

---

# Flatness and Gradient Alignment Are Both Necessary: Spectral-Aware Gradient-Aligned Exploration for Multi-Distribution Learning

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Sharpness-aware and gradient-alignment methods have been shown to improve  
2 generalization, however each family of methods targets a single geometric property  
3 of the loss landscape, while ignoring the other. In this paper, we show that this  
4 omission is structurally unavoidable and that both *flatness* and *gradient alignment*  
5 should be considered in multi-distribution learning settings. Specifically, we derive  
6 an excess-risk decomposition that yields two additive leading-order terms: (i)  
7 an alignment term, controlled by the trace of  $\bar{H}^{-1}\Sigma_g$  and (ii) a curvature term,  
8 controlled by  $\bar{H}$ , where  $\bar{H}$  is the average Hessian and  $\Sigma_g$  is the covariance of the  
9 gradient across distributions. Notably,  $\bar{H}$  appears inverted in one and non-inverted  
10 in the other. We further show, via a counterexample, that neither quantity bounds  
11 the other in general, so no algorithm targeting only one term can guarantee low  
12 excess risk. Motivated by this decomposition, we propose SAGE (Spectral-Aware  
13 Gradient-Aligned Exploration) that targets both terms. The curvature component  
14 replaces SAM’s gradient-scaled perturbation with the polar factor of each layer’s  
15 gradient matrix, computed via Newton–Schulz iteration, so that the ascent step  
16 probes all directions with similar magnitude. On the other hand, the alignment  
17 component injects isotropic noise at the descent step, the magnitude of which  
18 scales with cross-distribution gradient disagreement. Experiments on five domain-  
19 generalization and two multi-task learning benchmarks show that the proposed  
20 method establishes a new state-of-the-art on DomainBed and acts as a general-  
21 purpose improvement to base MTL solvers, remaining competitive with, or even  
22 surpassing, state-of-the-art methods<sup>1</sup>.

## 23 1 Introduction

24 Understanding the geometry of loss landscapes [1, 2] in overparameterized neural networks has  
25 become a central lens for interpreting optimization behavior and why certain solutions generalize  
26 better than others [3, 1, 2, 4–6]. Specifically, two geometric properties have emerged as particularly  
27 influential: the *flatness* [7–11] of the loss surface around a minimum, and the *alignment* [12–18]  
28 of gradients computed across different training distributions, e.g., domains in the case of Domain  
29 Generalization (DG) [19, 20], or tasks, in Multi-Task Learning (MTL) [21].

30 Even though both flatness and gradient alignment have been studied extensively in isolation, existing  
31 methods typically pursue each property as a standalone objective [7, 16, 12, 13, 22, 23, 15], while  
32 their interplay has received limited formal treatment [24, 9]. Intuitively, a flat minimum may sit in a  
33 region where distribution-specific gradients conflict, while an aligned minimum may lie in a sharp

---

<sup>1</sup>Code provided as supplemental material and will be published upon acceptance.

34 basin where small perturbations inflate the loss. We argue that neither property alone is sufficient  
 35 and that robust generalization benefits from minima that are simultaneously flat across directions and  
 36 exhibit improved gradient agreement across training sub-objectives.

37 To justify this claim, we derive an excess-risk decomposition (Theorem 1) that yields two additive  
 38 leading-order terms: (i) an alignment term scaling as  $\text{tr}(\bar{H}^{-1}\Sigma_g)$ , where  $\bar{H}$  is the average curvature  
 39 and  $\Sigma_g$  is the gradient covariance matrix, across distributions and (ii) a curvature term scaling as  
 40  $\text{tr}(\bar{H})$ .  $\bar{H}$  appears *inverted* in the first term and *non-inverted* in the other and, as a result, targeting  
 41 only one of the two during optimization can leave the other unconstrained. We formalize this  
 42 independence in Counterexample 1, which shows that neither quantity bounds the other, even in the  
 43 simple quadratic family.

44 Motivated by this decomposition, and based on the fact that first-order optimization methods (in-  
 45 cluding first-order flatness-aware methods) ignore curvature, we propose **SAGE** (Spectral-Aware  
 46 Gradient-Aligned Exploration) a method that jointly targets both curvature and gradient agreement  
 47 via principled heuristics. For the curvature term ( $\text{tr}(\bar{H})$ ), SAGE replaces SAM’s [7] gradient-scaled  
 48 ascent perturbation with the polar factor  $UV^T$  of each layer’s gradient matrix, computed efficiently  
 49 via Newton–Schulz iteration [25], scaled by the layer Frobenius norm so that all spectral directions  
 50 are probed with comparable magnitude <sup>2</sup> (details in Section 3). For the alignment term, SAGE  
 51 injects isotropic Gaussian noise at the descent step with magnitude scaling with cross-distribution  
 52 gradient disagreement, biasing the optimizer away from high-conflict regions. Empirical evaluation  
 53 on DG [19, 20] and MTL [21] benchmarks indicates that SAGE compares favorably to both  
 54 sharpness-aware and gradient-alignment baselines across seven datasets.

55 Our key contributions can be summarized as follows:

- 56 • We derive an excess-risk decomposition for multi-distribution learning (Theorem 1) that  
 57 splits the leading-order risk into two additive terms, an alignment term scaling as  $\text{tr}(\bar{H}^{-1}\Sigma_g)$   
 58 and a curvature term scaling as  $\text{tr}(\bar{H})$ , with  $\bar{H}$  appearing inverted in one and non-inverted  
 59 in the other.
- 60 • We provide a separation result (Counterexample 1) showing that the two terms are independ-  
 61 ently controllable for quadratic losses, showing that even in a simple scenario, targeting  
 62 each property in isolation may prove insufficient.
- 63 • Based on the above, we propose SAGE, a method targeting both terms of the decompo-  
 64 sition; a spectral perturbation (curvature) and a gradient alignment-driven noise injection  
 65 (alignment).
- 66 • We evaluate our proposed method on five DG and two MTL benchmarks, demonstrating  
 67 competitive results against both sharpness-aware and gradient-alignment baselines.

## 68 2 Theoretical Analysis: Flatness and Gradient Alignment in 69 Multi-Distribution Learning

### 70 2.1 Setup and Notation

71 Let  $\mathcal{X}$  be the input space and  $\mathcal{Y}$  the label space. Let  $e$  be a joint distribution  $P_{\mathcal{X}\mathcal{Y}}^e$  over  $\mathcal{X} \times \mathcal{Y}$   
 72 (denoted as *environment* in [30]). For example, different distributions/environments may correspond  
 73 to different *domains* in the domain generalization literature or *tasks* in the multi-task literature. Let  $\mathcal{E}$   
 74 be a set of distributions and  $\mathcal{P}$  a meta-distribution over  $\mathcal{E}$ . For each  $e \in \mathcal{E}$ , let  $\mathcal{L}_e : \Theta \rightarrow \mathbb{R}$  denote a  
 75 per-distribution loss, where  $\Theta \subseteq \mathbb{R}^d$  is an open parameter space. The *population risk* is:

$$R(\theta) := \mathbb{E}_{e \sim \mathcal{P}}[\mathcal{L}_e(\theta)], \quad (1)$$

76 and the *empirical multi-distribution risk* over  $K$  iid training distributions  $e_1, \dots, e_K \sim \mathcal{P}$  is

$$\hat{R}(\theta) := \frac{1}{K} \sum_{k=1}^K \mathcal{L}_{e_k}(\theta). \quad (2)$$

---

<sup>2</sup>We do not explicitly address the decomposition via second-order approximations [9, 26, 27] to avoid additional computational burden. In addition, the same orthogonalization in the gradient matrix singular basis has been used in recent optimizers [28, 29] and corresponds to the steepest-descent step under the spectral norm [25].

77 This is the standard objective minimized by ERM. At any base point  $\theta_0 \in \Theta$  we introduce the  
 78 per-distribution gradient and Hessian,

$$g_e := \nabla_{\theta} \mathcal{L}_e(\theta_0), \quad H_e := \nabla_{\theta}^2 \mathcal{L}_e(\theta_0), \quad (3)$$

79 and the following cross-distribution statistics,

$$\bar{g} := \mathbb{E}_e[g_e], \quad \bar{H} := \mathbb{E}_e[H_e], \quad \Sigma_g := \text{Cov}_e(g_e) = \mathbb{E}_e[(g_e - \bar{g})(g_e - \bar{g})^T]. \quad (4)$$

80  $\bar{H}$  is the average curvature targeted by sharpness-aware methods,  $\Sigma_g$  is the cross-distribution gradient  
 81 covariance targeted (implicitly or explicitly) by gradient-alignment methods, and  $\bar{g}$  is the expected  
 82 gradient direction, vanishing ( $\bar{g} = 0$ ) at stationary points of  $R$ .

83 **Assumption 1** (Regularity). Each  $\mathcal{L}_e$  is three times continuously differentiable (i.e.,  $C^3$ ) in  $\theta$ , with  
 84 third derivatives uniformly bounded in a neighborhood of  $\theta^* := \arg \min_{\theta} R(\theta)$ , and  $\bar{H}(\theta^*) \succ 0$   
 85 (Positive-Definite). The gradients  $g_e(\theta^*)$  have finite second moments under  $\mathcal{P}$ .

## 86 2.2 Excess-Risk Decomposition

87 The main result of our analysis is a decomposition of the expected excess risk of the first-order  
 88 empirical minimizer into two additive terms, one of which depends on  $\Sigma_g$  through  $\bar{H}^{-1}$  and the other  
 89 of which depends on  $\bar{H}$  directly.

**Theorem 1** (Multi-distribution excess-risk decomposition). *Under Assumption 1, let  $\hat{\theta}$  denote the minimizer of  $\hat{R}$  over  $K$  iid training distributions, and let  $\xi \sim \mathcal{N}(0, \sigma^2 I)$  denote isotropic Gaussian noise of scale  $\sigma > 0$ . Then, for a parameter  $\theta = \hat{\theta} + \xi$ :*

$$\mathbb{E}[\mathbb{E}_{\xi}[R(\theta)]] - R(\theta^*) = \underbrace{\frac{1}{2K} \text{tr}(\bar{H}^{-1} \Sigma_g)}_{\text{alignment term}} + \underbrace{\frac{\sigma^2}{2} \text{tr}(\bar{H})}_{\text{curvature term}} + O(K^{-3/2}) + O(\sigma^3), \quad (5)$$

where the outer expectation is over the iid sampling of training distributions,  $\bar{H} := \nabla^2 R(\theta^*)$ , and  $\Sigma_g := \text{Cov}_e(\nabla \mathcal{L}_e(\theta^*))$ .

90

91 **Proof sketch.** To prove Theorem 1, we expand  $R$  quadratically around  $\theta^*$ , where the linear term  
 92 vanishes since  $\bar{g}(\theta^*) = 0$ . We then express  $\hat{\theta} - \theta^* = -\bar{H}^{-1} \hat{g} + O_P(K^{-1})$  via the implicit-  
 93 function argument standard in M-estimator asymptotics [31]. Substituting into the quadratic and  
 94 taking expectation over the iid distribution sampling, with  $\text{Cov}(\hat{g}) = \frac{1}{K} \Sigma_g$ , yields the alignment  
 95 term. Separately, expanding  $R$  around  $\hat{\theta}$  under random perturbations  $\xi \sim \mathcal{N}(0, \sigma^2 I)$ , e.g. due to  
 96 optimization noise, produces the curvature term. Full proof in Appendix B.2.

97 **Remark 1.** The two terms exhibit a structural trade-off, i.e.,  $\bar{H}$  appears inverted in the alignment  
 98 term and non-inverted in the curvature term. As a result, minimizing one term by reshaping the  
 99 Hessian strictly inflates the other. Considered in isolation, the alignment term implies that excess risk  
 100 is minimized by maximizing  $\bar{H}$ . The curvature term mitigates this by explicitly penalizing excessive  
 101 sharpness, ensuring a balanced optimization objective.

