  $\mathbf{H}_{\mathbf{W}}^t := \nabla_{\mathbf{W}}^2 f(\mathbf{z}^t, \mathbf{y}) \in \mathbb{R}^{Vd \times Vd}$ , and  $\mu := \|\phi\|^2$ . In logit space, we denote the logit gradient by  $\mathbf{g}^t := \nabla_{\mathbf{z}} f(\mathbf{z}^t, \mathbf{y}) = \mathbf{p}^t - \mathbf{y} \in \mathbb{R}^V$ , and denote the logit Hessian by  $\mathbf{H}_{\mathbf{z}}^t := \nabla_{\mathbf{z}}^2 f(\mathbf{z}^t, \mathbf{y}) \in \mathbb{R}^{V \times V}$ .

#### 3.2 THEORY

The theoretical results of Ren & Sutherland (2024) demonstrate that the *squeezing effect* arises from the objective with a negative learning rate. Specifically, they prove that the probability of the ground-truth label necessarily decreases, while the probability of the model’s most confident incorrect class necessarily increases. In this work, we provide a finer-grained analysis of the learning dynamics under this setting. We establish a unified modeling framework for the residuals of all classes and derive the linear convergence rate up to higher-order remainders. Furthermore, we apply our framework to prior analyses and further establish a rigorous conclusion that SAM can effectively mitigate the squeezing effect.

For GD, first-order derivatives are sufficient to characterize its dynamics. However, the intrinsic curvature regularization effect of SAM motivates us to further investigate the geometric structure of the parameter space through the Hessian matrix. To this end, we develop a theoretical framework that connects the geometry of the parameter space and the logit space, via the link between the parameter Hessian and the logit Hessian.

**Proposition 3.1** (Geometry of the logit space; simplified version of Proposition A.1). *In coordinates,  $\mathbf{H}_{\mathbf{W}} = \mathbf{H}_{\mathbf{z}} \otimes (\phi\phi^\top)$ . Thus, if  $\phi \neq \mathbf{0}$ , then  $\text{rank}(\mathbf{H}_{\mathbf{W}}) = \text{rank}(\mathbf{H}_{\mathbf{z}})$ . Moreover, the second-order effect of any parameter perturbation depends only on the induced logits perturbation  $T_{\phi}(\Delta \mathbf{W}) := \Delta \mathbf{W} \phi$ .*

This proposition establishes that all second-order effects in the parameter space, whose Hessian  $\mathbf{H}_{\mathbf{W}}$  lies in  $\mathbb{R}^{Vd \times Vd}$ , can be equivalently studied through the logit Hessian  $\mathbf{H}_{\mathbf{z}}$  in  $\mathbb{R}^{V \times V}$ , thereby greatly simplifying the analysis of second-order dynamics. Next, we establish a unified framework

to track the SAM dynamics in both the parameter space and the logit space, thanks to their favorable geometric structures. Unlike prior work, our framework can simultaneously trace the evolution of all coordinates of the parameters, logits, and residuals, while providing precise control over the error terms.

**Theorem 3.2** (SAM dynamics in parameter and logit space; informal version of Theorem A.2). *Assume that we conduct the SAM update for  $\mathbf{W}$ . Under mild assumptions, there exists a constant  $C > 0$  such that the following expansions hold with  $O(\eta^2)$  remainders:*

$$\begin{aligned}
 (\text{parameters}) \quad \mathbf{W}^{t+1} &= \mathbf{W}^t - \eta \left( \mathbf{g}^t \phi^\top + \underbrace{\tilde{\rho}^t \mathbf{H}_z^t \mathbf{g}^t \phi^\top}_{\text{SAM's correction}} \right) + \mathbf{R}_W^t, \quad \|\mathbf{R}_W^t\| \leq C \eta^2, \\
 (\text{logits}) \quad \mathbf{z}^{t+1} &= \mathbf{z}^t - \eta \mu \left( \mathbf{g}^t + \underbrace{\tilde{\rho}^t \mathbf{H}_z^t \mathbf{g}^t}_{\text{SAM's correction}} \right) + \mathbf{r}_z^t, \quad \|\mathbf{r}_z^t\| \leq C \eta^2, \\
 (\text{residuals}) \quad \mathbf{g}^{t+1} &= \mathbf{p}^{t+1} - \mathbf{y} = \left( \mathbf{I} - \eta \mu \mathbf{H}_z^t - \underbrace{\eta \mu \tilde{\rho}^t (\mathbf{H}_z^t)^2}_{\text{SAM's correction}} \right) (\mathbf{p}^t - \mathbf{y}) + \mathbf{r}_g^t, \quad \|\mathbf{r}_g^t\| \leq C \eta^2,
 \end{aligned}$$

where  $\tilde{\rho}^t := \rho \sqrt{\mu} / \|\mathbf{g}^t\|$  is the equivalent perturbation coefficient.

It is worth noting that when  $\rho = 0$ , the dynamics reduce to the GD dynamics. This theorem, viewed through the lens of the logit Hessian, provides a precise theory for characterizing GD and SAM dynamics across spaces. In both parameter and logit space, GD amounts to scaling by the logit gradient, whereas SAM introduces an additional  $\mathbf{H}_z$  correction term that can be regarded as a preconditioning matrix. Moreover, the updates of the residual vector under GD and SAM are both preconditioned by  $\mathbf{H}_z$  (and, for SAM, by  $(\mathbf{H}_z)^2$ ). This implies that if we choose the eigenvectors of the logit Hessian as a basis, the curvature coupling effects of both the first-order and second-order terms can be unified. To formalize this intuition, we show that  $\mathbf{g}$  lies precisely in the column space of  $\mathbf{H}_z$ , thus we can select the nonzero eigenvectors of  $\mathbf{H}_z$  as a basis to obtain the coordinate representation of  $\mathbf{g}$ .

**Proposition 3.3.**  $\mathbf{H}_z$  is symmetric positive semidefinite with  $\ker(\mathbf{H}_z) = \text{span}\{\mathbf{1}\}$  and  $\text{rank}(\mathbf{H}_z) = V - 1$ . Moreover, for the residual  $\mathbf{g}$  we have  $\mathbf{1}^\top \mathbf{g} = 0$ , hence  $\mathbf{g} \in \mathbf{1}^\perp = \text{range}(\mathbf{H}_z)$ ; in particular, given any eigenbasis of  $\mathbf{H}_z$  restricted to  $\mathbf{1}^\perp$ ,  $\mathbf{g}$  admits a unique coordinate representation in that basis.

**Corollary 3.4** (Modal dynamics in the eigenbasis of  $\mathbf{H}_z^t$ ). *Under the same assumptions as Theorem 3.2. For each  $t$ , let the spectral decomposition of the symmetric positive-semidefinite matrix  $\mathbf{H}_z^t$  be*

$$\mathbf{H}_z^t = \sum_{k=1}^{V-1} \lambda_k^t \mathbf{v}_k^t (\mathbf{v}_k^t)^\top,$$

where  $\lambda_k^t > 0$ ,  $(\mathbf{v}_k^t)^\top \mathbf{v}_\ell^t = \delta_{k\ell}$  are the non-zero eigenvalues and eigenvectors. Define the modal coefficients of the residual  $\mathbf{g}^t = \mathbf{p}^t - \mathbf{y}$  by

$$\mathbf{e}_k^t := (\mathbf{v}_k^t)^\top \mathbf{g}^t, \quad \mathbf{e}_k^{t+1} := (\mathbf{v}_k^t)^\top \mathbf{g}^{t+1}, \quad k = 1, \dots, V-1. \quad (4)$$

Then there exists a constant  $C > 0$  such that for all nonzero modes  $k \geq 1$ ,

$$\mathbf{e}_k^{t+1} = \left( 1 - \eta \mu \left[ \lambda_k^t + \underbrace{\tilde{\rho}^t (\lambda_k^t)^2}_{\text{SAM's correction}} \right] \right) \mathbf{e}_k^t + \mathbf{r}_k^t, \quad |\mathbf{r}_k^t| \leq C \eta^2. \quad (5)$$

Proofs are deferred to Appendix A. The corollary diagonalizes the vector dynamics into coordinate-wise scalars in the eigenbasis of  $\mathbf{H}_z$ , making SAM's effect transparent. We now characterize the additional SAM correction in two regimes.

**Case 1: Positive  $\eta$** , corresponding to the  $\mathbf{y}^+$  objective in DPO. In this case, GD induces a stronger contraction of the residual  $\mathbf{g}$  along the high-curvature directions, i.e., those associated with large eigenvalues of  $\mathbf{H}_z$ . The additional correction term introduced by SAM has the same sign as that of GD, thereby amplifying this effect. **Case 2: Negative  $\eta$** , corresponding to the  $\mathbf{y}^-$  objective in DPO. Here, GD causes the residual  $\mathbf{g}$  to expand more rapidly along high-curvature directions. Furthermore, standard SAM with positive  $\rho$  exacerbates this phenomenon, causing the residual to expand

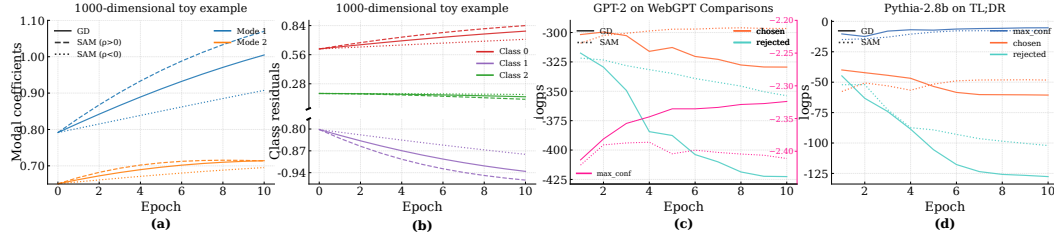


Figure 1: Training dynamics under different settings. (a–b) 1000-dimensional toy example with three classes, trained with a negative learning rate under GD, SAM ( $\rho > 0$ ), and SAM ( $\rho < 0$ ). Panel (a) shows the modal coefficients, and panel (b) shows the class residuals. (c) Real-data experiment on WebGPT Comparisons with GPT-2, comparing GD and SAM: the panel reports the log-probabilities of the chosen responses, the rejected responses, and `max_conf`, which denotes the model’s most confident response. (d) Real-data experiment on the TL;DR dataset with Pythia-2.8B, showing the same three curves (chosen, rejected, and `max_conf`).

even faster along high-curvature directions compared to GD. By contrast, choosing a negative  $\rho$  counteracts this expansion.

Next, we extend our theoretical framework to the result of Ren & Sutherland (2024), which introduced the squeezing effect. For consistency with their notation, we let  $y$  denote the ground-truth class index (with one-hot label  $\mathbf{y} = \mathbf{e}_y$ ).

**Lemma 3.5** (One-step confidence ratios under GD, Lemma 1 of Ren & Sutherland (2024)). *For each class  $i \in [V]$ , define the one-step confidence ratio  $\alpha_i := p_i^{t+1}/p_i^t$ . Consider the objective with a negative learning rate  $\eta < 0$ , and denote its ground-truth label by  $y^*$ . Let  $y^- = \arg \max_{j \neq y^*} p_j^t$  be the most confident incorrect class. Then*

$$\alpha_{y^*}^{\text{GD}} > 1, \quad \alpha_{y^-}^{\text{GD}} < 1.$$

Lemma 3.5 formalizes the squeezing effect under GD with a negative learning rate: the probability of the most confident incorrect class increases, while that of the ground-truth class decreases. Within our framework, we next analyze the ratio of these two probabilities after a one-step SAM update.

**Corollary 3.6** (One-step confidence ratios under SAM, informal version of Corollary A.5). *Under the same assumptions as Theorem 3.2, assume that  $\eta\rho > 0$ . Then, for sufficiently small step size  $|\eta|$ , the following inequalities hold:*

$$\alpha_{y^*}^{\text{SAM}} \leq \alpha_{y^*}^{\text{GD}}, \quad \alpha_{y^-}^{\text{SAM}} \geq \alpha_{y^-}^{\text{GD}}. \quad (6)$$

Here  $y \in \{y^+, y^-\}$  denotes the ground-truth label corresponding to the positive or negative learning rate, respectively. Moreover, the inequalities in equation 6 are strict whenever  $p_{y^*}^t \in (0, 1)$  and  $p_y^t \leq \frac{1}{2}$ .