## 102 2.3 Decoupling of Flatness and Alignment

103 Theorem 1 gives an additive decomposition into a flatness-controlled term and an alignment-controlled  
 104 term. The natural next question is whether the two terms are linked, i.e. an inequality of the form  
 105 “small  $\text{tr}(\bar{H})$  implies small  $\text{tr}(\bar{H}^{-1} \Sigma_g)$ ” that would allow a single criterion to bound both. The  
 106 following counterexample shows that such link may not exist, even in simple settings.

107 **Counterexample 1** (Decoupling of flatness and gradient alignment). Consider the family of multi-  
 108 distribution learning problems with per-distribution losses

$$\mathcal{L}_e(\theta) = \frac{1}{2} \theta^T A \theta + b_e^T \theta, \quad \mathbb{E}_e[b_e] = 0, \quad (6)$$

109 parameterized by a symmetric PD matrix  $A \in \mathbb{R}^{d \times d}$  and a collection of linear coefficients  $\{b_e\}$   
 110 satisfying  $\mathbb{E}_e[b_e] = 0$ . For every  $M > 0$ , there exist instances of Eq. (6) such that:

111 (i) **Flat but misaligned:**  $\text{tr}(\bar{H}) \leq M^{-1}$  and  $\text{tr}(\bar{H}^{-1}\Sigma_g) \geq M$ .

112 (ii) **Aligned but sharp:**  $\text{tr}(\bar{H}^{-1}\Sigma_g) \leq M^{-1}$  and  $\text{tr}(\bar{H}) \geq M$ .

113 In particular, neither  $\text{tr}(\bar{H})$  nor  $\text{tr}(\bar{H}^{-1}\Sigma_g)$  can be bounded above by any function of the other alone.

114 **Proof sketch.** For the quadratic family,  $\bar{H} = A$  depends only on  $A$  and  $\Sigma_g = \mathbb{E}_e[b_e b_e^T]$  depends  
 115 only on  $\{b_e\}$ , so the two quantities can be specified independently. Setting  $A = \lambda I$  and choosing  $\{b_e\}$   
 116 supported on a single direction with prescribed variance yields both constructions, by an appropriate  
 117 choice of  $\lambda$ . Full proof details in Appendix B.3.

118 **Remark 2.** This decoupling result in the quadratic setting indicates that an algorithm targeting only  
 119 one of the two quantities may increase the generalization error resulting from the other. Thus, both  
 120 flatness and alignment terms need to be considered during optimization.

121 To make Theorem 1 and Counterexample 1 concrete, we construct in Appendix C a two-dimensional  
 122 learning problem in which the support of  $\Sigma_g$  coincides with the flat eigendirection of  $\bar{H}$ . The example  
 123 realizes and exhibits two distinct failure modes along orthogonal eigendirections.

### 124 3 Spectral-Aware Gradient-Aligned Exploration

125 In this section, we describe SAGE, a single algorithm that targets the excess-risk decomposition terms  
 126 of Theorem 1 via two separate principled heuristic mechanisms. The two components introduced  
 127 below are tractable surrogates designed to probe each term and are introduced to avoid 2nd-order  
 128 computational overhead.

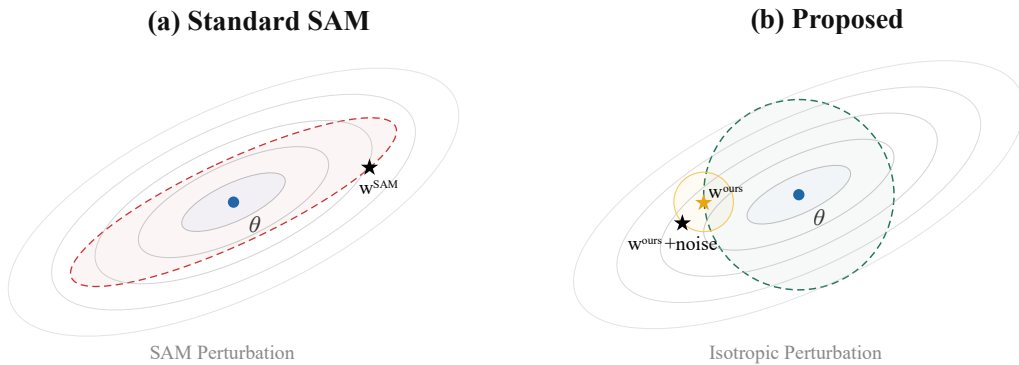


Figure 1: **Comparison of standard SAM and our proposed method.** (a) Standard SAM computes the perturbation as  $\epsilon_{\text{SAM}} = \rho \nabla_{\theta} \mathcal{L} / \|\nabla_{\theta} \mathcal{L}\|_2$ , determined solely by the gradient direction (dashed ellipse), and the descent step lands at  $w^{\text{SAM}}$  (★). (b) We replace the gradient with its orthogonal polar factor ( $G \rightarrow UV^T$ ), setting all singular values to one and scaling by the layer norm, yielding a perturbation  $\epsilon' = \rho \|W\|_F \cdot UV^T$  that probes all directions with equal magnitude (dashed circle). After the descent step lands at  $w^{\text{ours}}$  (★), isotropic noise scaled by cross-distribution gradient disagreement perturbs the update (orange circle), discouraging convergence to regions of gradient conflict.

#### 129 3.1 Targeting the Curvature Term

130 Motivated by the curvature term in Theorem 1 we follow the logic of SAM [7] and SAM-like [32, 24,  
 131 8, 22, 10, 11] methods, which have empirically proven to be effective in multi-distribution settings,  
 132 but propose an alternative for the ascent perturbation step. Specifically, SAM-like methods derive the  
 133 perturbation by approximating the inner maximization  $\max_{\|\epsilon\| \leq \rho} \mathcal{L}(\theta + \epsilon)$  with a first-order Taylor  
 134 expansion,  $\mathcal{L}(\theta + \epsilon) \approx \mathcal{L}(\theta) + \epsilon^T \nabla \mathcal{L}(\theta)$ , whose maximizer over the  $L_2$  ball admits the closed form  
 135  $\epsilon_{\text{SAM}} = \rho \nabla \mathcal{L} / \|\nabla \mathcal{L}\|_2$ . As apparent, the above perturbation direction ignores curvature information,  
 136 is determined by the gradient vector alone, and is therefore probed by dominating gradients during

137 training<sup>3</sup>. To address the above limitation, we propose replacing SAM’s perturbation with an isotropic  
 138 perturbation during the ascent step, inspired by preconditioning methods such as Shampoo [29]  
 139 and Muon [28]. These methods have empirically shown to achieve effective preconditioning by  
 140 orthogonalizing each  $2D$  weight (i.e., parameter) matrix (e.g., layer) of a network independently. We  
 141 follow the same heuristic to expand the exploration of additional perturbation directions of SAM.

142 **Spectral perturbation via orthogonalized gradients.** For each weight matrix  $W$  with gradient  
 143  $G = \nabla_W \mathcal{L}$ , we compute the orthogonal polar factor  $Q = UV^T$ , where  $G = U\Sigma V^T$  is the SVD of  $G$ .  
 144 This amounts to setting every singular value of the gradient to one, so that  $Q = \sum_{i=1}^r u_i v_i^T$  retains  
 145 the directional structure of  $G$ , i.e. the left and right singular vectors, while equalizing the magnitude  
 146 across all singular directions. A perturbation along  $Q$  therefore probes the loss surface with equal  
 147 magnitude in every spectral direction of the gradient. Since computing the exact polar factor via SVD  
 148 adds a computational burden of  $\mathcal{O}(\min\{m, n\} \cdot mn)$  per weight matrix, we instead approximate  $Q$   
 149 with the Newton–Schulz iteration [25, 35–37], which converges cubically to the orthogonal factor  
 150 using only matrix multiplications:

$$X_0 = \frac{G}{\|G\|_F}, \quad X_{k+1} = \frac{1}{2} X_k (3I - X_k^T X_k). \quad (7)$$

151 After  $T$  iterations ( $T=5$  suffices in practice),  $X_T \approx UV^T$ . Because each step involves only two  
 152 matrix products of the same shape as  $G$ , the overhead on modern GPU hardware is negligible. For  
 153 bias vectors and other one-dimensional parameter groups, where matrix polar decomposition is not  
 154 applicable, we fall back to the standard  $L_2$ -normalized perturbation  $\epsilon = \rho g / \|g\|_2$ , consistent with  
 155 vanilla SAM.

156 **Scale-adaptive radius.** Furthermore, to address the limitation from the well-documented interaction  
 157 between SAM’s fixed perturbation radius and the scale invariance of modern architectures [34], we  
 158 scale the orthogonalized radius by the Frobenius norm of each weight matrix  $\|W\|_F$ , following  
 159 the intuition of adaptive sharpness [8]. The final resulting spectral perturbation for weight matrix  
 160  $W$  is therefore  $\epsilon_W^{\text{spec}} = \rho \|W\|_F UV^T$ , where  $U, V$  are from the SVD of  $G = \nabla_W \mathcal{L}$ . Under this  
 161 formulation, every singular direction receives a perturbation of identical magnitude  $\rho \|W\|_F$ , and the  
 162 perturbation rescales with the current parameter norm, ensuring that the sharpness measure remains  
 163 invariant to reparameterizations of the form  $W \mapsto \alpha W$ . More details in Appendix D.

### 164 3.2 Targeting the Alignment term

165 Motivated by the alignment term in Theorem 1, we introduce a gradient-agreement-based regular-  
 166 izer at the descent step. Since the cross-distribution covariance  $\Sigma_g$  is intractable to estimate and  
 167 differentiate through at the scale of modern networks, we replace it with a tractable scalar proxy of  
 168 cross-distribution gradient *disagreement* and use it to modulate noise that biases the optimizer away  
 169 from high-conflict regions, rather than minimizing  $\text{tr}(\bar{H}^{-1} \Sigma_g)$  directly.

170 In several settings, such as domain generalization [19, 20, 38], multi-task learning [21], the training  
 171 signal decomposes into  $K$  distinct distributions, each corresponding to a different domain, data  
 172 source, or task. The overarching goal is to learn a model that performs well on all sub-tasks  
 173 and unseen OOD data. Intuitively, *directions of parameter update that are beneficial across all*  
 174 *distributions are more likely to generalize to new ones*. This intuition has been formalized under  
 175 gradient alignment [16, 12–14], where it has been shown that under certain assumptions [15] it  
 176 is a necessary condition for domain-invariance. Specifically, when gradients of different tasks or  
 177 domains point in similar directions, the shared gradient direction captures features that are invariant  
 178 across distributions. On the contrary, when they conflict the model is being pulled toward domain- or  
 179 task-specific shortcuts.

180 Formally, as defined in Section 2.1, given distribution loss functions  $\mathcal{L}_{e_k}(\theta)$  for  $K$  source distributions  
 181  $e_1, \dots, e_K \sim \mathcal{P}$ , where in our case a distribution or environment  $e_k$  may represent a domain, a task,  
 182 or any other meaningful partition of the training data, we measure the agreement among gradients  
 183  $g_{e_k} = \nabla_{\theta} \mathcal{L}_{e_k}(\theta)$  using the average pairwise cosine similarity,  $S(\theta) = \frac{2}{K(K-1)} \sum_{i < j} \frac{\langle g_{e_i}, g_{e_j} \rangle}{\|g_{e_i}\| \|g_{e_j}\|}$ .

<sup>3</sup>Even though iterated SAM dynamics have been shown to implicitly minimize  $\lambda_{\max}(\bar{H})$  [33], it is a different functional of the Hessian from the trace  $\text{tr}(\bar{H})$  identified by Theorem 1 as the relevant target.

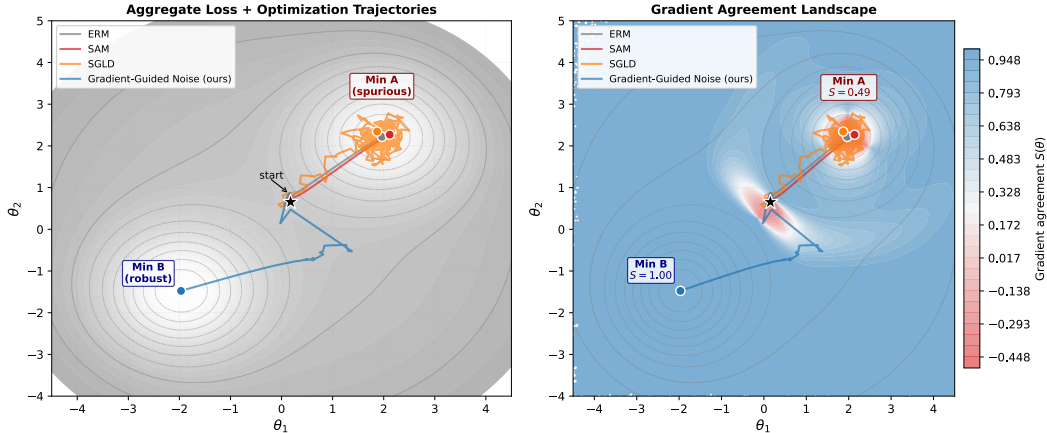


Figure 2: *Left*: Aggregate loss contours with optimization trajectories from a shared starting point ( $\star$ ). ERM (grey), SAM (red), and SGLD (orange) converge to Minimum A, whereas only gradient-aligned noise (blue, ours) escapes to Minimum B. *Right*: The same trajectories overlaid on the gradient agreement landscape  $S(\theta)$ , where blue and red indicate high and low cross-domain gradient agreement, respectively.

184 When  $S(\theta) \approx 1$ , all distributions agree on the update direction, whereas when  $S(\theta) \leq 0$ , gradients  
 185 conflict, indicating that the optimizer is navigating a region where domain- or objective-specific fea-  
 186 tures dominate. Rather than modifying the gradient direction directly (as in PCGrad [16]) and to avoid  
 187 further computational burden via second-order approximation techniques, we adopt a noise-injection  
 188 strategy inspired by the recently proposed Gradient-Guided Annealing (GGA) framework [15]. We  
 189 inject isotropic Gaussian noise into the gradient at the descent step, with magnitude inversely propor-  
 190 tional to the gradient agreement  $\beta = \gamma(1 - S(\theta))$ , where  $\gamma > 0$  is a hyperparameter controlling the  
 191 maximum noise scale. The updated gradient then becomes  $g' = \nabla_{\theta} \mathcal{L}(\theta + \epsilon) + \beta \xi$ ,  $\xi \sim \mathcal{N}(0, I)$ .

192 The intuition behind this design is twofold. When gradients are well-aligned ( $S \approx 1$ ) the additive  
 193 noise vanishes, whereas when gradients conflict ( $S \ll 1$ ) the noise discourages update directions  
 194 that benefit some distributions at the expense of others. We illustrate the behavior between ERM,  
 195 SAM [7], SGLD<sup>4</sup> [39, 40], and the proposed noise mechanism in Fig. 2.

### 196 3.3 Spectral-Aware Gradient-Aligned Exploration

197 We now describe the complete SAGE algorithm, which unifies the spectral perturbation and gradient-  
 198 alignment into a single training procedure. At each iteration, SAGE: (i) computes per-distribution  
 199 gradients  $g_{e_k}$  and aggregate gradient from a single forward pass, (ii) evaluates the gradient agreement  
 200  $S(\theta)$  and sets the noise scale  $\beta = \gamma(1 - S(\theta))$ , (iii) constructs the spectral perturbation by orthogo-  
 201 nalizing each weight matrix’s gradient via the Newton-Schulz iteration and scaling by  $\rho \|W\|_F$ , (iv)  
 202 performs a forward-backward pass at the perturbed point  $\theta + \epsilon$ , and (v) restores the clean weights,  
 203 adds isotropic gradient alignment-driven Gaussian noise  $\beta \mathcal{N}(0, I)$  to the adversarial gradient, and  
 204 finally applies the base optimizer update step. A detailed summary in pseudocode of the SAGE  
 205 algorithm is provided in Appendix H, whereas a conceptual illustration is provided in Fig. 1.

## 206 4 Related Work

207 In this section we situate SAGE within the literature of DG and MTL, and provide an overview of the  
 208 most related works. An extended discussion regarding the landscape of the two fields is provided in  
 209 Appendix A.

210 **Sharpness-aware minimization.** The observation that flat minima generalize better than sharp  
 211 ones [5] motivated Sharpness-Aware Minimization (SAM) [7], which performs a gradient ascent step

<sup>4</sup>SGLD or Stochastic Gradient Langevin Dynamics adds constant isotropic noise during optimization.

212 followed by a descent step so that the optimizer preferentially converges to flat regions of the loss  
 213 landscape. Several refinements have since been proposed, such as GSAM [32] which introduces a  
 214 surrogate gap objective, while ASAM [8] adapts perturbation magnitudes to the parameter scale, and  
 215 SAGM [24] additionally encourages the clean-loss and perturbed-loss gradients to remain aligned,  
 216 jointly minimizing sharpness and the surrogate gap. In the DG setting, DISAM [10] showed that  
 217 SAM’s perturbation is biased toward the domain with the largest gradient and proposed a domain-  
 218 loss variance constraint to counteract this effect, while SWAD [73] achieves flatness passively by  
 219 averaging model weights along the training trajectory. In MTL, SAMO [11] combines global and  
 220 per-task perturbation information through a zeroth-order estimator. An important note is the majority  
 221 of all sharpness-aware algorithms incorporate the perturbation method of the initially proposed  
 222 SAM [7].

223 **Gradient-alignment methods.** A complementary line of work uses agreement among per-  
 224 environment or per-task gradients as an indicator of domain invariance. In DG, Fish [13] ap-  
 225 proximately maximizes the inner product of per-domain gradients, while Mansilla et al. [12] adapt  
 226 gradient surgery [16] to mute gradient dimensions where domains disagree in sign. More recently,  
 227 GGA [15] searches for parameter configurations with well-aligned domain gradients via a simulated-  
 228 annealing-inspired procedure. In MTL, PCGrad [16] projects conflicting task gradients, CAGrad [17]  
 229 maximizes worst-case per-task improvement, and Nash-MTL [48] frames the problem as a bargaining  
 230 game. FAMO [47] circumvents the  $\mathcal{O}(K)$  per-iteration cost of gradient manipulation by working in  
 231 log-loss space with amortized weight updates. However, these methods enforce gradient alignment  
 232 without considering the curvature of the resulting solution.

233 **Positioning of SAGE.** The most closely related prior work is SAGM [24], which also augments  
 234 SAM with a gradient-matching term. However, SAGM’s matching aligns the clean-loss and perturbed-  
 235 loss gradients, both computed on the *aggregate* data, and does not involve per-environment gradient  
 236 information. SAGE differs in two respects: (i) its spectral perturbation replaces SAM’s gradient-  
 237 scaled ascent with the polar factor of each layer’s gradient, computed via Newton–Schulz iteration and  
 238 scaled by the layer’s Frobenius norm, yielding a reparameterization-invariant sharpness probe; and (ii)  
 239 its noise injection at the descent step is scaled by the degree of cross-environment gradient conflict,  
 240 targeting the alignment term  $\text{tr}(\bar{H}^{-1}\Sigma_g)$  identified in Theorem 1. Unlike methods that modify the  
 241 gradient direction (PCGrad, Fish) or solve auxiliary optimization problems (Nash-MTL, FAMO),  
 242 SAGE’s noise injection is lightweight and requires no per-environment gradient storage beyond what  
 243 is needed for computing pairwise cosine similarities. SAGE can also operate as a drop-in replacement  
 244 for the optimizer’s ascent and descent steps, making it compatible with existing gradient manipulation  
 245 methods in MTL (as demonstrated in Table 2) and with the standard DomainBed protocol in DG  
 246 (Table 1).

## 247 5 Experiments

### 248 5.1 Experimental setting

249 **Datasets and Protocol.** Our method is evaluated in two different experimental settings, namely  
 250 Domain Generalization (DG) [20] and Multi-Task learning [21]. For DG, we follow the protocol of  
 251 the widely adopted and challenging DomainBed [41] benchmark, and conduct experiments on five  
 252 image classification datasets. Specifically, we follow the leave-one-domain-out evaluation protocol  
 253 for PACS [42], VLCS [43], OfficeHome [44], TerraIncognita [45], and DomainNet [46], and report  
 254 the average top-1 accuracy over 3 runs based on training-domain split validation. For the MTL setting,  
 255 we follow the standard protocol in recent literature [11, 47, 48] and evaluate SAGE on Cityscapes [49]  
 256 (semantic segmentation and depth estimation), and NYU-v2 [50] (semantic segmentation, depth  
 257 estimation, and surface normal prediction).

258 **Implementation Details.** For the DG experiments, we follow recent literature [51, 10] and finetune  
 259 a CLIP [52] pretrained ViT-B/16 [53] model. For the MTL setting, we follow [54, 17, 47, 48] and  
 260 employ MTAN [55] as the shared backbone, with task-specific attention modules built on top of  
 261 SegNet [56]. Dataset and training hyperparameter details are provided in Appendix E.

262 **5.2 Main results**

263 Table 1 presents the results on the DomainBed benchmark. SAGE achieves an average accuracy of  
 264 78.9%, establishing a new state-of-the-art in this particular setting. Notably, SAGE even surpasses  
 265 recent sharpness-aware adaptations tailored for DG, including SAGM (76.4%), DISAM (76.5%),  
 266 and the BOA regularizer (78.1%), which is applied to SAM-pretrained backbones. Extending  
 267 beyond domain shift, we demonstrate SAGE’s versatility on the NYU-v2 and Cityscapes MTL  
 268 benchmarks (Table 2). The results illustrate that integrating SAGE improves the performance of  
 269 underlying base MTL solvers. For instance, augmenting FairGrad [54] with SAGE (SAGE-FairGrad)  
 270 yields improvements, notably reducing the per-task performance drop against the single-task (STL)  
 271 ( $\Delta m\%$  ↓) from  $-4.96$  to  $-5.63$  on NYU-v2, and from  $3.90$  to  $2.63$  on Cityscapes. Similar relative  
 272 gains are observed when SAGE is applied to Linear Scalarization (LS) and MGDA [57].

Table 1: **Comparison of SAGE on DomainBed.** The top out-of-domain accuracies on five domain generalization benchmarks averaged over three trials, are presented. Results of previously proposed methods are from [51].

Algorithm	PACS	VLCS	OfficeHome	TerraInc	DomainNet	Avg.
ERM [58]	95.9	81.9	84.1	56.1	59.5	75.5
IRM [30]	96.1	82.9	83.2	56.7	59.1	75.6
DANN [59]	96.3	81.8	83.0	56.0	58.4	75.1
CDANN [60]	96.5	82.4	82.9	55.6	58.4	75.2
MMD [61]	95.8	82.3	83.6	57.4	59.9	75.8
RSC [62]	96.5	82.2	83.2	58.2	59.0	75.8
CORAL [63]	96.4	82.6	83.8	57.5	59.8	76.0
IIB [64]	96.5	82.3	84.2	58.2	58.6	76.0
FISH [13]	96.9	82.7	85.0	58.0	59.1	75.0
SAM [7]	96.6	82.9	85.4	56.2	59.8	76.2
GSAM [32]	96.6	82.9	85.6	55.4	59.8	76.1
GAM [26]	96.4	83.6	85.5	55.3	59.5	76.1
SAGM [24]	96.8	82.8	85.2	58.0	59.1	76.4
GGa[15]	97.2	84.1	85.7	58.8	59.8	77.1
DISAM [10]	97.1	82.7	85.4	57.3	59.8	76.5
BOA (SAM pretrained) [51]	<b>97.4</b>	<b>86.5</b>	<b>86.0</b>	<b>60.3</b>	<b>60.2</b>	<b>78.1</b>
<b>SAGE (Ours)</b>	<b>97.4</b>	<b>85.5</b>	<b>87.1</b>	<b>61.1</b>	<b>63.4</b>	<b>78.9</b>

273 **Summary.** Across the two benchmark settings, SAGE establishes a new state-of-the-art on the  
 274 DomainBed DG benchmark (78.9% average) and acts as a general-purpose improvement for base  
 275 MTL solvers. On MTL specifically, specialized methods such as SAMO-FairGrad remain the  
 276 strongest on absolute  $\Delta m\%$ , while SAGE provides consistent gains over the underlying solvers  
 277 (LS, MGDA, FairGrad) without requiring task-routing or solver-specific design. In combination,

Table 2: Results on NYU-v2 (3-task) and Cityscapes (2-task) datasets. The best results are highlighted in **bold** and second-best are underlined. Green plus signs (+) and Red minuses (−) indicate whether the addition of SAGE on top of base solvers improves or hinders training, respectively.