The proof is deferred to Appendix A. Corollary 3.6 and Lemma 3.5 together imply that, when  $\eta < 0$ , using SAM with a negative  $\rho < 0$  moderates the growth of the most confident incorrect class and slows the decay of the ground-truth class, thereby preventing excessive expansion and premature collapse. Our analysis thus reveals a key, albeit somewhat counterintuitive, fact: for negative  $\eta$ , one should choose a *negative*  $\rho$  (interpreted as a perturbation along the gradient descent direction), which effectively alleviates the squeezing effect.

We empirically validate our theoretical findings using a 1000-dimensional toy example with three classes. Specifically, we first train for 10 epochs using class 0 as the label for initialization, mimicking the SFT process, and then switch to class 1 as the label while continuing training with a negative learning rate. As shown in Figure 1, this setup faithfully reproduces the squeezing effect observed in prior work (Ren & Sutherland, 2024): both modal coefficients expand rapidly, the probabilities of class 1 and class 2 decrease, and only the probability of class 0, the model’s most confident incorrect prediction, increases. Moreover, SAM with positive  $\rho$  exacerbates this effect, whereas SAM with negative  $\rho$  hinders this trend, exactly as predicted by our theory.

Additionally, Corollary 3.6 shows that for  $\eta > 0$ , SAM with  $\rho > 0$  likewise mitigates the squeezing effect: the contraction of  $y^*$  is accelerated, while the growth of the ground-truth  $y^+$  is enhanced.

Taken together, these results establish a simple rule: during training, choosing  $\rho$  with the *same sign* as the learning rate alleviates the squeezing effect—specifically, it restrains the growth of  $y^*$  and promotes (or reduces the suppression of)  $y^+$  and  $y^-$ . To validate this idea, we track the probability dynamics of the chosen responses, the rejected responses, [and the model’s most confident responses](#). We conduct two real-world experiments: fine-tuning a GPT-2 (Radford et al., 2019) model on a subset of the WebGPT Comparisons dataset (Nakano et al., 2022) (Figure 1c), and fine-tuning a Pythia-2.8B model (Biderman et al., 2023) on a subset of TL;DR dataset (Stiennon et al., 2020) (Figure 1d). In both settings, we observe the same trend: SAM increases the probability of the chosen responses, slows the decrease in the probability of the rejected responses, [and prevents the probability of the most confident responses from growing](#). These findings are fully consistent with our theoretical predictions.

### 3.3 FROM THEORY TO PRACTICE

In practice, an important challenge in applying SAM to DPO is that it requires an additional forward and backward pass, thereby nearly doubling the computational cost. However, our dynamical analysis shows that curvature regularization can still be achieved even when the perturbation is applied solely in the logit space (with an appropriate choice of the sign of  $\rho$ ), which also alleviates the squeezing effect. Motivated by this observation, we suggest using a computationally efficient SAM variant that perturbs only in the last layer, called *logits-SAM*, to improve the effectiveness and robustness of DPO. Its objective can be formulated as follows:

$$\mathcal{L}_{\text{DPO}}^{\text{logits-SAM}}(\mathbf{W}, \theta; \mathbf{x}, \mathbf{y}^+, \mathbf{y}^-) = \mathcal{L}_{\text{DPO}}\left(\mathbf{W} + \rho \frac{\nabla_{\mathbf{W}} \mathcal{L}_{\text{DPO}}(\mathbf{W}, \theta; \mathbf{x}, \mathbf{y}^+, \mathbf{y}^-)}{\|\nabla_{\mathbf{W}} \mathcal{L}_{\text{DPO}}(\mathbf{W}, \theta; \mathbf{x}, \mathbf{y}^+, \mathbf{y}^-)\|}, \theta; \mathbf{x}, \mathbf{y}^+, \mathbf{y}^-\right).$$

where  $\mathbf{W}$  denotes the parameters in the output layer, and  $\theta$  denotes the parameters except  $\mathbf{W}$ .

**Implementation.** Unlike our theoretical setting, common DPO implementations<sup>12</sup> typically encode the  $\mathbf{y}^-$  objective as negative while using a single positive learning rate, rather than assigning positive and negative rates to  $\mathbf{y}^+$  and  $\mathbf{y}^-$ , respectively. Accordingly, we adopt this convention in our logits-SAM implementation. Our dynamical analysis further indicates that  $\rho$  should share the sign of the learning rate; hence we consistently use a positive  $\rho$ . [In Appendix B, we provide the full update formulas and a derivation establishing the equivalence between the theoretical and practical settings, summarized in Table 5.](#)

*Remark.* This choice does not render our analysis of the negative learning rate redundant. For first-order methods such as GD, using a negative objective with a positive learning rate is equivalent to using a positive objective with a negative learning rate. Therefore, our analysis applies fully to the case of negative objectives.

The implementation pseudocode can be found in Algorithm 1 of Appendix B. We compute the perturbation manually using the hidden states from the penultimate layer and the parameters of the final layer, requiring only a single full forward-backward pass instead of the two full passes required in standard SAM. Since the parameters of the final layer typically constitute only a small fraction of all trainable parameters (e.g., 4.64% in Pythia-2.8B and 1.81% in Mistral-7B), the additional training overhead introduced by logits-SAM is negligible. A detailed comparison of wall-clock time and peak memory usage is provided in Section 4.3.

## 4 EXPERIMENTS

### 4.1 EXPERIMENTAL SETUP

**Datasets.** We conduct DPO training on three widely used datasets to evaluate our algorithm: Anthropic-HH (Bai et al., 2022), the Reddit TL;DR summarization dataset (Stiennon et al., 2020), and the UltraFeedback Binarized dataset (Cui et al., 2023).

<sup>1</sup><https://github.com/eric-mitchell/direct-preference-optimization>

<sup>2</sup><https://github.com/huggingface/trl>

Table 1: Evaluation results (WR %) on HH and TL;DR datasets using Pythia-2.8B. The judge is GPT-5-mini. The highest value within each method group (baseline vs. logits-SAM) is **bolded**.

Method	HH		TL;DR	
	vs SFT	vs chosen	vs SFT	vs chosen
DPO	70.52	56.35	84.21	34.78
DPO+logits-SAM	<b>72.28</b>	<b>60.51</b>	<b>89.58</b>	<b>36.57</b>
SLiC-HF	65.27	54.72	91.88	31.36
SLiC-HF+logits-SAM	<b>71.87</b>	<b>62.21</b>	<b>94.40</b>	<b>32.80</b>
CPO	66.60	58.19	90.99	39.38
CPO+logits-SAM	<b>70.24</b>	<b>59.90</b>	<b>93.29</b>	<b>45.41</b>

**Models.** Following common practice, we adopt SFT models as our base models. We use Pythia-2.8B (Biderman et al., 2023) for experiments on Anthropic-HH and Reddit TL;DR, and Mistral-7B-v0.1 (Jiang et al., 2023) for UltraFeedback. For Pythia-2.8B, we initialize from the Hugging Face open-source checkpoint<sup>3</sup>, which was SFT for one epoch on Anthropic-HH. For the TL;DR experiments, we use the checkpoint<sup>4</sup>, which was SFT for one epoch on Reddit TL;DR. For Mistral-7B-v0.1, we use the Alignment Handbook (Tunstall et al., 2023a) checkpoint Zephyr-7b<sup>5</sup> (Tunstall et al., 2023b), which was SFT for one epoch on UltraChat-200k (Ding et al., 2023).

**Evaluation.** For Pythia-2.8B, we evaluate model performance on Anthropic-HH and Reddit TL;DR by measuring win rates (WR) against both the SFT baseline and the human-preferred responses, using GPT-5-mini (version 2025-08-07) as the automatic judge. Following the DPO paper, we set the decoding temperature to 0 for HH and 1 for TL;DR. For Mistral-7B-v0.1, we conduct evaluation on three popular open-ended instruction-following benchmarks: AlpacaEval 2 (Dubois et al., 2024), Arena-Hard v0.1 (Li et al., 2024), and MT-Bench (Zheng et al., 2023). Details of each benchmark can be found in Appendix C. We adopt the default generation parameters provided by each benchmark. Specifically, we report both length-controlled win rates (LC) and raw WR for AlpacaEval 2, model WR for Arena-Hard v0.1, and averaged judge scores (1–10) for MT-Bench, all following the standard evaluation protocols, with default decoding configurations.

**Baselines.** We apply logits-SAM to DPO and two SOTA variants, SLiC-HF (Zhao et al., 2023) and CPO (Xu et al., 2024). We use AdamW optimizer (Loshchilov & Hutter, 2019) in all experiments. For Pythia-2.8B, we set batch size 64 and learning rate  $1 \times 10^{-6}$ , following the DPO paper; for Mistral-7B, we use batch size 128 and learning rate  $5 \times 10^{-7}$ , following the Alignment Handbook’s recommended settings.

**Hyperparameters.** For DPO, we adopt the recommended  $\beta$  values from the DPO paper and the Alignment Handbook, which are widely used and well tuned. For SLiC-HF and CPO, we select hyperparameters following the tuning protocol from Meng et al. (2024b). For logits-SAM, we keep all hyperparameters identical to each corresponding baseline to ensure fairness; the only additional hyperparameter is  $\rho$ , which we tune over  $\{1 \times 10^{-5}, 1 \times 10^{-4}, 1 \times 10^{-3}\}$ . Full hyperparameter settings are provided in Table 6 and Table 7 of Appendix C.

## 4.2 EXPERIMENTAL RESULTS

**Performance of summarization and dialogue generation tasks.** We present the results in Table 1. We find that logits-SAM consistently improves performance across both HH and TL;DR datasets. All three baselines (DPO, SLiC-HF, and CPO) achieve higher win rates against both SFT and chosen responses when augmented with logits-SAM. Notably, SLiC-HF shows the largest gains on HH (+6.60 pp vs SFT, +7.49 pp vs chosen), while CPO achieves strong improvements on TL;DR (+2.30 pp vs SFT, +6.03 pp vs chosen), demonstrating that logits-SAM provides stable and generalizable benefits across different optimization methods.

<sup>3</sup><https://huggingface.co/lomahony/eleuther-pythia2.8b-hh-sft>

<sup>4</sup><https://huggingface.co/trl-lib/pythia-2.8b-deduped-tldr-sft>

<sup>5</sup><https://huggingface.co/alignment-handbook/zephyr-7b-sft-full>



Table 2: Evaluation results on AlpacaEval 2 (LC and WR), Arena-Hard v0.1 (WR), and MT-Bench using Mistral-7B-v0.1. Judges are GPT-4 Turbo for AlpacaEval 2, and GPT-4.1 for Arena-Hard v0.1 and MT-Bench. The highest value within each method group (baseline vs. logits-SAM) is **bolded**.

Method	AlpacaEval 2		Arena-Hard v0.1	MT-Bench
	LC (%)	WR (%)	WR (%)	(score)
DPO	13.08	10.96	19.0	5.49
DPO+logits-SAM	<b>13.90</b>	<b>11.62</b>	<b>23.1</b>	<b>5.79</b>
SLiC-HF	8.92	8.97	19.1	5.05
SLiC-HF+logits-SAM	<b>10.63</b>	<b>9.23</b>	<b>21.1</b>	<b>5.22</b>
CPO	8.97	8.13	19.2	5.22
CPO+logits-SAM	<b>13.32</b>	<b>11.78</b>	<b>21.4</b>	<b>5.49</b>

**Performance on open-ended instruction-following benchmarks.** We present the results in Table 2. The results demonstrate that combining logits-SAM with different DPO variants consistently yields performance gains across all benchmarks. On open-ended instruction-following evaluations, logits-SAM improves both length-controlled and original win rates on AlpacaEval 2 (e.g., with CPO: +4.35 pp LC, +3.65 pp WR), increases head-to-head win rate on Arena-Hard v0.1 (e.g., with DPO: +4.1 pp WR), and provides steady gains on MT-Bench (e.g., DPO: +0.30, SLiC-HF: +0.17, CPO: +0.27). These findings indicate that logits-SAM is a generally effective and robust enhancement across diverse evaluation settings.