Method	NYU-v2										$\Delta m\%$ ↓	Cityscapes				
	Segmentation		Depth		Surface Normal					Segmentation		Depth		$\Delta m\%$ ↓		
	mIoU ↑	Pix Acc ↑	Abs Err ↓	Rel Err ↓	Angle Distance ↓		Within $\epsilon^\circ$ ↑			mIoU ↑		Pix Acc ↑	Abs Err ↓		Rel Err ↓	
					Mean	Median	<11.25	<22.5	<30							
STL	38.30	63.76	0.6754	0.2780	25.01	19.21	30.14	57.20	69.15		74.01	93.16	0.0125	27.77		
LS	39.29	65.33	0.5493	0.2263	28.15	23.96	22.09	47.50	61.08	5.59	75.18	93.49	0.0155	46.77	22.60	
SI	38.45	64.27	0.5354	0.2201	27.60	23.37	22.53	48.57	62.32	4.39	70.95	91.73	0.0161	33.83	14.11	
RLW [65]	37.17	63.77	0.5759	0.2410	28.27	24.18	22.26	47.05	60.62	7.78	74.57	93.41	0.0158	47.79	24.38	
DWA [55]	39.11	65.31	0.5510	0.2285	27.61	23.18	24.17	50.18	62.39	3.57	75.24	93.52	0.0160	44.37	21.45	
UW [66]	36.87	63.17	0.5446	0.2260	27.04	22.61	23.54	49.05	63.65	4.05	72.02	92.85	0.0140	30.13	5.89	
MGDA [57]	30.47	59.90	0.6070	0.2555	24.88	19.45	29.18	56.88	69.36	1.38	68.84	91.54	0.0309	33.50	44.14	
FCGrad [16]	38.06	64.64	0.5550	0.2325	27.41	22.80	23.86	49.83	63.14	3.97	75.13	93.48	0.0154	42.07	18.29	
GradDrop [67]	39.39	65.12	0.5455	0.2279	27.48	22.96	23.38	49.44	62.87	3.58	75.27	93.53	0.0157	47.54	23.73	
IMTL-G [68]	39.35	65.60	0.5426	0.2256	26.02	21.19	26.20	53.13	66.24	-0.76	75.33	93.49	0.0135	38.41	11.10	
CaGrad [17]	39.79	65.49	0.5486	0.2250	26.31	21.58	25.61	52.36	65.58	0.20	75.16	93.48	0.0141	37.60	11.64	
MoCo [69]	40.30	<b>66.07</b>	0.5575	0.2135	26.67	21.83	25.61	51.78	64.85	0.16	75.42	93.55	0.0149	34.19	9.90	
Nash-MTL [48]	40.13	65.93	<b>0.5261</b>	0.2171	25.26	20.08	28.40	55.47	68.15	-4.04	75.41	93.66	0.0129	35.02	6.82	
FAMO [47]	38.88	64.90	0.5474	0.2194	25.06	19.57	29.21	56.61	68.98	-4.10	74.54	93.29	0.0145	32.59	8.13	
FairGrad [54]	38.80	65.29	0.5572	0.2322	24.55	18.97	30.50	57.94	70.14	<b>-4.96</b>	74.10	93.03	0.0135	29.92	3.90	
F-MTL [70]	<b>40.42</b>	65.61	0.5389	0.2121	25.03	19.75	28.90	56.19	68.72	-4.77	76.63	93.76	0.0124	31.17	1.87	
SAMO-LS [11]	39.59	<u>65.72</u>	0.5514	0.2246	27.38	22.78	24.09	49.82	63.01	2.88	<b>76.46</b>	<b>93.76</b>	0.0147	39.85	14.30	
SAMO-MGDA [11]	29.85	60.83	0.6111	0.2388	<b>24.11</b>	<b>18.18</b>	<b>32.16</b>	<b>59.59</b>	<b>71.15</b>	-2.19	73.28	93.26	0.0133	<b>30.57</b>	4.30	
SAMO-FairGrad [11]	39.05	65.06	0.5359	0.2137	24.43	18.79	30.98	58.35	70.42	<b>-6.55</b>	74.37	93.14	0.0129	<b>26.30</b>	<b>-0.62</b>	
SAGE-LS	39.67 +	65.03 −	0.5331 +	0.2180 +	27.40 +	23.06 +	23.62 +	49.28 +	62.87 +	3.81 +	75.69 +	93.57 +	0.0146 +	42.84 +	19.92 +	
SAGE-MGDA	33.24 +	61.88 +	0.5707 +	0.2266 +	24.55 +	18.99 +	30.15 +	57.91 +	70.07 +	-1.77 +	73.46 +	93.28 +	<b>0.0123 +</b>	29.59 +	3.82 +	
SAGE-FairGrad	<u>40.32 +</u>	65.34 +	<u>0.5301 +</u>	<b>0.2082 +</b>	24.53 +	18.97 ±	30.28 −	57.44 −	69.66 −	<b>-5.63 +</b>	75.40 +	93.63 +	0.0129 +	29.33 +	2.63 +	
SAGE-SAMO-LS	40.15 +	65.59 −	0.5508 +	0.2251 −	27.40 −	22.76 +	23.34 −	49.86 +	63.56 +	3.82 −	<u>76.07 −</u>	<u>93.73 −</u>	0.0149 −	42.15 −	18.71 −	
SAGE-SAMO-MGDA	32.32 +	62.78 +	0.5651 +	0.2201 +	<u>24.23 −</u>	<b>18.10 +</b>	<b>32.41 +</b>	<b>59.68 +</b>	<b>71.16 +</b>	-3.66 +	73.92 +	93.20 −	0.0133 ±	32.46 −	9.85 −	
SAGE-SAMO-FairGrad	39.76 +	65.80 +	0.5413 −	0.2164 −	24.34 +	18.69 +	31.01 +	58.38 +	70.46 +	<b>-6.19 −</b>	75.08 +	93.52 +	0.0131 −	32.02 −	4.88 −	

278 the DG experiments demonstrate that targeting both curvature and gradient alignment can yield  
 279 state-of-the-art results, while the MTL results demonstrate that the same mechanisms can be added to  
 280 existing gradient-balancing methods and yield improved results.

### 281 5.3 Ablations

282 Figure 3 (left) reports OfficeHome accuracy over selections of the perturbation radius  $\rho$  and the  
 283 noise scale  $\gamma$ . Performance is stable across selection, with the exception of the selection of either  
 284 relatively low or high values. In the same, figure on the right we compare the performance of SAGE  
 285 when removing each of its components and against the Muon [28] optimizer. As apparent from the  
 286 results, each component of SAGE, i.e. orthogonalization via Newton-Schulz, weight scaling in each  
 287 layer and the additive noise, contributes positively during training. Further ablations, such as cosine  
 288 similarity plots between the average training-domain and unseen target domain gradients are provided  
 289 in Appendix F.

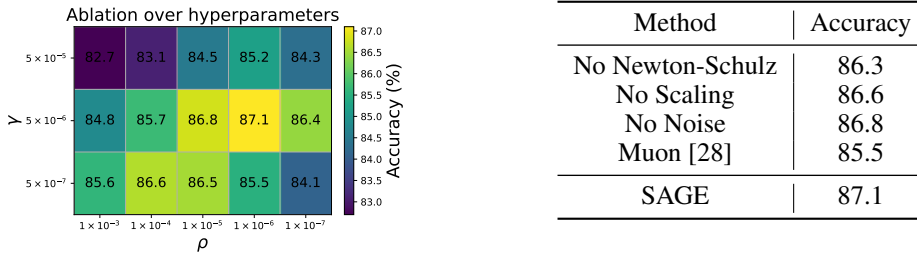


Figure 3: *Left*: Ablation results for the  $\rho$  and  $\gamma$  hyperparameters and *right*: Ablations regarding each component of SAGE, on the OfficeHome dataset.

## 290 6 Conclusions, Limitations & Future Work

291 In this paper, we presented an excess-risk decomposition for multi-distribution learning that exposes  
 292 a structural trade-off between the curvature of the loss surface and the alignment of per-distribution  
 293 gradients, along with a separation result showing that neither quantity bounds the other. Motivated  
 294 by this decomposition, we proposed SAGE, whose two components target the two terms via: a  
 295 spectral perturbation that probes all singular directions of each layer’s gradient with comparable  
 296 magnitude, and a gradient-disagreement-scaled noise injection that biases the optimizer away from  
 297 high-conflict regions. Across five DG and two MTL benchmarks, SAGE compares favorably with  
 298 both sharpness-aware and gradient-alignment baselines.

299 One limitation of SAGE is that it incurs additional per-step cost relative to vanilla SAM, as New-  
 300 ton-Schulz orthogonalization adds matrix multiplications on every weight matrix, and the noise  
 301 mechanism requires per-distribution gradient computation (Appendix G). Additionally, the empirical  
 302 gains are not uniformly state-of-the-art across all benchmarks. In the MTL setting in particular, when  
 303 paired with base solvers SAGE improves performance in almost all cases, but does not dominate every  
 304 metric compared to methods that have been specifically developed for MTL. More fundamentally,  
 305 the connection between SAGE’s mechanisms and the terms of Theorem 1 is motivational rather than  
 306 formal: we have not shown that either component provably reduces a tractable surrogate of its target  
 307 term, and Counterexample 1 is an argument in the quadratic family rather than a guarantee for the  
 308 non-convex deep-learning regime. These limitations point to natural directions for future work. For  
 309 example, methods that directly estimate or bound the two terms of Theorem 1 could replace the  
 310 spectral perturbation with a target that more closely matches the curvature term, while estimators  
 311 of  $\Sigma_g$  could be coupled with regularizers that penalize the alignment term more directly. Finally,  
 312 whether the curvature-alignment trade-off identified also applies to other settings such as continual  
 313 learning or federated optimization, where multi-distribution structure arises from different generative  
 314 processes than DG and MTL, is an interesting open question which we aim to research in future  
 315 work.

## References

- 316
- 317 [1] Hao Li, Zheng Xu, Gavin Taylor, Christoph Studer, and Tom Goldstein. Visualizing the loss  
318 landscape of neural nets. *Advances in neural information processing systems*, 31, 2018.
- 319 [2] Stanislav Fort and Stanislaw Jastrzebski. Large scale structure of neural network loss landscapes.  
320 *Advances in Neural Information Processing Systems*, 32, 2019.
- 321 [3] Lei Wu, Zhanxing Zhu, et al. Towards understanding generalization of deep learning: Perspec-  
322 tive of loss landscapes. *arXiv preprint arXiv:1706.10239*, 2017.
- 323 [4] Akshay Rangamani et al. *Loss landscapes and generalization in neural networks: Theory and*  
324 *applications*. PhD thesis, Johns Hopkins University, 2020.
- 325 [5] Nitish Shirish Keskar, Dheevatsa Mudigere, Jorge Nocedal, Mikhail Smelyanskiy, and Ping  
326 Tak Peter Tang. On large-batch training for deep learning: Generalization gap and sharp minima.  
327 *arXiv preprint arXiv:1609.04836*, 2016.
- 328 [6] Yiding Jiang, Behnam Neyshabur, Hossein Mobahi, Dilip Krishnan, and Samy Bengio. Fantastic  
329 generalization measures and where to find them. *arXiv preprint arXiv:1912.02178*, 2019.
- 330 [7] Pierre Foret, Ariel Kleiner, Hossein Mobahi, and Behnam Neyshabur. Sharpness-aware min-  
331 imization for efficiently improving generalization. In *International Conference on Learning*  
332 *Representations*, 2021.
- 333 [8] Jungmin Kwon, Jeongseop Kim, Hyunseo Park, and In Kwon Choi. Asam: Adaptive sharpness-  
334 aware minimization for scale-invariant learning of deep neural networks. In *International*  
335 *conference on machine learning*, pages 5905–5914. PMLR, 2021.
- 336 [9] Minyoung Kim, Da Li, Shell X Hu, and Timothy Hospedales. Fisher sam: Information geometry  
337 and sharpness aware minimisation. In *International Conference on Machine Learning*, pages  
338 11148–11161. PMLR, 2022.
- 339 [10] Ruipeng Zhang, Ziqing Fan, Jiangchao Yao, Ya Zhang, and Yanfeng Wang. Domain-inspired  
340 sharpness aware minimization under domain shifts. In *The Twelfth International Conference on*  
341 *Learning Representations*, 2024. URL <https://openreview.net/forum?id=I4wB3HA3dJ>.
- 342 [11] Hao Ban, Gokul Ram Subramani, and Kaiyi Ji. Samo: A lightweight sharpness-aware approach  
343 for multi-task optimization with joint global-local perturbation. In *Proceedings of the IEEE/CVF*  
344 *International Conference on Computer Vision*, pages 785–795, 2025.
- 345 [12] Lucas Mansilla, Rodrigo Echeveste, Diego H Milone, and Enzo Ferrante. Domain generalization  
346 via gradient surgery. In *Proceedings of the IEEE/CVF international conference on computer*  
347 *vision*, pages 6630–6638, 2021.
- 348 [13] Yuge Shi, Jeffrey Seely, Philip Torr, Siddharth N, Awni Hannun, Nicolas Usunier, and Gabriel  
349 Synnaeve. Gradient matching for domain generalization. In *International Conference on*  
350 *Learning Representations*, 2022. URL <https://openreview.net/forum?id=vDwBW49Hm0>.
- 351 [14] Binh M Le and Simon S Woo. Gradient alignment for cross-domain face anti-spoofing. In  
352 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages  
353 188–199, 2024.
- 354 [15] Aristotelis Ballas and Christos Diou. Gradient-guided annealing for domain generalization. In  
355 *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 20558–20568,  
356 2025.
- 357 [16] Tianhe Yu, Saurabh Kumar, Abhishek Gupta, Sergey Levine, Karol Hausman, and Chelsea Finn.  
358 Gradient surgery for multi-task learning. *Advances in Neural Information Processing Systems*,  
359 33:5824–5836, 2020.
- 360 [17] Bo Liu, Xingchao Liu, Xiaojie Jin, Peter Stone, and Qiang Liu. Conflict-averse gradient descent  
361 for multi-task learning. *Advances in neural information processing systems*, 34:18878–18890,  
362 2021.

- 363 [18] Dmitry Senushkin, Nikolay Patakin, Arseny Kuznetsov, and Anton Konushin. Independent  
364 component alignment for multi-task learning. In *Proceedings of the IEEE/CVF Conference on*  
365 *Computer Vision and Pattern Recognition (CVPR)*, pages 20083–20093, June 2023.
- 366 [19] Jindong Wang, Cuiling Lan, Chang Liu, Yidong Ouyang, Tao Qin, Wang Lu, Yiqiang Chen,  
367 Wenjun Zeng, and S Yu Philip. Generalizing to unseen domains: A survey on domain general-  
368 ization. *IEEE transactions on knowledge and data engineering*, 35(8):8052–8072, 2022.
- 369 [20] Kaiyang Zhou, Ziwei Liu, Yu Qiao, Tao Xiang, and Chen Change Loy. Domain Generalization:  
370 A Survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(4):4396–4415,  
371 April 2023. ISSN 1939-3539. doi: 10.1109/TPAMI.2022.3195549. Conference Name: IEEE  
372 Transactions on Pattern Analysis and Machine Intelligence.
- 373 [21] Yu Zhang and Qiang Yang. A survey on multi-task learning. *IEEE Transactions on Knowledge*  
374 *and Data Engineering*, 34(12):5586–5609, 2022. doi: 10.1109/TKDE.2021.3070203.
- 375 [22] Bingcong Li and Georgios Giannakis. Enhancing sharpness-aware optimization through variance  
376 suppression. *Advances in Neural Information Processing Systems*, 36:70861–70879, 2023.
- 377 [23] Yilang Zhang, Bingcong Li, and Georgios B Giannakis. Preconditioned sharpness-aware  
378 minimization: Unifying analysis and a novel learning algorithm. In *ICASSP 2025-2025 IEEE*  
379 *International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1–5.  
380 IEEE, 2025.
- 381 [24] Pengfei Wang, Zhaoxiang Zhang, Zhen Lei, and Lei Zhang. Sharpness-aware gradient matching  
382 for domain generalization. In *Proceedings of the IEEE/CVF Conference on Computer Vision*  
383 *and Pattern Recognition*, pages 3769–3778, 2023.
- 384 [25] Jeremy Bernstein and Laker Newhouse. Old optimizer, new norm: An anthology. *arXiv preprint*  
385 *arXiv:2409.20325*, 2024.
- 386 [26] Xingxuan Zhang, Renzhe Xu, Han Yu, Hao Zou, and Peng Cui. Gradient norm aware minimiza-  
387 tion seeks first-order flatness and improves generalization. In *Proceedings of the IEEE/CVF*  
388 *Conference on Computer Vision and Pattern Recognition*, pages 20247–20257, 2023.
- 389 [27] Dahun Shin, Dongyeop Lee, Jinseok Chung, and Namhoon Lee. Sassa: Sharpness-aware  
390 adaptive second-order optimization with stable hessian approximation. In *Forty-second Inter-*  
391 *national Conference on Machine Learning, 2025*. URL [https://openreview.net/forum?](https://openreview.net/forum?id=7bgqx50oVe)  
392 [id=7bgqx50oVe](https://openreview.net/forum?id=7bgqx50oVe).
- 393 [28] Keller Jordan, Yuchen Jin, Vlado Boza, You Jiacheng, Franz Cesista, Laker Newhouse, and  
394 Jeremy Bernstein. Muon: An optimizer for hidden layers in neural networks, 2024. URL  
395 <https://kellerjordan.github.io/posts/muon/>.
- 396 [29] Vineet Gupta, Tomer Koren, and Yoram Singer. Shampoo: Preconditioned stochastic tensor  
397 optimization. In *International Conference on Machine Learning*, pages 1842–1850. PMLR,  
398 2018.
- 399 [30] Martin Arjovsky, Léon Bottou, Ishaan Gulrajani, and David Lopez-Paz. Invariant risk mini-  
400 mization. *arXiv preprint arXiv:1907.02893*, 2019.
- 401 [31] Aad W Van der Vaart. *Asymptotic statistics*, volume 3. Cambridge university press, 2000.
- 402 [32] Juntang Zhuang, Boqing Gong, Liangzhe Yuan, Yin Cui, Hartwig Adam, Nicha C Dvornek,  
403 sekhar tatikonda, James s Duncan, and Ting Liu. Surrogate gap minimization improves  
404 sharpness-aware training. In *International Conference on Learning Representations, 2022*. URL  
405 <https://openreview.net/forum?id=ed0NMANhLu->.
- 406 [33] Kaiyue Wen, Tengyu Ma, and Zhiyuan Li. How does sharpness-aware minimization minimize  
407 sharpness?, 2023. URL <https://arxiv.org/abs/2211.05729>.
- 408 [34] Laurent Dinh, Razvan Pascanu, Samy Bengio, and Yoshua Bengio. Sharp minima can generalize  
409 for deep nets. In *International Conference on Machine Learning*, pages 1019–1028. PMLR,  
410 2017.

- 411 [35] Chun-Hua Guo and Nicholas J Higham. A schur–newton method for the matrix  $p$  th  
412 root and its inverse. *SIAM Journal on Matrix Analysis and Applications*, 28(3):788–804, 2006.
- 413 [36] Åke Björck and Clazett Bowie. An iterative algorithm for computing the best estimate of an  
414 orthogonal matrix. *SIAM Journal on Numerical Analysis*, 8(2):358–364, 1971.
- 415 [37] Zdislav Kovarik. Some iterative methods for improving orthonormality. *SIAM Journal on*  
416 *Numerical Analysis*, 7(3):386–389, 1970.
- 417 [38] Aristotelis Ballas and Christos Diou. Towards domain generalization for ecg and eeg classifica-  
418 tion: Algorithms and benchmarks. *IEEE Transactions on Emerging Topics in Computational*  
419 *Intelligence*, 8(1):44–54, 2024. doi: 10.1109/TETCI.2023.3306253.
- 420 [39] Max Welling and Yee W Teh. Bayesian learning via stochastic gradient langevin dynamics.  
421 In *Proceedings of the 28th international conference on machine learning (ICML-11)*, pages  
422 681–688, 2011.
- 423 [40] Yee Whye Teh, Alexandre Thiéry, and Sebastian J Vollmer. Consistency and fluctuations for  
424 stochastic gradient langevin dynamics. *Journal of Machine Learning Research*, 17(7), 2016.
- 425 [41] Ishaan Gulrajani and David Lopez-Paz. In search of lost domain generalization. In *International*  
426 *Conference on Learning Representations*, 2021.
- 427 [42] Da Li, Yongxin Yang, Yi-Zhe Song, and Timothy M Hospedales. Deeper, broader and artier  
428 domain generalization. In *Proceedings of the IEEE international conference on computer vision*,  
429 pages 5542–5550, 2017.
- 430 [43] Chen Fang, Ye Xu, and Daniel N Rockmore. Unbiased metric learning: On the utilization of  
431 multiple datasets and web images for softening bias. In *Proceedings of the IEEE International*  
432 *Conference on Computer Vision*, pages 1657–1664, 2013.
- 433 [44] Hemanth Venkateswara, Jose Eusebio, Shayok Chakraborty, and Sethuraman Panchanathan.  
434 Deep hashing network for unsupervised domain adaptation. In *Proceedings of the IEEE*  
435 *conference on computer vision and pattern recognition*, pages 5018–5027, 2017.
- 436 [45] Sara Beery, Grant Van Horn, and Pietro Perona. Recognition in terra incognita. In *Proceedings*  
437 *of the European conference on computer vision (ECCV)*, pages 456–473, 2018.
- 438 [46] Xingchao Peng, Qinxun Bai, Xide Xia, Zijun Huang, Kate Saenko, and Bo Wang. Moment  
439 matching for multi-source domain adaptation. In *Proceedings of the IEEE/CVF international*  
440 *conference on computer vision*, pages 1406–1415, 2019.
- 441 [47] Bo Liu, Yihao Feng, Peter Stone, and Qiang Liu. Famo: Fast adaptive multitask optimization.  
442 *Advances in Neural Information Processing Systems*, 36:57226–57243, 2023.
- 443 [48] Aviv Navon, Aviv Shamsian, Idan Achituve, Haggai Maron, Kenji Kawaguchi, Gal Chechik,  
444 and Ethan Fetaya. Multi-task learning as a bargaining game. In *International Conference on*  
445 *Machine Learning*, pages 16428—16446. PMLR, 2022.
- 446 [49] Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo  
447 Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The cityscapes dataset for semantic  
448 urban scene understanding. In *Proceedings of the IEEE conference on computer vision and*  
449 *pattern recognition*, pages 3213–3223, 2016.
- 450 [50] Nathan Silberman, Derek Hoiem, Pushmeet Kohli, and Rob Fergus. Indoor segmentation  
451 and support inference from rgb-d images. In *European conference on computer vision*, pages  
452 746–760. Springer, 2012.
- 453 [51] Aodi Li, Liansheng Zhuang, Xiao Long, Houqiang Li, and Shafei Wang. Exploring mode  
454 connectivity in krylov subspace for domain generalization. In *The Fourteenth International*  
455 *Conference on Learning Representations*, 2026.