#### 4.3 ADDITIONAL ANALYSIS

**Computational overhead.** Compared to vanilla SAM, logits-SAM minimizes additional computational overhead. We report wall-clock training time and peak memory on Pythia-2.8B trained on the Reddit TL;DR dataset (Figure 2), using data-parallel training (DDP) across two NVIDIA A100 GPUs with a per-device batch size of 4. The results show that logits-SAM adds only  $\sim 2\text{--}3\%$  extra time, with negligible peak-memory overhead. By contrast, vanilla SAM is practically infeasible for Pythia-2.8B on A100s with DDP: it nearly doubles the step time (due to an extra full forward-backward pass) and requires a perturbation buffer comparable to the model size (for billion-parameter models, this entails more than 10 GB of additional GPU memory), which leads to out-of-memory even with batch size 1. These observations highlight the clear computational cost advantage of logits-SAM.

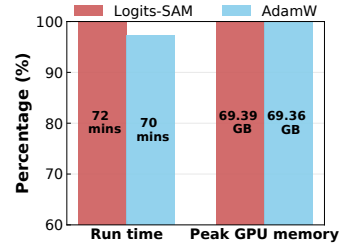


Figure 2: Efficiency comparison.

**Sensitivity analysis.** We present a sensitivity analysis of the additional hyperparameter  $\rho$  for logits-SAM in Table 3. The results indicate that, within a reasonable range of  $\rho$ , performance is typically improved, whereas further enlarging  $\rho$  leads to a marked degradation. Notably, unlike original SAM, logits-SAM perturbs only the output layer, so the appropriate scale of  $\rho$  is much smaller than the range (0.01–0.5) recommended in the SAM paper. We recommend starting the search for logits-SAM’s  $\rho$  at  $10^{-5}$  or  $10^{-4}$  and, if resources permit, performing a finer sweep in this neighborhood.

Table 3: Performance on HH and TL;DR datasets under different  $\rho$  values. Each entry reports win rate vs SFT (left) and vs chosen (right).

Dataset	$\rho = 0$ (AdamW)	$\rho = 10^{-5}$	$\rho = 10^{-4}$	$\rho = 10^{-3}$	$\rho = 10^{-2}$
HH	70.52 / 56.35	69.47 / 58.27	72.28 / 60.51	68.49 / 59.52	65.49 / 56.31
TL;DR	84.21 / 34.78	87.79 / 33.97	89.58 / 36.57	84.25 / 29.93	81.56 / 29.31

**Learning dynamics.** In Figure 3, we compare the learning dynamics of AdamW and logits-SAM when training Mistral-7B on the UltraFeedback dataset. The figure reports training loss, evaluation



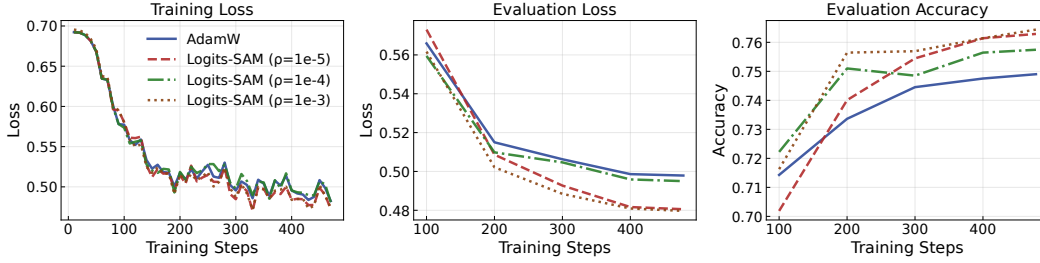


Figure 3: Learning dynamics of Mistral-7B on UltraFeedback. We compare AdamW and logits-SAM in terms of training loss, evaluation loss, and evaluation accuracy, and report curves for logits-SAM under different values of  $\rho$ .

loss, and evaluation accuracy across training steps, and includes curves for logits-SAM under multiple choices of  $\rho$  ( $1 \times 10^{-5}$ ,  $1 \times 10^{-4}$ ,  $1 \times 10^{-3}$ ). We observe that for all three values of  $\rho$ , logits-SAM achieves training loss that is fairly close to that of AdamW, yet consistently attains lower evaluation loss and higher evaluation accuracy. This indicates that logits-SAM provides better generalization than AdamW, and the fact that a range of  $\rho$  values yields consistent improvements suggests that its benefits are robust to the choice of this hyperparameter.

**Sharpness.** To further probe the reasons underlying the generalization gains of logits-SAM, we measure the traces of the parameter Hessian and the logit Hessian at the final checkpoint of Mistral-7B. For AdamW, the traces are  $1.337 \times 10^4 / 2.732 \times 10^2$  (parameter / logit Hessian), while for logits-SAM they are reduced to  $1.186 \times 10^4 / 2.586 \times 10^2$ . This reduction indicates that logits-SAM converges to a flatter solution, which is widely believed to be beneficial for generalization.

**Extension to AI safety and on-policy setting.** Razin et al. (2024) refer to the squeezing effect as *likelihood displacement* and find that, in AI safety scenarios, it can reduce the model’s harmful-response refusal rate, which leads to severe safety concerns. We evaluate the performance of logits-SAM in the same on-policy setting and AI safety scenario as in their work. The reference model is the instruction-tuned Gemma-2B-IT (Team et al., 2024), and the evaluation is conducted using SorryBench (Xie et al., 2024). We train for one epoch with a learning rate of  $1 \times 10^{-6}$  and a batch size of 16. We compare the performance of the reference model, DPO, DPO with logits-SAM, CHES (Razin et al., 2024), and CHES with logits-SAM. For CHES, we filter 50% of samples using the CHES score. Performance is measured by the harmful-response refusal rate and is reported in Table 4. The results show that logits-SAM significantly improves performance in this setting. In particular, DPO with logits-SAM avoids the degradation in refusal rate and performs better than the reference model. Combining logits-SAM with the CHES method of Razin et al. (2024) further increases the refusal rate, with an absolute improvement of approximately 9% on both the training and test sets. These findings indicate that logits-SAM can be effectively transferred to other settings and tasks.

Table 4: Train and test refusal rates for different methods on SorryBench (higher is better).

	Ref model	DPO	DPO+logits-SAM	CHES	CHES+logits-SAM
Train Refusal	0.8054	0.7703	<b>0.8135</b>	0.8459	<b>0.9324</b>
Test Refusal	0.7231	0.7077	<b>0.7538</b>	0.7846	<b>0.8769</b>

## 5 RELATED WORK

**Reinforcement learning from human feedback.** RLHF has emerged as the de facto post-training recipe for aligning large language models (Christiano et al., 2017; Ziegler et al., 2019; Ouyang et al., 2022; Bai et al., 2022), typically combining supervised fine-tuning (Zhou et al., 2023; Taori et al., 2023; Conover et al., 2023; Wang et al., 2023b), reward modeling (Gao et al., 2023; Luo et al., 2023; Lambert et al., 2024), and policy optimization (Schulman et al., 2017; Anthony et al., 2017).

To reduce the complexity and instability of online preference optimization, offline methods such as SLiC-HF (Zhao et al., 2023) and RRHF (Yuan et al., 2023) learn policies from comparisons using closed-form objectives. DPO (Rafailov et al., 2024b) is a central example that maximizes the log-probability margin between preferred and rejected responses relative to a reference policy. Thanks to its simplicity and training stability, DPO has rapidly gained popularity, spurring a line of variants aimed at improving performance. For example, Azar et al. (2024) propose IPO, a more theoretically grounded variant; CPO (Xu et al., 2024) approximates the reference policy as uniform to eliminate the reference term; f-DPO (Wang et al., 2023a) generalizes DPO via a family of  $f$ -divergences; SimPO (Meng et al., 2024a) uses length-normalized scores that better reflect generation-time preferences; and Cal-DPO (Xiao et al., 2024) aligns the implicit reward scale with likelihoods.

**Squeezing effect (likelihood displacement).** The squeezing effect (Ren & Sutherland, 2024), also known as likelihood displacement (Razin et al., 2024), refers to the recently identified phenomenon in which the probability of the ground-truth label is unintentionally reduced during DPO training. This effect has been widely observed and can lead to performance degradation, reduced safety, and even alignment failure (Pal et al., 2024; Yuan et al., 2024; Rafailov et al., 2024a; Tajwar et al., 2024; Pang et al., 2024). Several studies have attempted to mitigate this issue. Asadi et al. (2025) constrain the shift of probability mass between preferred and rejected responses in the reference and target policies. Liu et al. (2025) introduce a KL-divergence-based policy drift constraint to dynamically regularize policy updates. Razin et al. (2024) strengthen safety alignment by filtering samples that are likely to induce likelihood displacement based on the CHES score between token embeddings. Unlike existing approaches, which either focus on designing alternative objective functions or filter the training data, our method takes a pure optimization-based perspective. It is therefore conceptually orthogonal to these techniques and can be used in combination with them.

**Kernel and fixed feature regime of LLMs.** In the context of LLMs, there is a growing body of work that investigates model dynamics through the lens of kernels. A pioneering line of work by Malladi et al. (2023) uses Neural Tangent Kernel-based dynamics (Jacot et al., 2018) to accurately characterize the behavior of LLM fine-tuning and achieves performance comparable to fine-tuning through kernel methods, under the fixed feature assumption. Afzal et al. (2024) leverage the spectrum of the NTK to predict the generalization performance of LLM fine-tuning, and Jang et al. (2024) study the training dynamics of low-rank adaptation from an NTK perspective.

**Sharpness-aware minimization.** A widely held belief in the deep learning community is that flatter solutions typically generalize better (Hochreiter & Schmidhuber, 1997; Keskar et al., 2016; Dinh et al., 2017; Jiang et al., 2019; Xie et al., 2020; Liu et al., 2023). Motivated by this view, SAM (Foret et al., 2021) is a bilevel optimization method that explicitly seeks flatter minima, and it has gained popularity for delivering consistent improvements across a wide range of supervised learning tasks (Foret et al., 2021; Kwon et al., 2021; Kaddour et al., 2022; Liu et al., 2022; Kim et al., 2022; Li & Giannakis, 2023). Most relevant to our work are its recent applications in LLMs. Singh et al. (2025) propose Functional-SAM for LLM pretraining and demonstrate strong performance, while Lee & Yoon (2025) apply SAM to Proximal Policy Optimization to improve robustness in both the reward and action spaces. Logits-SAM is a byproduct mentioned in recent studies, yet it is often overlooked. Baek et al. (2024) analyze the effect of label noise on SAM in linear regression and argue that Jacobian-SAM, the counterpart of logits-SAM, plays the dominant role. Similarly, Singh et al. (2025) identify Jacobian-SAM, also referred to as Functional-SAM, as more important and show that it can effectively improve the generalization performance of LLM pretraining.

## 6 CONCLUSION

We analyzed the squeezing effect in DPO via coordinate-wise dynamics in parameter and logit spaces. Our framework shows that GD with negative  $\eta$  drives residuals to expand along high-curvature directions, and that SAM suppresses this behavior via curvature regularization; in particular, negative  $\eta$  calls for negative  $\rho$ . Motivated by this, we adopt *logits-SAM*, which perturbs only the output layer and adds negligible overhead, and demonstrate consistent gains in effectiveness and robustness across models and datasets. We expect these insights to inform curvature-aware preference optimization going forward.

## REPRODUCIBILITY STATEMENT

All theoretical results presented in this paper are accompanied by complete proofs, which can be found in Appendix A. To further facilitate reproducibility, we will release the source code upon publication, allowing the community to verify and build upon our results.