- 456 [52] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,  
457 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual  
458 models from natural language supervision. In *International conference on machine learning*,  
459 pages 8748–8763. PmLR, 2021.
- 460 [53] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai,  
461 Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly,  
462 Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image  
463 recognition at scale. In *International Conference on Learning Representations*, 2021. URL  
464 <https://openreview.net/forum?id=YicbFdNTTy>.
- 465 [54] Hao Ban and Kaiyi Ji. Fair resource allocation in multi-task learning. In *Proceedings of the*  
466 *41st International Conference on Machine Learning, ICML'24*. JMLR.org, 2024.
- 467 [55] Shikun Liu, Edward Johns, and Andrew J Davison. End-to-end multi-task learning with attention.  
468 In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages  
469 1871–1880, 2019.
- 470 [56] Vijay Badrinarayanan, Alex Kendall, and Roberto Cipolla. Segnet: A deep convolutional  
471 encoder-decoder architecture for image segmentation. *IEEE transactions on pattern analysis*  
472 *and machine intelligence*, 39(12):2481–2495, 2017.
- 473 [57] Ozan Sener and Vladlen Koltun. Multi-task learning as multi-objective optimization. *Advances*  
474 *in neural information processing systems*, 31, 2018.
- 475 [58] V Vapnik. Statistical learning theory. *NY: Wiley*, 1998.
- 476 [59] Yaroslav Ganin, Evgeniya Ustinova, Hana Ajakan, Pascal Germain, Hugo Larochelle, François  
477 Laviolette, Mario March, and Victor Lempitsky. Domain-adversarial training of neural networks.  
478 *Journal of machine learning research*, 17(59):1–35, 2016.
- 479 [60] Ya Li, Mingming Gong, Xinmei Tian, Tongliang Liu, and Dacheng Tao. Domain generaliza-  
480 tion via conditional invariant representations. In *AAAI Conference on Artificial Intelligence*,  
481 volume 32, 2018.
- 482 [61] Haoliang Li, Sinno Jialin Pan, Shiqi Wang, and Alex C Kot. Domain generalization with  
483 adversarial feature learning. In *Computer Vision and Pattern Recognition*, 2018.
- 484 [62] Zeyi Huang, Haohan Wang, Eric P Xing, and Dong Huang. Self-challenging improves cross-  
485 domain generalization. *European Conference on Computer Vision*, 2020.
- 486 [63] Baochen Sun and Kate Saenko. Deep coral: Correlation alignment for deep domain adaptation.  
487 In *European Conference on Computer Vision*, 2016.
- 488 [64] Bo Li, Yifei Shen, Yezhen Wang, Wenzhen Zhu, Dongsheng Li, Kurt Keutzer, and Han  
489 Zhao. Invariant information bottleneck for domain generalization. In *Proceedings of the AAAI*  
490 *Conference on Artificial Intelligence*, volume 36, pages 7399–7407, 2022.
- 491 [65] Baijiong Lin, Feiyang Ye, Yu Zhang, and Ivor W Tsang. Reasonable effectiveness of random  
492 weighting: A litmus test for multi-task learning. *arXiv preprint arXiv:2111.10603*, 2021.
- 493 [66] Alex Kendall, Yarin Gal, and Roberto Cipolla. Multi-task learning using uncertainty to weigh  
494 losses for scene geometry and semantics. In *Proceedings of the IEEE conference on computer*  
495 *vision and pattern recognition*, pages 7482–7491, 2018.
- 496 [67] Zhao Chen, Jiquan Ngiam, Yanping Huang, Thang Luong, Henrik Kretschmar, Yuning Chai,  
497 and Dragomir Anguelov. Just pick a sign: Optimizing deep multitask models with gradient sign  
498 dropout. *Advances in Neural Information Processing Systems*, 33:2039–2050, 2020.
- 499 [68] Liyang Liu, Yi Li, Zhanghui Kuang, Jing-Hao Xue, Yimin Chen, Wenming Yang, Qingmin Liao,  
500 and Wayne Zhang. Towards impartial multi-task learning. In *International Conference on Learn-*  
501 *ing Representations*, 2021. URL <https://openreview.net/forum?id=IMPnRXEWpvr>.

- 502 [69] Heshan Devaka Fernando, Han Shen, Miao Liu, Subhajit Chaudhury, Keerthiram Murugesan,  
503 and Tianyi Chen. Mitigating gradient bias in multi-objective learning: A provably convergent  
504 approach. In *The eleventh international conference on learning representations*, 2023.
- 505 [70] Hoang Phan, Lam Tran, Quyen Tran, Ngoc Tran, Tuan Truong, Qi Lei, Nhat Ho, Dinh Phung,  
506 and Trung Le. Beyond losses reweighting: Empowering multi-task learning via the general-  
507 ization perspective. In *Proceedings of the IEEE/CVF International Conference on Computer*  
508 *Vision*, pages 2440–2450, 2025.
- 509 [71] David Krueger, Ethan Caballero, Joern-Henrik Jacobsen, Amy Zhang, Jonathan Binas, Dinghui  
510 Zhang, Remi Le Priol, and Aaron Courville. Out-of-distribution generalization via risk ex-  
511 trapolation (rex). In *International conference on machine learning*, pages 5815–5826. PMLR,  
512 2021.
- 513 [72] Alexandre Rame, Corentin Dancette, and Matthieu Cord. Fishr: Invariant gradient variances for  
514 out-of-distribution generalization. In *International Conference on Machine Learning*, pages  
515 18347–18377. PMLR, 2022.
- 516 [73] Junbum Cha, Sanghyuk Chun, Kyungjae Lee, Han-Cheol Cho, Seunghyun Park, Yunsung Lee,  
517 and Sungrae Park. Swad: Domain generalization by seeking flat minima. *Advances in Neural*  
518 *Information Processing Systems*, 34:22405–22418, 2021.
- 519 [74] Jean-Antoine Désidéri. Multiple-gradient descent algorithm (mgda) for multiobjective opti-  
520 mization. *Comptes Rendus. Mathématique*, 350(5-6):313–318, 2012.
- 521 [75] Hoang Phan, Tung Lam Tran, Ngoc N Tran, Nhat Ho, Dinh Phung, and Trung Le. Improving  
522 multi-task learning via seeking task-based flat regions. 2022.

523

524

# Appendix

525

526

## Table of Contents

---

527	<b>A Extended Related Work</b>	<b>16</b>
528	A.1 Domain Generalization . . . . .	16
529	A.2 Multi-Task Learning . . . . .	17
530	A.3 Positioning of SAGE . . . . .	18
531	<b>B Proofs</b>	<b>18</b>
532	B.1 Setup Reminder: . . . . .	18
533	B.2 Proof of Theorem 1: Multi-distribution excess-risk decomposition . . . . .	19
534	B.3 Proof of Counterexample 1: Decoupling of flatness and gradient alignment . . .	22
535	B.4 Limitations . . . . .	23
536	<b>C Motivating Example: Why Both Properties Are Necessary</b>	<b>23</b>
537	<b>D Scale Invariance and Scale-Dependent Flatness in SAM</b>	<b>24</b>
538	D.1 Illustrative Example . . . . .	25
539	<b>E Experimental details &amp; Hyperparameters</b>	<b>25</b>
540	E.1 Details & Hyperparameters . . . . .	25
541	E.2 Infrastructure . . . . .	26
542	<b>F Additional Ablations</b>	<b>26</b>
543	<b>G Computational Analysis</b>	<b>26</b>
544	<b>H Pseudo code of the proposed SAGE algorithm</b>	<b>26</b>

---

545

546

547

## 548 A Extended Related Work

### 549 A.1 Domain Generalization

550 Domain generalization (DG) aims to train a model on data from multiple source domains such that it  
551 generalizes to previously unseen target domains whose distributions may differ substantially from  
552 those encountered during training [19, 20]. The core challenge is that Empirical Risk Minimization  
553 (ERM), which minimizes the average loss over all training data, may converge to solutions that  
554 exploit domain-specific shortcuts rather than domain-invariant features, resulting in poor out-of-  
555 distribution performance. In this section, we will focus on the methodological groups that are closest  
556 in nature with our proposed method, namely domain-invariant representation learning, sharpness-  
557 aware optimization, gradient-based alignment, and loss-landscape exploration.

558 **Domain-invariant and robust optimization methods.** A central line of work found in the literature  
559 seeks representations that are invariant across training domains. In one of the most fundamental  
560 works, Invariant Risk Minimization (IRM) [30] constrains the optimal classifier to be identical across  
561 all data domains, while V-REx [71] penalizes the variance of per-domain risks. Fishr [72] on the other  
562 hand enforces consistency of per-domain gradient variances, and CORAL [63] aligns second-order  
563 feature statistics across domains. Although these methods directly modify the training objective  
564 to encourage domain-invariant solutions, they do not directly consider the geometry of the loss  
565 landscape.

566 **Sharpness-aware minimization for DG.** As already mentioned in the main text, the observation  
567 that flat minima tend to generalize better than sharp ones [7, 5] has motivated the application of SAM  
568 and its variants to domain generalization. Specifically, GSAM [32] refines the sharpness estimate  
569 via a surrogate gap objective, ASAM [8] introduces adaptive, scale-aware perturbation magnitudes,  
570 GAM [26] regularizes the gradient norm to seek first-order flatness, and VaSSO [22] addresses the  
571 noise introduced by mini-batch sampling in the perturbation direction. Furthermore, SAGM [24]  
572 simultaneously minimizes the empirical loss, the perturbed loss, and the surrogate gap between them,  
573 encouraging the clean-loss gradient  $\nabla\mathcal{L}(\theta)$  and the perturbed-loss gradient  $\nabla\mathcal{L}_p(\theta)$  to remain aligned  
574 so that the optimizer converges to regions that are both flat and low-loss. Importantly, the gradient  
575 matching is between the clean and perturbed objectives at the *aggregate* level and does not involve  
576 per-domain gradient information. Finally, Fisher SAM [9] re-parameterizes the perturbation using  
577 the Fisher information metric to better capture the intrinsic geometry of the model.

578 A key limitation of vanilla SAM under domain shifts was identified by Zhang et al. [10], who observed  
579 that when domains differ in difficulty or data quantity, SAM’s gradient-based perturbation is biased  
580 toward whichever domain has the largest gradient, disrupting training for the others. To address  
581 this, they proposed DISAM (Domain-Inspired SAM), which incorporates a domain-loss variance  
582 minimization constraint into the perturbation generation step. The resulting perturbation adaptively  
583 up-weights well-converged domains and down-weights under-converged ones, yielding faster overall  
584 convergence and improved generalization. DISAM is compatible with other SAM variants (GSAM,  
585 SAGM) and achieves consistent improvements on the DomainBed benchmark, particularly on datasets  
586 with large domain gaps such as TerraIncognita. Finally, complementary ensemble-based approaches  
587 such as SWAD [73] achieve flatness passively by averaging model weights densely along the training  
588 trajectory, rather than explicitly optimizing a sharpness objective.

589 **Gradient-alignment methods for DG.** A parallel family of approaches leverages the agreement of  
590 per-domain gradients as a signal for domain invariance. The intuition is straightforward: if a model’s  
591 internal representation does not depend on the domain, then the expected gradients of the loss with  
592 respect to different domains should point in similar directions. Shi et al. [13] proposed Fish, which  
593 approximately maximizes the inner product of per-domain gradients. Mansilla et al. [12] adapted  
594 gradient surgery from multi-task learning [16] to DG by muting or randomizing gradient dimensions  
595 where domains disagree in sign. Le and Woo [14] applied gradient alignment to cross-domain face  
596 anti-spoofing, demonstrating the generality of the approach beyond standard image classification  
597 benchmarks. More recently, Gradient-Guided Annealing (GGA) [15] was proposed, for searching  
598 parameter configurations via a Simulated Annealing-inspired process, where domain gradients are  
599 well-aligned before continuing standard optimization.

600 **Beyond local flatness: mode connectivity.** While the methods above focus on properties local to  
601 a single minimum (flatness, gradient alignment), Li et al. [51] recently shifted attention to a global  
602 property of the loss landscape: *mode connectivity*, the phenomenon whereby distinct local minima are  
603 connected by continuous low-loss pathways. They demonstrated empirically that models with poor  
604 and strong out-of-domain generalization can be connected via such pathways and proposed the Billiard  
605 Optimization Algorithm (BOA), which navigates these pathways by alternating line search (to locate  
606 loss contour boundaries) and reflection (to redirect the trajectory along the contour). To overcome  
607 the curse of dimensionality in high-dimensional parameter spaces, BOA operates within a low-  
608 dimensional Krylov subspace constructed from training gradients and Hessian–vector products, which  
609 they showed to be well-aligned with test gradients across diverse datasets and architectures. BOA  
610 achieves state-of-the-art results on DomainBed with ViT backbones, outperforming all sharpness-  
611 aware baselines, and represents a qualitatively different approach to loss-landscape exploitation for  
612 DG.