## REFERENCES

- Zahra Rahimi Afzal, Tara Esmailbeig, Mojtaba Soltanalian, and Mesrob I Ohannessian. Can the spectrum of the neural tangent kernel anticipate fine-tuning performance? In *Adaptive Foundation Models: Evolving AI for Personalized and Efficient Learning*, 2024.
- Thomas Anthony, Zheng Tian, and David Barber. Thinking fast and slow with deep learning and tree search. *Advances in neural information processing systems*, 30, 2017.
- Kavosh Asadi, Julien Han, Idan Pipano, Xingzi Xu, Dominique Perrault-Joncas, Shoham Sabach, Karim Bouyarmane, and Mohammad Ghavamzadeh. C2-dpo: Constrained controlled direct preference optimization, 2025. URL <https://arxiv.org/abs/2502.17507>.
- Mohammad Gheshlaghi Azar, Zhaohan Daniel Guo, Bilal Piot, Remi Munos, Mark Rowland, Michal Valko, and Daniele Calandriello. A general theoretical paradigm to understand learning from human preferences. In *International Conference on Artificial Intelligence and Statistics*, pp. 4447–4455. PMLR, 2024.
- Christina Baek, Zico Kolter, and Aditi Raghunathan. Why is sam robust to label noise? *arXiv preprint arXiv:2405.03676*, 2024.
- Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, et al. Training a helpful and harmless assistant with reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*, 2022.
- Stella Biderman, Hailey Schoelkopf, Quentin Gregory Anthony, Herbie Bradley, Kyle O’Brien, Eric Hallahan, Mohammad Aflah Khan, Shivanshu Purohit, USVSN Sai Prashanth, Edward Raff, et al. Pythia: A suite for analyzing large language models across training and scaling. In *International Conference on Machine Learning*, pp. 2397–2430. PMLR, 2023.
- Ralph Allan Bradley and Milton E Terry. Rank analysis of incomplete block designs: I. the method of paired comparisons. *Biometrika*, 39(3/4):324–345, 1952.
- Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep reinforcement learning from human preferences. *Advances in neural information processing systems*, 30, 2017.
- Mike Conover, Matt Hayes, Ankit Mathur, Jianwei Xie, Jun Wan, Sam Shah, Ali Ghodsi, Patrick Wendell, Matei Zaharia, and Reynold Xin. Free dolly: Introducing the world’s first truly open instructiontuned llm. 2023.
- Ganqu Cui, Lifan Yuan, Ning Ding, Guanming Yao, Wei Zhu, Yuan Ni, Guotong Xie, Zhiyuan Liu, and Maosong Sun. Ultrafeedback: Boosting language models with high-quality feedback, 2023.
- Josef Dai, Xuehai Pan, Ruiyang Sun, Jiaming Ji, Xinbo Xu, Mickel Liu, Yizhou Wang, and Yaodong Yang. Safe rlhf: Safe reinforcement learning from human feedback. *arXiv preprint arXiv:2310.12773*, 2023.
- Ning Ding, Yulin Chen, Bokai Xu, Yujia Qin, Zhi Zheng, Shengding Hu, Zhiyuan Liu, Maosong Sun, and Bowen Zhou. Enhancing chat language models by scaling high-quality instructional conversations, 2023.
- Laurent Dinh, Razvan Pascanu, Samy Bengio, and Yoshua Bengio. Sharp minima can generalize for deep nets. In *International Conference on Machine Learning*, pp. 1019–1028. PMLR, 2017.

- Yann Dubois, Balázs Galambosi, Percy Liang, and Tatsunori B Hashimoto. Length-controlled alpaca-eval: A simple way to debias automatic evaluators. *arXiv preprint arXiv:2404.04475*, 2024.
- Pierre Foret, Ariel Kleiner, Hossein Mobahi, and Behnam Neyshabur. Sharpness-aware minimization for efficiently improving generalization. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net, 2021. URL <https://openreview.net/forum?id=6Tmlmposlrm>.
- Leo Gao, John Schulman, and Jacob Hilton. Scaling laws for reward model overoptimization. In *International Conference on Machine Learning*, pp. 10835–10866. PMLR, 2023.
- Sepp Hochreiter and Jürgen Schmidhuber. Flat minima. *Neural computation*, 9(1):1–42, 1997.
- Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and generalization in neural networks. *Advances in neural information processing systems*, 31, 2018.
- Uijeong Jang, Jason D Lee, and Ernest K Ryu. Lora training in the ntk regime has no spurious local minima. *arXiv preprint arXiv:2402.11867*, 2024.
- Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Léo Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. Mistral 7b, 2023. URL <https://arxiv.org/abs/2310.06825>.
- Yiding Jiang, Behnam Neyshabur, Hossein Mobahi, Dilip Krishnan, and Samy Bengio. Fantastic generalization measures and where to find them. *arXiv preprint arXiv:1912.02178*, 2019.
- Jean Kaddour, Linqing Liu, Ricardo Silva, and Matt J Kusner. When do flat minima optimizers work? *Advances in Neural Information Processing Systems*, 35:16577–16595, 2022.
- Nitish Shirish Keskar, Dheevatsa Mudigere, Jorge Nocedal, Mikhail Smelyanskiy, and Ping Tak Peter Tang. On large-batch training for deep learning: Generalization gap and sharp minima. *arXiv preprint arXiv:1609.04836*, 2016.
- Minyoung Kim, Da Li, Shell X Hu, and Timothy Hospedales. Fisher sam: Information geometry and sharpness aware minimisation. In *International Conference on Machine Learning*, pp. 11148–11161. PMLR, 2022.
- Jungmin Kwon, Jeongseop Kim, Hyunseo Park, and In Kwon Choi. Asam: Adaptive sharpness-aware minimization for scale-invariant learning of deep neural networks. In *International Conference on Machine Learning*, pp. 5905–5914. PMLR, 2021.
- Nathan Lambert, Valentina Pyatkin, Jacob Morrison, LJ Miranda, Bill Yuchen Lin, Khyathi Chandu, Nouha Dziri, Sachin Kumar, Tom Zick, Yejin Choi, et al. Rewardbench: Evaluating reward models for language modeling. *arXiv preprint arXiv:2403.13787*, 2024.
- Hyun Kyu Lee and Sung Whan Yoon. Flat reward in policy parameter space implies robust reinforcement learning. In *The Thirteenth International Conference on Learning Representations*, 2025.
- Bingcong Li and Georgios B. Giannakis. Enhancing sharpness-aware optimization through variance suppression, 2023. URL <https://arxiv.org/abs/2309.15639>.
- Tianle Li, Wei-Lin Chiang, Evan Frick, Lisa Dunlap, Tianhao Wu, Banghua Zhu, Joseph E. Gonzalez, and Ion Stoica. From crowdsourced data to high-quality benchmarks: Arena-hard and benchbuilder pipeline, 2024. URL <https://arxiv.org/abs/2406.11939>.
- Hong Liu, Sang Michael Xie, Zhiyuan Li, and Tengyu Ma. Same pre-training loss, better downstream: Implicit bias matters for language models. In *International Conference on Machine Learning*, pp. 22188–22214. PMLR, 2023.
- Yong Liu, Siqi Mai, Xiangning Chen, Cho-Jui Hsieh, and Yang You. Towards efficient and scalable sharpness-aware minimization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 12360–12370, 2022.

- Zongkai Liu, Fanqing Meng, Lingxiao Du, Zhixiang Zhou, Chao Yu, Wenqi Shao, and Qiaosheng Zhang. Cpgd: Toward stable rule-based reinforcement learning for language models. *arXiv preprint arXiv:2505.12504*, 2025.
- Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization, 2019. URL <https://arxiv.org/abs/1711.05101>.
- Haipeng Luo, Qingfeng Sun, Can Xu, Pu Zhao, Jianguang Lou, Chongyang Tao, Xiubo Geng, Qingwei Lin, Shifeng Chen, and Dongmei Zhang. Wizardmath: Empowering mathematical reasoning for large language models via reinforced evol-instruct. *arXiv preprint arXiv:2308.09583*, 2023.
- Sadhika Malladi, Alexander Wettig, Dingli Yu, Danqi Chen, and Sanjeev Arora. A kernel-based view of language model fine-tuning. In *International Conference on Machine Learning*, pp. 23610–23641. PMLR, 2023.
- Yu Meng, Mengzhou Xia, and Danqi Chen. Simpo: Simple preference optimization with a reference-free reward. *Advances in Neural Information Processing Systems*, 37:124198–124235, 2024a.
- Yu Meng, Mengzhou Xia, and Danqi Chen. Simpo: Simple preference optimization with a reference-free reward, 2024b. URL <https://arxiv.org/abs/2405.14734>.
- Reiichiro Nakano, Jacob Hilton, Suchir Balaji, Jeff Wu, Long Ouyang, Christina Kim, Christopher Hesse, Shantanu Jain, Vineet Kosaraju, William Saunders, Xu Jiang, Karl Cobbe, Tyna Eloundou, Gretchen Krueger, Kevin Button, Matthew Knight, Benjamin Chess, and John Schulman. Webgpt: Browser-assisted question-answering with human feedback, 2022. URL <https://arxiv.org/abs/2112.09332>.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35: 27730–27744, 2022.
- Arka Pal, Deep Karkhanis, Samuel Dooley, Manley Roberts, Siddhartha Naidu, and Colin White. Smaug: Fixing failure modes of preference optimisation with dpo-positive. *arXiv preprint arXiv:2402.13228*, 2024.
- Richard Yuanzhe Pang, Weizhe Yuan, He He, Kyunghyun Cho, Sainbayar Sukhbaatar, and Jason Weston. Iterative reasoning preference optimization. *Advances in Neural Information Processing Systems*, 37:116617–116637, 2024.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- Rafael Rafailov, Joey Hejna, Ryan Park, and Chelsea Finn. From  $r$  to  $q^*$ : Your language model is secretly a  $q$ -function. *arXiv preprint arXiv:2404.12358*, 2024a.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D. Manning, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model, 2024b. URL <https://arxiv.org/abs/2305.18290>.
- Jeff Rasley, Samyam Rajbhandari, Olatunji Ruwase, and Yuxiong He. Deepspeed: System optimizations enable training deep learning models with over 100 billion parameters. In *Proceedings of the 26th ACM SIGKDD international conference on knowledge discovery & data mining*, pp. 3505–3506, 2020.
- Noam Razin, Sadhika Malladi, Adithya Bhaskar, Danqi Chen, Sanjeev Arora, and Boris Hanin. Unintentional unalignment: Likelihood displacement in direct preference optimization. *arXiv preprint arXiv:2410.08847*, 2024.
- Yi Ren and Danica J Sutherland. Learning dynamics of llm finetuning. *arXiv preprint arXiv:2407.10490*, 2024.

- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- Sidak Pal Singh, Hossein Mobahi, Atish Agarwala, and Yann Dauphin. Avoiding spurious sharpness minimization broadens applicability of sam. *arXiv preprint arXiv:2502.02407*, 2025.
- Nisan Stiennon, Long Ouyang, Jeff Wu, Daniel M. Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford, Dario Amodei, and Paul Christiano. Learning to summarize from human feedback. In *NeurIPS*, 2020.
- Fahim Tajwar, Anikait Singh, Archit Sharma, Rafael Rafailov, Jeff Schneider, Tengyang Xie, Stefano Ermon, Chelsea Finn, and Aviral Kumar. Preference fine-tuning of llms should leverage suboptimal, on-policy data. *arXiv preprint arXiv:2404.14367*, 2024.
- Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B Hashimoto. Stanford alpaca: An instruction-following llama model, 2023.
- Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, et al. Gemma: Open models based on gemini research and technology. *arXiv preprint arXiv:2403.08295*, 2024.
- Lewis Tunstall, Edward Beeching, Nathan Lambert, Nazneen Rajani, Shengyi Huang, Kashif Rasul, Alvaro Bartolome, Carlos M. Patiño, Alexander M. Rush, and Thomas Wolf. The Alignment Handbook, 2023a. URL <https://github.com/huggingface/alignment-handbook>.
- Lewis Tunstall, Edward Beeching, Nathan Lambert, Nazneen Rajani, Kashif Rasul, Younes Belkada, Shengyi Huang, Leandro von Werra, Clémentine Fourrier, Nathan Habib, Nathan Sarrazin, Omar Sanseviero, Alexander M. Rush, and Thomas Wolf. Zephyr: Direct distillation of lm alignment, 2023b. URL <https://arxiv.org/abs/2310.16944>.
- Chaoqi Wang, Yibo Jiang, Chenghao Yang, Han Liu, and Yuxin Chen. Beyond reverse kl: Generalizing direct preference optimization with diverse divergence constraints, 2023a. URL <https://arxiv.org/abs/2309.16240>.
- Guan Wang, Sijie Cheng, Xianyu Zhan, Xiangang Li, Sen Song, and Yang Liu. Openchat: Advancing open-source language models with mixed-quality data. *arXiv preprint arXiv:2309.11235*, 2023b.
- Teng Xiao, Yige Yuan, Huaisheng Zhu, Mingxiao Li, and Vasant G Honavar. Cal-dpo: Calibrated direct preference optimization for language model alignment. *Advances in Neural Information Processing Systems*, 37:114289–114320, 2024.
- Tinghao Xie, Xiangyu Qi, Yi Zeng, Yangsibo Huang, Udari Madhushani Sehwag, Kaixuan Huang, Luxi He, Boyi Wei, Dacheng Li, Ying Sheng, et al. Sorry-bench: Systematically evaluating large language model safety refusal. *arXiv preprint arXiv:2406.14598*, 2024.
- Zeke Xie, Issei Sato, and Masashi Sugiyama. A diffusion theory for deep learning dynamics: Stochastic gradient descent exponentially favors flat minima. *arXiv preprint arXiv:2002.03495*, 2020.
- Haoran Xu, Amr Sharaf, Yunmo Chen, Weiting Tan, Lingfeng Shen, Benjamin Van Durme, Kenton Murray, and Young Jin Kim. Contrastive preference optimization: Pushing the boundaries of llm performance in machine translation, 2024. URL <https://arxiv.org/abs/2401.08417>.
- Lifan Yuan, Ganqu Cui, Hanbin Wang, Ning Ding, Xingyao Wang, Jia Deng, Boji Shan, Huimin Chen, Ruobing Xie, Yankai Lin, et al. Advancing llm reasoning generalists with preference trees. *arXiv preprint arXiv:2404.02078*, 2024.
- Zheng Yuan, Hongyi Yuan, Chuanqi Tan, Wei Wang, Songfang Huang, and Fei Huang. Rrhf: Rank responses to align language models with human feedback without tears. *arXiv preprint arXiv:2304.05302*, 2023.

Yao Zhao, Rishabh Joshi, Tianqi Liu, Misha Khalman, Mohammad Saleh, and Peter J. Liu. Slic-hf: Sequence likelihood calibration with human feedback, 2023. URL <https://arxiv.org/abs/2305.10425>.

Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. Judging llm-as-a-judge with mt-bench and chatbot arena, 2023. URL <https://arxiv.org/abs/2306.05685>.

Chunting Zhou, Pengfei Liu, Puxin Xu, Srinivasan Iyer, Jiao Sun, Yuning Mao, Xuezhe Ma, Avia Efrat, Ping Yu, Lili Yu, et al. Lima: Less is more for alignment. *Advances in Neural Information Processing Systems*, 36:55006–55021, 2023.

Daniel M Ziegler, Nisan Stiennon, Jeffrey Wu, Tom B Brown, Alec Radford, Dario Amodei, Paul Christiano, and Geoffrey Irving. Fine-tuning language models from human preferences. *arXiv preprint arXiv:1909.08593*, 2019.

## A FORMAL THEOREMS AND PROOFS

**Proposition A.1** (Geometry of the logit space and the parameter–logit correspondence). *Let  $\ell : \mathbb{R}^V \rightarrow \mathbb{R}$  be  $C^2$ . Fix an input  $\mathbf{x}$  and a feature map  $\phi(\mathbf{x}) \in \mathbb{R}^d$ . For  $\mathbf{W} \in \mathbb{R}^{V \times d}$  set*

$$\mathbf{z} = \mathbf{W} \phi \in \mathbb{R}^V, \quad F(\mathbf{W}) = \ell(\mathbf{z}).$$

*Denote  $\mathbf{H}_z := \nabla_z^2 \ell(\mathbf{z}) \in \mathbb{R}^{V \times V}$  and  $\mathbf{H}_W := \nabla_W^2 F(\mathbf{W}) \in \mathbb{R}^{Vd \times Vd}$ .*

*Equip  $\mathbb{R}^{V \times d}$  with the Frobenius inner product  $\langle \mathbf{A}, \mathbf{B} \rangle_F = \text{tr}(\mathbf{A}^\top \mathbf{B})$  and  $\mathbb{R}^V$  with the Euclidean inner product. Let*

$$T_\phi : \mathbb{R}^{V \times d} \rightarrow \mathbb{R}^V, \quad T_\phi(\Delta \mathbf{W}) = \Delta \mathbf{W} \phi$$

*be the differential of the map  $\mathbf{W} \mapsto \mathbf{W} \phi$ , and let  $T_\phi^* : \mathbb{R}^V \rightarrow \mathbb{R}^{V \times d}$  be its adjoint with respect to these inner products, i.e.,  $\langle T_\phi(\Delta \mathbf{W}), \mathbf{v} \rangle = \langle \Delta \mathbf{W}, T_\phi^*(\mathbf{v}) \rangle_F$  for all  $\Delta \mathbf{W}, \mathbf{v}$ . Then  $T_\phi^*(\mathbf{v}) = \mathbf{v} \phi^\top$ . The following statements hold.*

(1) **Pullback identity (operator form).**

$$\mathbf{H}_W = T_\phi^* \mathbf{H}_z T_\phi$$

*as linear operators on  $\mathbb{R}^{V \times d}$ . Equivalently, in coordinates,*

$$\nabla_W F(\mathbf{W}) = (\nabla_z \ell(\mathbf{z})) \phi^\top, \quad \mathbf{H}_W = \mathbf{H}_z \otimes (\phi \phi^\top).$$

*Consequently, if  $\phi \neq \mathbf{0}$ , then*

$$\text{rank}(\mathbf{H}_W) = \text{rank}(\mathbf{H}_z).$$

(2) **Pullback of the bilinear form.** *For every  $\Delta \mathbf{W}, \Delta \mathbf{W}' \in \mathbb{R}^{V \times d}$ ,*

$$\langle \Delta \mathbf{W}, \mathbf{H}_W[\Delta \mathbf{W}'] \rangle_F = \langle T_\phi(\Delta \mathbf{W}), \mathbf{H}_z T_\phi(\Delta \mathbf{W}') \rangle$$

*and,*

$$\mathbf{H}_W[\Delta \mathbf{W}] = T_\phi^*(\mathbf{H}_z T_\phi(\Delta \mathbf{W})) = \mathbf{H}_z \Delta \mathbf{W} (\phi \phi^\top).$$

*Thus the second-order effect of any parameter perturbation depends only on the induced logits perturbation  $T_\phi(\Delta \mathbf{W}) = \Delta \mathbf{W} \phi$ .*

(3) **Surjectivity, kernel, and quotient-space view.** *If  $\phi \neq \mathbf{0}$ , then  $T_\phi$  is surjective. For any  $\Delta \mathbf{z} \in \mathbb{R}^V$ , a minimum-Frobenius-norm preimage is*

$$\Delta \mathbf{W}_* = \frac{\Delta \mathbf{z} \phi^\top}{\|\phi\|^2} \quad \text{with} \quad T_\phi(\Delta \mathbf{W}_*) = \Delta \mathbf{z}.$$

*The kernel is*

$$\ker(T_\phi) = \{ \Delta \mathbf{W} \in \mathbb{R}^{V \times d} : \Delta \mathbf{W} \phi = \mathbf{0} \},$$

*of dimension  $V(d-1)$ . Consequently,  $\mathbf{H}_W$  descends to the quotient  $\mathbb{R}^{V \times d} / \ker(T_\phi) \cong \mathbb{R}^V$ .*



*Proof.* A direct computation gives

$$\langle T_\phi(\Delta \mathbf{W}), \mathbf{v} \rangle = \text{tr}((\Delta \mathbf{W} \phi)^\top \mathbf{v}) = \text{tr}(\Delta \mathbf{W}^\top \mathbf{v} \phi^\top) = \langle \Delta \mathbf{W}, \mathbf{v} \phi^\top \rangle_F,$$

hence

$$T_\phi^*(\mathbf{v}) = \mathbf{v} \phi^\top.$$

**(1) Pullback identity and coordinate forms.** Let  $F(\mathbf{W}) = \ell(\mathbf{W} \phi)$ . The first differential of  $F$  is

$$dF[\Delta \mathbf{W}] = \langle \nabla_z \ell(z), T_\phi(\Delta \mathbf{W}) \rangle = \langle T_\phi^*(\nabla_z \ell(z)), \Delta \mathbf{W} \rangle_F,$$

so

$$\nabla_{\mathbf{W}} F(\mathbf{W}) = T_\phi^*(\nabla_z \ell(z)) = (\nabla_z \ell(z)) \phi^\top.$$

Differentiating once more and using  $d(\nabla_z \ell(z))[\Delta z] = \mathbf{H}_z \Delta z$  with  $\Delta z = T_\phi(\Delta \mathbf{W})$  yields, for all  $\Delta \mathbf{W}, \Delta \mathbf{W}'$ ,

$$d^2 F[\Delta \mathbf{W}, \Delta \mathbf{W}'] = \langle T_\phi(\Delta \mathbf{W}), \mathbf{H}_z T_\phi(\Delta \mathbf{W}') \rangle.$$

By the Riesz representation on  $(\mathbb{R}^{V \times d}, \langle \cdot, \cdot \rangle_F)$ , this means

$$\mathbf{H}_{\mathbf{W}} = T_\phi^* \mathbf{H}_z T_\phi.$$

Using  $T_\phi^*(\mathbf{v}) = \mathbf{v} \phi^\top$  and  $T_\phi(\Delta \mathbf{W}) = \Delta \mathbf{W} \phi$ ,

$$\mathbf{H}_{\mathbf{W}}[\Delta \mathbf{W}] = T_\phi^*(\mathbf{H}_z(\Delta \mathbf{W} \phi)) = (\mathbf{H}_z(\Delta \mathbf{W} \phi)) \phi^\top = \mathbf{H}_z \Delta \mathbf{W} (\phi \phi^\top),$$

which is the coordinate (Kronecker) form used in the main text.

For the rank statement, assume  $\phi \neq \mathbf{0}$ . Then  $T_\phi$  is surjective and  $T_\phi^*$  is injective. Hence

$$\text{rank}(\mathbf{H}_{\mathbf{W}}) = \text{rank}(T_\phi^* \mathbf{H}_z T_\phi) = \text{rank}(\mathbf{H}_z T_\phi) = \text{rank}(\mathbf{H}_z),$$

because  $\text{range}(T_\phi) = \mathbb{R}^V$ .

**(2) Pullback of the bilinear form.** By the operator identity above,

$$\langle \Delta \mathbf{W}, \mathbf{H}_{\mathbf{W}}[\Delta \mathbf{W}'] \rangle_F = \langle \Delta \mathbf{W}, T_\phi^* \mathbf{H}_z T_\phi(\Delta \mathbf{W}') \rangle_F = \langle T_\phi(\Delta \mathbf{W}), \mathbf{H}_z T_\phi(\Delta \mathbf{W}') \rangle.$$

Equivalently,  $\mathbf{H}_{\mathbf{W}}[\Delta \mathbf{W}] = T_\phi^*(\mathbf{H}_z T_\phi(\Delta \mathbf{W})) = \mathbf{H}_z \Delta \mathbf{W} (\phi \phi^\top)$ . Thus the bilinear form on parameter space is the pullback of the bilinear form induced by  $\mathbf{H}_z$  on logit space.

**(3) Surjectivity, kernel and quotient view.** If  $\phi \neq \mathbf{0}$ , then for any  $\Delta z \in \mathbb{R}^V$

$$\Delta \mathbf{W}_* = \frac{\Delta z \phi^\top}{\|\phi\|^2} \quad \text{satisfies} \quad T_\phi(\Delta \mathbf{W}_*) = \Delta z,$$

so  $T_\phi$  is surjective. The same choice minimizes the Frobenius norm among all preimages (row-wise Cauchy–Schwarz). The kernel is  $\ker(T_\phi) = \{\Delta \mathbf{W} : \Delta \mathbf{W} \phi = \mathbf{0}\}$ , and rank–nullity gives  $\dim \ker(T_\phi) = V(d-1)$ . Finally, if  $\Delta \mathbf{W}_1 - \Delta \mathbf{W}_2 \in \ker(T_\phi)$ , then  $T_\phi(\Delta \mathbf{W}_1) = T_\phi(\Delta \mathbf{W}_2)$  and

$$\begin{aligned} \langle \Delta \mathbf{W}_1, \mathbf{H}_{\mathbf{W}}[\Delta \mathbf{W}_1] \rangle_F &= \langle T_\phi(\Delta \mathbf{W}_1), \mathbf{H}_z T_\phi(\Delta \mathbf{W}_1) \rangle = \langle T_\phi(\Delta \mathbf{W}_2), \mathbf{H}_z T_\phi(\Delta \mathbf{W}_2) \rangle \\ &= \langle \Delta \mathbf{W}_2, \mathbf{H}_{\mathbf{W}}[\Delta \mathbf{W}_2] \rangle_F, \end{aligned}$$

so the bilinear form descends to the quotient  $\mathbb{R}^{V \times d} / \ker(T_\phi) \cong \mathbb{R}^V$ .

If  $\phi = \mathbf{0}$  then  $T_\phi \equiv 0$  and  $\mathbf{H}_{\mathbf{W}} \equiv \mathbf{0}$ , the degenerate case.  $\square$

**Theorem A.2** (Dynamics of SAM). *Fix a SAMple  $\mathbf{x}$  and set  $\mu = \|\phi\|^2 < \infty$ . Assume:*

- (1)  $f(\mathbf{z}, \mathbf{y})$  is  $C^3$  in  $\mathbf{z}$  and there exists  $L < \infty$  such that  $\sup_{\mathbf{z}} \|\nabla_{\mathbf{z}}^3 f(\mathbf{z}, \mathbf{y})\| \leq L$ .
- (2) The step size  $|\eta| \in (0, 1]$  and the SAM radius satisfies  $|\rho| \leq \kappa \sqrt{|\eta|}$  with a constant  $\kappa \geq 0$ .
- (3) If  $\|\mathbf{g}^t\| = 0$ , set the inner perturbation to 0 and define  $\tilde{\rho}^t = 0$ ; otherwise  $\tilde{\rho}^t := \rho \sqrt{\mu} / \|\mathbf{g}^t\|$ .

Consider standard SAM:

$$\Delta \mathbf{W}_{\text{adv}}^t = \rho \frac{\nabla_{\mathbf{W}} f(\mathbf{W}^t)}{\|\nabla_{\mathbf{W}} f(\mathbf{W}^t)\|_F}, \quad \widetilde{\mathbf{W}}^t = \mathbf{W}^t + \Delta \mathbf{W}_{\text{adv}}^t, \quad \mathbf{W}^{t+1} = \mathbf{W}^t - \eta \nabla_{\mathbf{W}} f(\widetilde{\mathbf{W}}^t).$$

Then, there exists a constant  $C > 0$  (depending only on  $L, \mu, \kappa$ ) such that the following expansions hold with  $O(\eta^2)$  remainders:

$$(\text{parameters}) \quad \mathbf{W}^{t+1} = \mathbf{W}^t - \eta \left( \mathbf{g}^t \phi^\top + \tilde{\rho}^t \mathbf{H}_z^t \mathbf{g}^t \phi^\top \right) + \mathbf{R}_{\mathbf{W}}^t, \quad \|\mathbf{R}_{\mathbf{W}}^t\|_F \leq C \eta^2,$$

$$(\text{logits}) \quad \mathbf{z}^{t+1} = \mathbf{z}^t - \eta \mu \left( \mathbf{g}^t + \tilde{\rho}^t \mathbf{H}_z^t \mathbf{g}^t \right) + \mathbf{r}_z^t, \quad \|\mathbf{r}_z^t\| \leq C \eta^2,$$

$$(\text{logit gradient}) \quad \mathbf{g}^{t+1} = \left( \mathbf{I} - \eta \mu \mathbf{H}_z^t - \eta \mu \tilde{\rho}^t (\mathbf{H}_z^t)^2 \right) \mathbf{g}^t + \mathbf{r}_g^t, \quad \|\mathbf{r}_g^t\| \leq C \eta^2.$$

In particular, for softmax cross-entropy where  $\mathbf{g}^t = \mathbf{p}^t - \mathbf{y}$  and

$$(\text{residual}) \quad \mathbf{p}^{t+1} - \mathbf{y} = \left( \mathbf{I} - \eta \mu \mathbf{H}_z^t - \eta \mu \tilde{\rho}^t (\mathbf{H}_z^t)^2 \right) (\mathbf{p}^t - \mathbf{y}) + \mathbf{r}_g^t, \quad \|\mathbf{r}_g^t\| \leq C \eta^2.$$

*Proof.* Write  $F(\mathbf{W}) := f(\mathbf{W}\phi, \mathbf{y})$  and  $\mathbf{z} = \mathbf{W}\phi$ . By Proposition A.1 (Pullback/Kronecker and operator forms),

$$\nabla_{\mathbf{W}} F(\mathbf{W}) = \mathbf{g} \phi^\top, \quad \mathbf{H}_{\mathbf{W}}[\Delta \mathbf{W}] = \mathbf{H}_z \Delta \mathbf{W} (\phi \phi^\top),$$

and  $T_\phi(\Delta \mathbf{W}) = \Delta \mathbf{W} \phi$  with  $\|T_\phi\| \leq \|\phi\| = \sqrt{\mu}$ . Moreover, by the multilinear chain rule applied to  $F(\mathbf{W}) = f(\mathbf{W}\phi, \mathbf{y})$ ,

$$\nabla_{\mathbf{W}}^3 F(\mathbf{W})[\Delta_1, \Delta_2, \Delta_3] = \nabla_{\mathbf{z}}^3 f(\mathbf{z}, \mathbf{y}) [T_\phi(\Delta_1), T_\phi(\Delta_2), T_\phi(\Delta_3)], \quad (7)$$

hence the operator norm satisfies

$$\sup_{\mathbf{W}} \|\nabla_{\mathbf{W}}^3 F(\mathbf{W})\| \leq \left( \sup_{\mathbf{z}} \|\nabla_{\mathbf{z}}^3 f(\mathbf{z}, \mathbf{y})\| \right) \|T_\phi\|^3 \leq L \mu^{3/2}. \quad (8)$$

(i) **Parameter update.** Let

$$\Delta \mathbf{W}_{\text{adv}}^t = \rho \frac{\nabla_{\mathbf{W}} F(\mathbf{W}^t)}{\|\nabla_{\mathbf{W}} F(\mathbf{W}^t)\|_F} \quad \text{and} \quad \widetilde{\mathbf{W}}^t = \mathbf{W}^t + \Delta \mathbf{W}_{\text{adv}}^t.$$

If  $\|\mathbf{g}^t\| > 0$ , then  $\nabla_{\mathbf{W}} F(\mathbf{W}^t) = \mathbf{g}^t \phi^\top$  and  $\|\mathbf{g}^t \phi^\top\|_F = \|\mathbf{g}^t\| \|\phi\| = \|\mathbf{g}^t\| \sqrt{\mu}$ , so

$$\Delta \mathbf{W}_{\text{adv}}^t = \rho \frac{\mathbf{g}^t \phi^\top}{\|\mathbf{g}^t\| \sqrt{\mu}}, \quad \|\Delta \mathbf{W}_{\text{adv}}^t\|_F = |\rho| \leq \kappa \sqrt{|\eta|}.$$

(If  $\|\mathbf{g}^t\| = 0$ , our convention sets  $\Delta \mathbf{W}_{\text{adv}}^t = \mathbf{0}$ .) A second-order Taylor expansion of  $\nabla_{\mathbf{W}} F$  at  $\mathbf{W}^t$  gives, for some  $\theta \in (0, 1)$ ,

$$\nabla_{\mathbf{W}} F(\widetilde{\mathbf{W}}^t) = \nabla_{\mathbf{W}} F(\mathbf{W}^t) + \mathbf{H}_{\mathbf{W}}^t[\Delta \mathbf{W}_{\text{adv}}^t] + \frac{1}{2} \nabla_{\mathbf{W}}^3 F(\mathbf{W}^t + \theta \Delta \mathbf{W}_{\text{adv}}^t) [\Delta \mathbf{W}_{\text{adv}}^t, \Delta \mathbf{W}_{\text{adv}}^t].$$

By equation 8 and  $\|\Delta \mathbf{W}_{\text{adv}}^t\|_F \leq \kappa \sqrt{|\eta|}$ ,

$$\left\| \frac{1}{2} \nabla_{\mathbf{W}}^3 F(\cdot) [\Delta \mathbf{W}_{\text{adv}}^t, \Delta \mathbf{W}_{\text{adv}}^t] \right\| \leq \frac{1}{2} L \mu^{3/2} \|\Delta \mathbf{W}_{\text{adv}}^t\|_F^2 \leq C_0 |\eta|,$$

for a constant  $C_0 = C_0(L, \mu, \kappa)$ . Using the operator identity from Proposition A.1,

$$\mathbf{H}_{\mathbf{W}}^t[\Delta \mathbf{W}_{\text{adv}}^t] = \mathbf{H}_z^t \Delta \mathbf{W}_{\text{adv}}^t (\phi \phi^\top) = \frac{\rho \sqrt{\mu}}{\|\mathbf{g}^t\|} \mathbf{H}_z^t \mathbf{g}^t \phi^\top = \tilde{\rho}^t \mathbf{H}_z^t \mathbf{g}^t \phi^\top.$$

Therefore

$$\mathbf{W}^{t+1} = \mathbf{W}^t - \eta \left( \mathbf{g}^t \phi^\top + \tilde{\rho}^t \mathbf{H}_z^t \mathbf{g}^t \phi^\top \right) - \eta \mathbf{R}_{\nabla}^t,$$

where  $\|\mathbf{R}_{\nabla}^t\|_F \leq C_0 |\eta|$ . Setting  $\mathbf{R}_{\mathbf{W}}^t := -\eta \mathbf{R}_{\nabla}^t$  yields  $\|\mathbf{R}_{\mathbf{W}}^t\|_F \leq C \eta^2$  with  $C = C(L, \mu, \kappa)$ , proving the parameter expansion.

(ii) **Logit update.** Right-multiplying by  $\phi$  and using  $\mu = \|\phi\|^2$ ,

$$\mathbf{z}^{t+1} - \mathbf{z}^t = (\mathbf{W}^{t+1} - \mathbf{W}^t)\phi = -\eta\mu(\mathbf{g}^t + \tilde{\rho}^t \mathbf{H}_z^t \mathbf{g}^t) + \mathbf{r}_z^t,$$

with  $\|\mathbf{r}_z^t\| \leq \|\mathbf{R}_W^t\|_F \|\phi\| \leq C\eta^2$  (absorbing  $\sqrt{\mu}$  into  $C$ ). This proves the logits expansion.