## 613 A.2 Multi-Task Learning

614 Multi-task learning (MTL) trains a single model to perform multiple tasks simultaneously, with the  
615 goal of exploiting shared structure across tasks to improve data efficiency and generalization [21].  
616 The standard approach minimizes the average loss across all tasks, however, this simple objective  
617 often leads to suboptimal solutions because the gradients of different tasks may differ substantially  
618 in magnitude or direction, a phenomenon commonly referred to as *task conflict* or *conflicting*  
619 *gradients* [16, 17]. When one task’s gradient dominates the shared update, other tasks may stagnate  
620 or even regress, resulting in performance that falls short of what independent single-task models can  
621 achieve. Addressing this fundamental optimization challenge has motivated a rich body of work,  
622 which can be broadly organized into loss-weighting methods, gradient manipulation methods, and  
623 more recent approaches that integrate sharpness-aware minimization into the MTL framework.

624 **Loss weighting methods.** A natural strategy for balancing tasks is to assign scalar weights to  
625 each task’s loss and adapt those weights during training. Kendall et al. [66] proposed Uncertainty  
626 Weighting (UW), which uses learned task-dependent homoscedastic uncertainty as a proxy for task  
627 difficulty. Liu et al. [55] introduced Dynamic Weight Average (DWA), a heuristic that adjusts task  
628 weights based on the rate of change of each task’s loss over recent iterations. Lin et al. [65] showed  
629 that randomly sampling task weights at each step, i.e., Random Loss Weighting (RLW), can be  
630 surprisingly effective. A scale-invariant (SI) baseline that minimizes the sum of log-losses has also  
631 been considered, as it is invariant to multiplicative rescaling of individual task losses [47]. While  
632 computationally lightweight, these methods adjust only the loss aggregation and do not directly  
633 resolve geometric conflicts in the gradient space.

634 **Gradient manipulation methods.** A more principled family of approaches directly modifies the  
635 update direction to ensure balanced progress across tasks. Désidéri [74] introduced the Multiple  
636 Gradient Descent Algorithm (MGDA), which finds a convex combination of task gradients lying  
637 in the common descent cone, and Sener and Koltun [57] adapted MGDA for deep networks. PC-  
638 Grad [16] projects each task’s gradient onto the normal plane of conflicting gradients from other tasks.  
639 GradDrop [67] randomly zeroes gradient dimensions where tasks disagree in sign. CAGrad [17]  
640 constrains the update within a ball around the average gradient while maximizing worst-case per-task  
641 improvement. IMTL-G [68] finds an update direction with equal cosine similarity to all task gradients,  
642 while ICA [18] decomposes task gradients into independent components. Nash-MTL [48] formulates  
643 the problem as a bargaining game, and FairGrad [54] extends the bargaining framework with a  
644 fairness-aware conic combination. MoCo [69] corrects gradient bias through a momentum-based  
645 mechanism. A common limitation of all these methods is their  $\mathcal{O}(K)$  time and space overhead per  
646 iteration, which becomes prohibitive when both the number of tasks and the model size are large.

647 **Efficient multi-task optimization.** To address this scalability bottleneck, Liu et al. [47] proposed  
648 FAMO, which decreases all task losses at approximately equal *rates* by working in log-loss space  
649 and amortizing the weight update using only the change in task losses between consecutive iterations.  
650 This achieves  $\mathcal{O}(1)$  space and time overhead per iteration, matching the complexity of standard  
651 average-loss training while remaining competitive with gradient manipulation methods on standard  
652 benchmarks.

653 **Sharpness-aware minimization in MTL.** Ban et al. [11] provided empirical evidence that SAM [7]  
654 mitigates task conflicts by guiding optimization toward flatter loss regions where changes in one  
655 task’s loss do not substantially affect others. Their analysis revealed that applying SAM to the  
656 average loss (global SAM) and to each task individually (local SAM) both contribute to conflict  
657 reduction, but neither is consistently superior. Phan et al. [70] proposed F-MTL, which applies SAM  
658 independently to each task but incurs  $K$  additional backpropagations. Ban et al. [11] introduced  
659 SAMO, which combines global and local perturbation information via a weighted average and  
660 approximates expensive local gradients using a zeroth-order SPSA estimator requiring only forward  
661 passes. A layerwise normalization scheme stabilizes the estimates, achieving the benefits of joint  
662 global–local perturbations at a cost comparable to global-only SAM.

### 663 A.3 Positioning of SAGE

664 The methods reviewed above target individual geometric properties of the loss landscape, i.e., flatness,  
665 gradient alignment, mode connectivity, or task balancing, but limited methods address more than one  
666 simultaneously. Sharpness-aware methods (SAM, GSAM, SAGM, DISAM) seek flat minima without  
667 explicitly considering whether per-distribution gradients agree at those minima. Gradient-alignment  
668 methods (Fish, GGA, PCGrad in DG; PCGrad, CAGrad, Nash-MTL in MTL) enforce gradient  
669 agreement without considering the curvature of the resulting solution. The closest prior work is that  
670 of SAGM [24], which also augments SAM with a gradient-matching term. However, the two methods  
671 differ in three respects. First, SAGM’s gradient matching aligns the clean-loss and perturbed-loss  
672 gradients, which are both computed on the *aggregate* data, to ensure joint descent toward flat, low-loss  
673 regions and does not compute per-distribution gradients.

674 SAGE addresses the curvature through a spectral perturbation that replaces SAM’s gradient-scaled  
675 ascent with the polar factor of each layer’s gradient matrix, computed via Newton–Schulz iteration  
676 and scaled by the layer’s Frobenius norm. Regarding gradient alignment term, SAGE injects  
677 isotropic Gaussian noise at the descent step whose magnitude scales with the degree of cross-  
678 environment gradient conflict,  $\beta = \gamma(1 - S(\theta))$ , discouraging convergence to regions of high  
679 gradient disagreement. Unlike methods that directly modify the gradient direction (PCGrad, Fish) or  
680 solve auxiliary optimization problems (Nash-MTL, FAMO), SAGE’s noise injection is lightweight  
681 and requires no per-environment gradient storage beyond what is needed for computing the pairwise  
682 cosine similarity. This noise injection was inspired from the recent GGA [15] work.

683 Importantly, SAGE operates as a drop-in replacement for the optimizer’s ascent and descent steps,  
684 making it compatible with existing gradient manipulation methods in MTL (as demonstrated by the  
685 SAGE-MGDA, SAGE-FairGrad, and SAGE-SAMO combinations in Table 2) and with the standard  
686 DomainBed protocol in DG (Table 1).

## 687 B Proofs

### 688 B.1 Setup Reminder:

689 For each  $e \in \mathcal{E}$ , let  $\mathcal{L}_e : \Theta \rightarrow \mathbb{R}$  denote a per-distribution loss, where  $\Theta \subseteq \mathbb{R}^d$  is an open parameter  
690 space. The *population risk* is:

$$R(\theta) := \mathbb{E}_{e \sim \mathcal{P}}[\mathcal{L}_e(\theta)], \quad (8)$$

691 and the *empirical multi-distribution risk* over  $K$  iid training distributions  $e_1, \dots, e_K \sim \mathcal{P}$  is

$$\hat{R}(\theta) := \frac{1}{K} \sum_{k=1}^K \mathcal{L}_{e_k}(\theta). \quad (9)$$

692 At any base point  $\theta_0 \in \Theta$  we introduce the per-distribution gradient and Hessian,

$$g_e := \nabla_{\theta} \mathcal{L}_e(\theta_0), \quad H_e := \nabla_{\theta}^2 \mathcal{L}_e(\theta_0), \quad (10)$$

693 and the following cross-distribution statistics,

$$\bar{g} := \mathbb{E}_e[g_e], \quad \bar{H} := \mathbb{E}_e[H_e], \quad \Sigma_g := \text{Cov}_e(g_e) = \mathbb{E}_e[(g_e - \bar{g})(g_e - \bar{g})^T]. \quad (11)$$

694 **B.2 Proof of Theorem 1: Multi-distribution excess-risk decomposition**

695 *Proof of Theorem 1.* We restate the Theorem and Assumption included in the main text for complete-  
696 ness.

**Theorem 1** (Multi-distribution excess-risk decomposition). *Under Assumption 1, let  $\hat{\theta}$  denote the minimizer of  $\hat{R}$  over  $K$  iid training distributions, and let  $\xi \sim \mathcal{N}(0, \sigma^2 I)$  denote isotropic Gaussian noise of scale  $\sigma > 0$ . Then, for a parameter  $\theta = \hat{\theta} + \xi$ :*

$$\mathbb{E}[\mathbb{E}_\xi[R(\theta)]] - R(\theta^*) = \underbrace{\frac{1}{2K} \text{tr}(\bar{H}^{-1} \Sigma_g)}_{\text{alignment term}} + \underbrace{\frac{\sigma^2}{2} \text{tr}(\bar{H})}_{\text{curvature term}} + O(K^{-3/2}) + O(\sigma^3), \quad (12)$$

where the outer expectation is over the iid sampling of training distributions,  $\bar{H} := \nabla^2 R(\theta^*)$ , and  $\Sigma_g := \text{Cov}_e(\nabla \mathcal{L}_e(\theta^*))$ .

697

698 **Assumption 1** (Regularity). Each  $\mathcal{L}_e$  is three times continuously differentiable in  $\theta$ , with third  
699 derivatives uniformly bounded in a neighborhood of  $\theta^* := \arg \min_\theta R(\theta)$ , and  $\bar{H}(\theta^*) \succ 0$  (Positive-  
700 Definite). The gradients  $g_e(\theta^*)$  have finite second moments under  $\mathcal{P}$ .

701 Assumption 1 is the standard regularity condition under which quadratic excess-risk expansions are  
702 valid, which is however strictly weaker than the convexity assumptions used in classical M-estimation  
703 theory [31]. The  $C^3$  requirement is used only to calculate and control the remainder in a second-order  
704 Taylor expansion, while the PD-Hessian requirement ensures the quadratic form  $\bar{H}^{-1}$  is well defined.

705 The proof proceeds in three steps: (i) expand  $R$  around  $\theta^*$  to obtain the excess risk as a quadratic  
706 form in  $\hat{\theta} - \theta^*$ , (ii) characterize  $\hat{\theta} - \theta^*$  as a sample-mean quantity and take expectations to produce  
707 the alignment term, (iii) expand  $R$  around  $\hat{\theta}$  under random perturbations  $\xi \sim \mathcal{N}(0, \sigma^2 I)$  to produce  
708 the curvature term.

709 **Step 1: Quadratic expansion of  $R$  at  $\theta^*$ .** Because  $\theta^*$  minimizes the population risk  $R$ , we have  
710  $\nabla R(\theta^*) = \bar{g} = 0$ . Under Assumption 1, Taylor's theorem gives, for any  $\delta \in \mathbb{R}^d$  sufficiently small,

$$R(\theta^* + \delta) = R(\theta^*) + \nabla R(\theta^*)\delta + \frac{1}{2}\delta^T \nabla^2 R(\theta^*)\delta + r_3(\delta), \quad (13)$$

711 where  $|r_3(\delta)| \leq C\|\delta\|^3$  for some constant  $C$  depending on the uniform bound on third derivatives,  
712 by definition  $\nabla R(\theta^*) = \bar{g} = 0$ , and  $\nabla^2 R(\theta^*) = \bar{H}$ . Therefore, subtracting  $R(\theta^*)$  yields:

$$R(\theta^* + \delta) - R(\theta^*) = \frac{1}{2}\delta^T \bar{H}\delta + O(\|\delta\|^3). \quad (14)$$

713 Equation (14) is the standard quadratic excess-risk formula, which holds for any  $\delta$ . We now look at  
714  $\delta = \hat{\theta} - \theta^*$  and analyze the distance between the empirical minimizer ( $\hat{\theta}$ ) and population minimizer  
715 ( $\theta^*$ ).