(iii) **logit gradient update.** Since  $\mathbf{g} = \nabla_z f(\mathbf{z}, \mathbf{y})$ , a first-order Taylor expansion at  $\mathbf{z}^t$  gives

$$\mathbf{g}^{t+1} = \mathbf{g}^t + \mathbf{H}_z^t(\mathbf{z}^{t+1} - \mathbf{z}^t) + \frac{1}{2} \nabla_z^3 f(\mathbf{z}^t + \xi^t, \mathbf{y})[\Delta \mathbf{z}^t, \Delta \mathbf{z}^t], \quad \Delta \mathbf{z}^t = \mathbf{z}^{t+1} - \mathbf{z}^t.$$

By assumption  $\|\nabla_z^3 f\| \leq L$  and  $\|\Delta \mathbf{z}^t\| = O(\eta)$  from the previous step, hence the remainder has norm  $\leq C_1 \eta^2$ . Substituting the logits expansion from step (ii) yields

$$\mathbf{g}^{t+1} = \left( \mathbf{I} - \eta\mu \mathbf{H}_z^t - \eta\mu \tilde{\rho}^t (\mathbf{H}_z^t)^2 \right) \mathbf{g}^t + \mathbf{r}_g^t, \quad \|\mathbf{r}_g^t\| \leq C\eta^2,$$

after absorbing constants into  $C$ . This proves the logit gradient statement.

Combining (i)–(iii) completes the proof, with a constant  $C$  depending only on  $(L, \mu, \kappa)$ , and the bounds holding for all  $|\eta| \in (0, 1]$  and  $|\rho| \leq \kappa\sqrt{|\eta|}$ .

For softmax cross-entropy,

$$\nabla_z f(\mathbf{z}, \mathbf{y}) = \mathbf{p}(\mathbf{z}) - \mathbf{y}, \quad \mathbf{H}_z(\mathbf{z}) = \nabla_z^2 f(\mathbf{z}, \mathbf{y}) = \text{Diag}(\mathbf{p}(\mathbf{z})) - \mathbf{p}(\mathbf{z})\mathbf{p}(\mathbf{z})^\top.$$

Since  $\mathbf{p}(\mathbf{z}) \in \Delta^{V-1} \subset [0, 1]^V$  for all  $\mathbf{z}$ , every entry of the third derivative tensor  $\nabla_z^3 f(\mathbf{z}, \mathbf{y})$  is a bounded polynomial in  $\mathbf{p}(\mathbf{z})$  (hence in  $[0, 1]$ ). Therefore there exists a finite constant  $L_{\text{sm}}(V)$  depending only on  $V$  such that

$$\sup_{\mathbf{z}} \|\nabla_z^3 f(\mathbf{z}, \mathbf{y})\| \leq L_{\text{sm}}(V).$$

In particular,  $f$  is  $C^\infty$  and Assumption (1) of the theorem holds with  $L = L_{\text{sm}}(V)$ .  $\square$

**Proposition A.3.**  $\mathbf{H}_z$  is symmetric positive semidefinite with  $\ker(\mathbf{H}_z) = \text{span}\{\mathbf{1}\}$  and  $\text{rank}(\mathbf{H}_z) = V - 1$ . Moreover, for the residual  $\mathbf{g}$  we have  $\mathbf{1}^\top \mathbf{g} = 0$ , hence  $\mathbf{g} \in \mathbf{1}^\perp = \text{range}(\mathbf{H}_z)$ ; in particular, given any eigenbasis of  $\mathbf{H}_z$  restricted to  $\mathbf{1}^\perp$ ,  $\mathbf{g}$  admits a unique coordinate representation in that basis.

*Proof.* Let  $\mathbf{p} = \text{softmax}(\mathbf{z}) \in (0, 1)^V$  so that  $\mathbf{1}^\top \mathbf{p} = 1$ , and recall

$$\mathbf{H}_z = \text{Diag}(\mathbf{p}) - \mathbf{p}\mathbf{p}^\top.$$

For any  $\mathbf{v} \in \mathbb{R}^V$ ,

$$\mathbf{v}^\top \mathbf{H}_z \mathbf{v} = \sum_{i=1}^V p_i v_i^2 - \left( \sum_{i=1}^V p_i v_i \right)^2 = \text{Var}_{\mathbf{p}}(v) \geq 0,$$

hence  $\mathbf{H}_z$  is symmetric positive semidefinite. Moreover,  $\mathbf{v}^\top \mathbf{H}_z \mathbf{v} = 0$  iff  $\text{Var}_{\mathbf{p}}(v) = 0$ , i.e.,  $v_i$  is constant across  $i$ . Since  $p_i > 0$  for all  $i$ , this means  $\mathbf{v} = c\mathbf{1}$ , thus

$$\ker(\mathbf{H}_z) = \text{span}\{\mathbf{1}\} \Rightarrow \text{rank}(\mathbf{H}_z) = V - \dim \ker(\mathbf{H}_z) = V - 1.$$

Then  $\mathbf{1}^\top \mathbf{g} = \mathbf{1}^\top \mathbf{p} - \mathbf{1}^\top \mathbf{y} = 0$ , so  $\mathbf{g} \in \mathbf{1}^\perp$ . For any symmetric matrix,  $\text{range}(\mathbf{H}_z) = (\ker(\mathbf{H}_z))^\perp$ ; using  $\ker(\mathbf{H}_z) = \text{span}\{\mathbf{1}\}$  yields  $\mathbf{1}^\perp = \text{range}(\mathbf{H}_z)$ , hence  $\mathbf{g} \in \text{range}(\mathbf{H}_z)$ .

Restrict  $\mathbf{H}_z$  to the invariant subspace  $\mathbf{1}^\perp$ . Being symmetric,  $\mathbf{H}_z|_{\mathbf{1}^\perp}$  admits an orthonormal eigenbasis  $\{\mathbf{v}_k\}_{k=1}^{V-1}$  associated with its positive eigenvalues. Since  $\mathbf{g} \in \mathbf{1}^\perp$ , it has the unique expansion  $\mathbf{g} = \sum_{k=1}^{V-1} e_k \mathbf{v}_k$ , with  $e_k = (\mathbf{v}_k)^\top \mathbf{g}$ .  $\square$

**Corollary A.4** (Modal dynamics in the eigenbasis of  $\mathbf{H}_z^t$ ). *Under the same assumptions as Theorem A.2. For each  $t$ , let the spectral decomposition of the symmetric positive-semidefinite matrix  $\mathbf{H}_z^t$  be*

$$\mathbf{H}_z^t = \sum_{k=1}^{V-1} \lambda_k^t \mathbf{v}_k^t (\mathbf{v}_k^t)^\top,$$

where  $\lambda_k^t > 0$ ,  $(\mathbf{v}_k^t)^\top \mathbf{v}_\ell^t = \delta_{k\ell}$  are the non-zero eigenvalues and eigenvectors. Define the modal coefficients of the residual  $\mathbf{g}^t = \mathbf{p}^t - \mathbf{y}$  by

$$e_k^t := (\mathbf{v}_k^t)^\top \mathbf{g}^t, \quad k = 1, \dots, V-1.$$

Then there exists a constant  $C > 0$  such that for all nonzero modes  $k \geq 1$ ,

$$(\mathbf{v}_k^t)^\top \mathbf{g}^{t+1} = \left(1 - \eta \mu [\lambda_k^t + \tilde{\rho}^t (\lambda_k^t)^2]\right) e_k^t + r_k^t, \quad |r_k^t| \leq C \eta^2. \quad (9)$$

*Proof.* By Theorem A.2 (residual expansion), we have

$$\mathbf{g}^{t+1} = \left(\mathbf{I} - \eta \mu \mathbf{H}_z^t - \eta \mu \tilde{\rho}^t (\mathbf{H}_z^t)^2\right) \mathbf{g}^t + \mathbf{r}_g^t, \quad \|\mathbf{r}_g^t\| \leq C \eta^2. \quad (10)$$

Fix  $t$  and let the eigendecomposition of  $\mathbf{H}_z^t$  be  $\mathbf{H}_z^t = \sum_{k=1}^{V-1} \lambda_k^t \mathbf{v}_k^t (\mathbf{v}_k^t)^\top$  with  $\lambda_k^t > 0$  and  $\{\mathbf{v}_k^t\}_{k=1}^{V-1}$  orthonormal. (The zero mode corresponding to  $\lambda = 0$  is orthogonal to  $\mathbf{g}^t$  in the softmax-CE case and is therefore omitted.)

Project equation 10 onto the eigenvector  $\mathbf{v}_k^t$ :

$$(\mathbf{v}_k^t)^\top \mathbf{g}^{t+1} = (\mathbf{v}_k^t)^\top \left(\mathbf{I} - \eta \mu \mathbf{H}_z^t - \eta \mu \tilde{\rho}^t (\mathbf{H}_z^t)^2\right) \mathbf{g}^t + (\mathbf{v}_k^t)^\top \mathbf{r}_g^t.$$

Using the eigen-relations  $\mathbf{H}_z^t \mathbf{v}_k^t = \lambda_k^t \mathbf{v}_k^t$  and  $(\mathbf{H}_z^t)^2 \mathbf{v}_k^t = (\lambda_k^t)^2 \mathbf{v}_k^t$  and the definition  $e_k^t = (\mathbf{v}_k^t)^\top \mathbf{g}^t$ , we obtain

$$(\mathbf{v}_k^t)^\top \mathbf{g}^{t+1} = \left(1 - \eta \mu \lambda_k^t - \eta \mu \tilde{\rho}^t (\lambda_k^t)^2\right) e_k^t + r_k^t, \quad r_k^t := (\mathbf{v}_k^t)^\top \mathbf{r}_g^t.$$

Finally, since  $\|\mathbf{v}_k^t\| = 1$  we have  $|r_k^t| \leq \|\mathbf{r}_g^t\| \leq C \eta^2$ , which is exactly equation 9. This completes the proof.  $\square$

**Corollary A.5** (One-step confidence ratios under SAM). *Under the assumptions of Theorem A.2. Fix an iteration  $t$  and write  $\mathbf{p}^t = \text{softmax}(\mathbf{z}^t)$ ,  $\mathbf{g}^t = \mathbf{p}^t - \mathbf{e}_y$ , and  $\mathbf{H}_z^t = \text{diag}(\mathbf{p}^t) - \mathbf{p}^t (\mathbf{p}^t)^\top$ . For each class  $i \in [V]$ , define the one-step confidence ratio*

$$\alpha_i^\bullet := \frac{p_i^{t+1}(\bullet)}{p_i^t}, \quad \bullet \in \{\text{GD}, \text{SAM}\}.$$

Then  $\alpha_i^\bullet$  admits the representation

$$\alpha_i^\bullet = \frac{\sum_{j=1}^V e^{z_j^t}}{\sum_{j=1}^V \beta_j^\bullet e^{z_j^t}}, \quad \beta_j^{\text{GD}} = \exp\{-\eta'[(p_j^t - y_j) - (p_i^t - y_i)]\},$$

and the SAM correction appears multiplicatively as

$$\beta_j^{\text{SAM}} = \beta_j^{\text{GD}} \exp\left\{-\eta' \tilde{\rho}^t [(\mathbf{H}_z^t \mathbf{g}^t)_j - (\mathbf{H}_z^t \mathbf{g}^t)_i]\right\},$$

where  $\eta' = \eta \mu$  and, when  $\|\mathbf{g}^t\| > 0$ ,  $\tilde{\rho}^t = \rho \sqrt{\mu} / \|\mathbf{g}^t\|$  (otherwise  $\tilde{\rho}^t = 0$  by convention).

Let  $y$  be the ground-truth label and  $y^* = \arg \max_{j \neq y} p_j^t$  the most confident incorrect class.

Assume the sign condition  $\eta' \tilde{\rho}^t > 0$  and the radius scaling  $|\rho| \leq \kappa \sqrt{|\eta|}$ . Then there exists  $\eta_0 = \eta_0(\mathbf{p}^t, \mathbf{H}_z^t, \|\mathbf{g}^t\|, \mu, \kappa, L) > 0$  such that, for all step sizes  $0 < |\eta| \leq \eta_0$ , the following one-step inequalities hold without remainder terms:

$$\alpha_{y^*}^{\text{SAM}} \leq \alpha_{y^*}^{\text{GD}}, \quad \alpha_y^{\text{SAM}} \geq \alpha_y^{\text{GD}}.$$

Here  $y \in \{y^+, y^-\}$  denotes the ground-truth label corresponding to the positive or negative learning rate, respectively. Moreover, the inequalities are strict whenever  $p_{y^*}^t \in (0, 1)$  and  $p_y^t \leq \frac{1}{2}$ . In particular, when  $\tilde{\rho}^t = 0$  (no SAM), the two equalities hold.

*Proof.* Fix  $t$  and a class  $i \in [V]$ . Set  $\eta' = \eta \mu$ . By Theorem A.2 (logits line),

$$\Delta \mathbf{z} := \mathbf{z}^{t+1} - \mathbf{z}^t = -\eta' (\mathbf{g}^t + \tilde{\rho}^t \mathbf{H}_z^t \mathbf{g}^t) + \mathbf{r}_z^t, \quad \|\mathbf{r}_z^t\|_\infty \leq C_1 \eta^2,$$

where  $C_1$  depends only on  $(L, \mu, \kappa)$  and the hypothesis  $|\rho| \leq \kappa\sqrt{|\eta|}$  is in force.

For any increment  $\Delta z$ ,

$$\alpha_i = \frac{p_i(z^t + \Delta z)}{p_i(z^t)} = \frac{\sum_j e^{z_j^t}}{\sum_j \exp\{\Delta z_j - \Delta z_i\} e^{z_j^t}} = \frac{\sum_j e^{z_j^t}}{\sum_j \beta_j e^{z_j^t}}.$$

With the above  $\Delta z$ ,

$$\beta_j^{\text{SAM}} = \underbrace{\exp\{-\eta'(g_j^t - g_i^t)\}}_{\beta_j^{\text{GD}}} \underbrace{\exp\{-\eta'\tilde{\rho}^t \Delta_{j,i}^t\}}_{\text{curvature factor}} \underbrace{\exp\{r_j^t - r_i^t\}}_{\text{remainder factor}}, \quad \Delta_{j,i}^t := (\mathbf{H}_z^t \mathbf{g}^t)_j - (\mathbf{H}_z^t \mathbf{g}^t)_i.$$

From  $\|\mathbf{r}_z^t\|_\infty \leq C_1 \eta^2$ , we have  $e^{-2C_1 \eta^2} \leq \exp\{r_j^t - r_i^t\} \leq e^{2C_1 \eta^2}$  for all  $i, j$ .

With  $\mathbf{H}_z^t = \text{diag}(\mathbf{p}^t) - \mathbf{p}^t(\mathbf{p}^t)^\top$  and  $\mathbf{g}^t = \mathbf{p}^t - \mathbf{e}_y$ ,

$$(\mathbf{H}_z^t \mathbf{g}^t)_i = p_i^t(p_i^t - y_i - C^t), \quad C^t := \sum_k (p_k^t)^2 - p_y^t.$$

Let  $y^* = \arg \max_{j \neq y} p_j^t$ . Then  $C^t \leq p_{y^*}^t$  and one checks: (i) for  $i = y^*$ ,  $\Delta_{j,y^*}^t \leq 0$  for all  $j$ , and  $\Delta_{j,y^*}^t < 0$  for some  $j$  whenever  $p_{y^*}^t \in (0, 1)$ ; (ii) for  $i = y$ ,  $\Delta_{j,y}^t \geq 0$  for all  $j$  whenever  $p_y^t \leq \frac{1}{2}$ , and  $\Delta_{j,y}^t > 0$  for some  $j$  if  $p_y^t \in (0, \frac{1}{2}]$ .

Define

$$D_i^{\text{GD}} := \sum_j \beta_j^{\text{GD}} e^{z_j^t}, \quad \tilde{D}_i := \sum_j \beta_j^{\text{GD}} e^{z_j^t} \exp\{-\eta'\tilde{\rho}^t \Delta_{j,i}^t\}, \quad D_i^{\text{SAM}} := \sum_j \beta_j^{\text{SAM}} e^{z_j^t}.$$

By the remainder bounds,  $e^{-2C_1 \eta^2} \tilde{D}_i \leq D_i^{\text{SAM}} \leq e^{2C_1 \eta^2} \tilde{D}_i$ . Next, by  $e^x \geq 1 + x$  and the sign structure of  $\Delta_{j,i}^t$ ,

$$\frac{\tilde{D}_{y^*}}{D_{y^*}^{\text{GD}}} = \sum_j w_j^{(y^*)} e^{-\eta'\tilde{\rho}^t \Delta_{j,y^*}^t} \geq 1 + \eta'\tilde{\rho}^t \sum_j w_j^{(y^*)} (-\Delta_{j,y^*}^t) \geq 1 + c_{y^*} \eta'\tilde{\rho}^t,$$

for some  $c_{y^*} > 0$  whenever  $p_{y^*}^t \in (0, 1)$ ; here  $w_j^{(i)} := \beta_j^{\text{GD}} e^{z_j^t} / D_i^{\text{GD}}$  are positive weights. Similarly, for  $p_y^t \leq \frac{1}{2}$ ,

$$\frac{\tilde{D}_y}{D_y^{\text{GD}}} = \sum_j w_j^{(y)} e^{-\eta'\tilde{\rho}^t \Delta_{j,y}^t} \leq 1 - \eta'\tilde{\rho}^t \sum_j w_j^{(y)} \Delta_{j,y}^t \leq 1 - c_y \eta'\tilde{\rho}^t$$

for some  $c_y > 0$  (strict in the stated nondegenerate case).

Now use the scaling  $|\rho| \leq \kappa\sqrt{|\eta|}$ : then  $\eta'\tilde{\rho}^t = \Theta(\eta^{3/2})$ , whereas  $e^{\pm 2C_1 \eta^2} = 1 \pm O(\eta^2)$ . Hence there exists  $\eta_0 > 0$  (depending only on  $(\mathbf{p}^t, \mathbf{H}_z^t, \|\mathbf{g}^t\|, \mu, \kappa, L)$ ) such that for  $0 < |\eta| \leq \eta_0$ ,

$$D_{y^*}^{\text{SAM}} \geq e^{-2C_1 \eta^2} \tilde{D}_{y^*} \geq D_{y^*}^{\text{GD}} (1 + \frac{1}{2} c_{y^*} \eta' \tilde{\rho}^t), \quad D_y^{\text{SAM}} \leq e^{2C_1 \eta^2} \tilde{D}_y \leq D_y^{\text{GD}} (1 - \frac{1}{2} c_y \eta' \tilde{\rho}^t).$$

Since  $\alpha_i = (\sum_j e^{z_j^t}) / D_i$ , we obtain for  $0 < |\eta| \leq \eta_0$ :

$$\alpha_{y^*}^{\text{SAM}} \leq \alpha_{y^*}^{\text{GD}}, \quad \alpha_y^{\text{SAM}} \geq \alpha_y^{\text{GD}},$$

with strict inequalities under the stated nondegeneracy conditions (because then  $c_{y^*}, c_y > 0$ ). If  $\|\mathbf{g}^t\| = 0$  (so  $\tilde{\rho}^t = 0$  by convention), both become equalities. This completes the proof.  $\square$

## B IMPLEMENTATION

---

### Algorithm 1 Logits-SAM pseudocode

---

**Require:** model, batch,  $\rho$

- 1: Let  $W \leftarrow \text{lm.head.weight}$
  - 2: Run forward to get `loss_pre` and hidden states  $H$
  - 3:  $g \leftarrow \text{grad}(\text{loss\_pre}, W)$
  - 4:  $e \leftarrow \rho g / \|g\|_2$
  - 5: `logits_perturbed`  $\leftarrow \text{linear}(H, W + e)$
  - 6: Compute `loss_post` with `logits_perturbed`
  - 7: Backward `loss_post`
-

Table 5: Comparison between theoretical and practical settings of DPO with SAM. Although the signs differ for  $y^-$ , the resulting dynamics are equivalent. For  $y^+$ , the settings coincide.

Class	Objective	Learning rate	$\rho$	Setting
$y^+$	Positive objective $f = -\log p$	Positive ( $\eta > 0$ )	Positive	Theory = Practice
$y^-$ (Theory)	Positive objective $f = -\log p$	Negative ( $\eta < 0$ )	Negative	Theory
$y^-$ (Practice)	Negative objective $f = \log p$	Positive ( $\eta > 0$ )	Positive	Practice

**Equivalence of sign conventions for  $y^-$ .** Theoretical setting ( $f = -\log p$ ,  $\eta^- < 0$ ,  $\rho^- < 0$ ):

$$\theta_{t+1}^{\text{theory}} = \theta_t - \eta^- \nabla f(\theta_t).$$

Practical setting uses  $\tilde{f} = -f$ ,  $\tilde{\eta} = -\eta^- > 0$ :

$$\theta_{t+1}^{\text{prac}} = \theta_t - \tilde{\eta} \nabla \tilde{f}(\theta_t) = \theta_t - (-\eta^-)(-\nabla f(\theta_t)) = \theta_{t+1}^{\text{theory}}.$$

For SAM, theoretical perturbation and update:

$$\epsilon_t^{\text{theory}} = \rho^- \frac{\nabla f(\theta_t)}{\|\nabla f(\theta_t)\|}, \quad \theta_{t+1}^{\text{theory}} = \theta_t - \eta^- \nabla f(\theta_t + \epsilon_t^{\text{theory}}).$$

Practical setting uses  $\tilde{f} = -f$ ,  $\tilde{\rho} = -\rho^- > 0$ ,  $\tilde{\eta} = -\eta^- > 0$ :

$$\epsilon_t^{\text{prac}} = \tilde{\rho} \frac{\nabla \tilde{f}(\theta_t)}{\|\nabla \tilde{f}(\theta_t)\|} = \rho^- \frac{\nabla f(\theta_t)}{\|\nabla f(\theta_t)\|} = \epsilon_t^{\text{theory}},$$

$$\theta_{t+1}^{\text{prac}} = \theta_t - \tilde{\eta} \nabla \tilde{f}(\theta_t + \epsilon_t^{\text{prac}}) = \theta_{t+1}^{\text{theory}}.$$

## C ADDITIONAL EXPERIMENTAL DETAILS

**Benchmark details.** **AlpacaEval 2** (Dubois et al., 2024) is a large-scale preference benchmark for open-ended instruction following that uses LLM-as-a-judge calibrated to human preferences; its evaluation set contains 805 single-turn instructions, and models are typically compared in pairwise settings against a baseline. **Arena-Hard v0.1** (Li et al., 2024) is a challenging subset of difficult user instructions mined from Chatbot Arena; it enables fine-grained, head-to-head comparisons between models via pairwise judging and comprises 500 hard prompts. **MT-Bench** (Zheng et al., 2023) is a multi-turn dialogue benchmark that tests a model’s ability to handle diverse conversational tasks across several categories; the standard evaluation set consists of 80 multi-turn questions.

**Training details.** For experiments on Pythia-2.8B, we use two NVIDIA A100 GPUs with data-parallel training under DDP; for Mistral-7B, we use four NVIDIA A100 GPUs with parallel training via DeepSpeed ZeRO-3 (Rasley et al., 2020).

## D LLM USAGE STATEMENT

In preparing this manuscript, we employed a large language model (LLM) as an auxiliary tool. Specifically, the LLM was used to assist with proofreading, formatting, and grammar checking of the text.

Method	Objective	Hyperparameter
SLiC-HF	$\max(0, \delta - \log \pi_\theta(y_w   x) + \log \pi_\theta(y_l   x)) - \lambda \log \pi_\theta(y_w   x)$	$\lambda \in \{0.1, 0.5, 1.0, 10.0\};$ $\delta \in \{0.5, 1.0, 2.0, 10.0\}$
CPO	$-\log \sigma(\beta \log \pi_\theta(y_w   x) - \beta \log \pi_\theta(y_l   x)) - \lambda \log \pi_\theta(y_w   x)$	$\lambda = 1.0; \beta \in \{0.01, 0.05, 0.1\}$

Table 6: Objectives and hyperparameters for SLiC-HF and CPO.

Method	Pythia-2.8B	Mistral-7B
DPO	$1 \times 10^{-4}$	$1 \times 10^{-5}$
SLiC-HF	$1 \times 10^{-3}$	$1 \times 10^{-4}$
CPO	$1 \times 10^{-4}$	$1 \times 10^{-5}$

Table 7: Choice of  $\rho$  for logits-SAM.