716 **Step 2: Characterization of  $\hat{\theta} - \theta^*$ .** Note that the empirical minimizer satisfies  $\nabla \hat{R}(\hat{\theta}) = 0$  across  
717  $K$  distributions. Expanding  $\nabla \hat{R}$  around  $\theta^*$  using Taylor's theorem yields,

$$\nabla \hat{R}(\hat{\theta}) = \nabla \hat{R}(\theta^*) + \nabla^2 \hat{R}(\theta^*)(\hat{\theta} - \theta^*) + O(\|\hat{\theta} - \theta^*\|^2), \quad (15)$$

718 or,

$$0 = \nabla \hat{R}(\hat{\theta}) = \hat{g} + \hat{H}(\hat{\theta} - \theta^*) + O(\|\hat{\theta} - \theta^*\|^2), \quad (16)$$

719 where  $\hat{g} := \frac{1}{K} \sum_{k=1}^K g_{e_k}(\theta^*)$  and  $\hat{H} := \frac{1}{K} \sum_{k=1}^K H_{e_k}(\theta^*)$ .

720 Solving for  $\hat{\theta} - \theta^*$ :

$$\begin{aligned} \hat{H}(\hat{\theta} - \theta^*) &= -\hat{g} - O(\|\hat{\theta} - \theta^*\|^2) \\ \hat{\theta} - \theta^* &= -\hat{H}^{-1}\hat{g} - \hat{H}^{-1}O(\|\hat{\theta} - \theta^*\|^2) \\ \hat{\theta} - \theta^* &= -\hat{H}^{-1}\hat{g} + O(\|\hat{\theta} - \theta^*\|^2) \end{aligned} \quad (17)$$

721 Where from Assumption 1,  $\bar{H}$  is PD, so assuming enough data so will  $\hat{H}$ . To find the “size” of the  
722  $(\hat{\theta} - \theta^*) \approx -\hat{H}^{-1}\hat{g}$ , we take the vector norm (magnitude) of both sides. Using the standard property  
723 of matrix norms and the triangle inequality, we can separate the matrix from the vector to create an  
724 upper bound:

$$\|\hat{\theta} - \theta^*\| \approx \|\hat{H}^{-1}\hat{g}\| \leq \|\hat{H}^{-1}\| \cdot \|\hat{g}\|. \quad (18)$$

725 The eigenvalues of an inverse matrix are simply  $1/\lambda$ . Because the original eigenvalues ( $\lambda$ ) of  $\hat{H}$  are  
726 bounded away from zero, the eigenvalues of the inverse matrix cannot blow up to infinity. As the  
727 those eigenvalues are capped,  $\|\hat{H}^{-1}\|$  is just some finite constant number  $c$  (based on the largest  
728 eigenvalue). Therefore:

$$\begin{aligned} \|\hat{\theta} - \theta^*\| &\leq c\|\hat{g}\| \\ \|\hat{\theta} - \theta^*\| &\leq O(\|\hat{g}\|) \\ \|\hat{\theta} - \theta^*\|^2 &\leq O(\|\hat{g}\|^2) \end{aligned} \quad (19)$$

729 And substituting 19 to 17, yields:

$$\hat{\theta} - \theta^* = -\hat{H}^{-1}\hat{g} + O(\|\hat{g}\|^2) \quad (20)$$

730 This is the standard implicit-function argument used in M-estimator asymptotics [31]. We now  
731 continue to replace  $\hat{H}$  with  $\bar{H}$ , which concentrates to  $\|\hat{H} - \bar{H}\|$ , as  $K \rightarrow \infty$ . By the law of  
732 large numbers, as  $K$  increases, the empirical average converges to the true expectation:  $\hat{H} \rightarrow \bar{H}$ .  
733 Recall that  $\hat{H}$  is the empirical Hessian, calculated by averaging the actual Hessians observed across  
734  $K$  randomly sampled training distributions:  $\hat{H} = \frac{1}{K} \sum_{k=1}^K H_{e_k}$ . Because the training distributions  
735  $e_1, \dots, e_K$  are sampled independently and identically (iid) from the meta-distribution  $\mathcal{P}$ , each  
736 observed Hessian  $H_{e_k}$  is an independent random matrix. This means  $\hat{H}$  is simply a sample mean,  
737 and we want to know how far this sample mean deviates from the true mean  $\bar{H}$ . Let  $\hat{H}_{ij}$  denote the  
738 scalar entry in the  $i$ -th row and  $j$ -th column of the empirical Hessian matrix, and  $\bar{H}_{ij}$  denote the  
739 corresponding entry in the population Hessian. By definition, the empirical entry is a simple sample  
740 mean of the individual distribution Hessians:

$$\hat{H}_{ij} = \frac{1}{K} \sum_{k=1}^K (H_{e_k})_{ij}.$$

741 Since the third derivatives of the loss are uniformly bounded in a neighborhood of  $\theta^*$ , as per 1, each  
742 entry  $(H_{e_k})_{ij}$  is itself bounded, and there exists some finite variance  $\sigma_{ij}^2 < \infty$  such that:

$$\text{Var}((H_{e_k})_{ij}) = \sigma_{ij}^2.$$

743 Since each distribution is sampled i.i.d, from the Central Limit Theorem (or Chebyshev’s Inequality)  
744 for each element; as  $K \rightarrow \infty$ :

$$\sqrt{K}(\hat{H}_{ij} - \bar{H}_{ij}) \rightarrow \mathcal{N}(0, \sigma_{ij}^2).$$

745 Where by calculating the Frobenius norm of the difference, it can be easily shown that  $\|\hat{H} - \bar{H}\|_F =$   
746  $O_P(K^{-1/2})$  and we can formally bound the random error matrix:

$$\|\hat{H} - \bar{H}\| = O_P(K^{-1/2}) \quad (21)$$

747 Therefore the distance between the empirical sample  $\hat{H}$  and the true population  $\bar{H}$  shrinks at a rate  
748 proportional to  $1/\sqrt{K}$ .

749 Again, following the same logic, by the Central Limit Theorem  $\sqrt{K}\hat{g} \rightarrow \mathcal{N}(0, \Sigma_g)$ , and the magni-  
 750 tude of the average gradient  $\|\hat{g}\|$  shrinks at a rate of  $O_P(K^{-1/2})$  and in turn the squared gradient  
 751 term becomes  $O(\|\hat{g}\|^2) = O_P(K^{-1})$ .

752 Thus, Equation 20 becomes:

$$\hat{\theta} - \theta^* = -\bar{H}^{-1}\hat{g} + O_P(K^{-1}). \quad (22)$$

753 **Taking expectations.** Substitute Eq. (22) into Eq. (14):

$$\begin{aligned} R(\hat{\theta}) - R(\theta^*) &= \frac{1}{2}(\bar{H}^{-1}\hat{g})^T \bar{H}(\bar{H}^{-1}\hat{g}) + O_P(K^{-3/2}) \\ &= \frac{1}{2}\hat{g}^T \bar{H}^{-1}\hat{g} + O_P(K^{-3/2}), \end{aligned} \quad (23)$$

754 where we used  $\bar{H}^{-1}\bar{H}\bar{H}^{-1} = \bar{H}^{-1}$  and the symmetry  $(\bar{H}^{-1})^T = \bar{H}^{-1}$ . The  $\|\delta\|^3$  term in Eq. (14)  
 755 becomes  $O_P(K^{-3/2})$  because  $\|\hat{g}\| = O_P(K^{-1/2})$  by the central limit theorem. Now take expectation  
 756 over the iid sampling of the  $K$  training distributions. Using the matrix identity

$$\mathbb{E}[x^T A x] = \text{tr}(A \text{Cov}(x)) + \mathbb{E}[x]^T A \mathbb{E}[x] \quad (24)$$

757 (valid for any random vector  $x$  with finite second moments and any symmetric  $A$ ) with  $x = \hat{g}$  and  
 758  $A = \bar{H}^{-1}$ , we compute:

- 759 •  $\mathbb{E}[\hat{g}] = \mathbb{E}[\frac{1}{K} \sum_k g_{e_k}] = \bar{g} = 0$  (since  $\theta^*$  is stationary for  $R$ ), so the mean term vanishes.
- 760 •  $\text{Cov}(\hat{g}) = \frac{1}{K^2} \sum_k \text{Cov}(g_{e_k}) = \frac{1}{K} \Sigma_g$  by independence of the distributions and due to the  
 761 fact that each distribution is drawn from  $\mathcal{P}$ .

762 Substituting:

$$\mathbb{E}[\hat{g}^T \bar{H}^{-1} \hat{g}] = \text{tr}(\bar{H}^{-1} \cdot \frac{1}{K} \Sigma_g) = \frac{1}{K} \text{tr}(\bar{H}^{-1} \Sigma_g). \quad (25)$$

763 Combining with Eq. (23):

$$\mathbb{E}[R(\hat{\theta}) - R(\theta^*)] = \frac{1}{2K} \text{tr}(\bar{H}^{-1} \Sigma_g) + O(K^{-3/2}). \quad (26)$$

764 This establishes the alignment term.

765 **Step 3: Expansion of  $R$  at  $\hat{\theta}$  under random perturbations.** We want to account for the fact that  
 766 the learner does not exactly realize  $\hat{\theta}$  and to capture the sharpness of the local landscape, we evaluate  
 767 the expected risk of the empirical minimizer under random perturbations, e.g. due to optimization  
 768 noise or finite-precision arithmetic. Let the deployed parameter be  $\theta = \hat{\theta} + \xi$ , where  $\xi \sim \mathcal{N}(0, \sigma^2 I)$   
 769 is an isotropic Gaussian perturbation of scale  $\sigma > 0$ .

770 By expanding  $R$  at  $\hat{\theta}$ , from Taylor's theorem:

$$R(\hat{\theta} + \xi) = R(\hat{\theta}) + \nabla R(\hat{\theta})\xi + \frac{1}{2}\xi^T \nabla^2 R(\hat{\theta})\xi + O(\|\xi\|^3), \quad (27)$$

771 and taking expectation with respect to the noise  $\xi$ ,

$$\mathbb{E}_\xi[R(\hat{\theta} + \xi)] = R(\hat{\theta}) + \nabla R(\hat{\theta})\mathbb{E}[\xi] + \frac{1}{2}\mathbb{E}[\xi^T \nabla^2 R(\hat{\theta})\xi] + O(\sigma^3) \quad (28)$$

772 The linear term vanishes because the perturbation is zero-mean ( $\mathbb{E}[\xi] = 0$ ). For the quadratic term,  
 773 apply the identity (24) with  $x = \xi$  (zero mean) and  $A = \nabla^2 R(\hat{\theta})$ , using  $\text{Cov}(\xi) = \sigma^2 I$ :

$$\frac{1}{2}\mathbb{E}[\xi^T \nabla^2 R(\hat{\theta})\xi] = \frac{\sigma^2}{2} \text{tr}(\nabla^2 R(\hat{\theta})) \quad (29)$$

774 Since  $\nabla^2 R(\hat{\theta}) = \bar{H} + O_P(K^{-1/2})$  by continuity of the Hessian and the fact that  $\hat{\theta} \rightarrow \theta^*$ , we can  
 775 replace  $\nabla^2 R(\hat{\theta})$  with the population Hessian  $\bar{H}$  up to a correction that multiplies the existing  $\sigma^2/2$   
 776 factor and contributes at order  $O(\sigma^2 K^{-1/2})$ , which is absorbed into the stated error. Thus, the  
 777 expected risk under perturbation  $\xi$  is:

$$\mathbb{E}_\xi[R(\hat{\theta})] = R(\hat{\theta}) + \frac{\sigma^2}{2} \text{tr}(\bar{H}) + O(\sigma^3) + O(\sigma^2 K^{-1/2}) \quad (30)$$

778 **Combining.** Taking expectation of Eq. (30) over the distribution sampling and adding Eq. (26):

$$\mathbb{E}[\mathbb{E}_\xi[R(\hat{\theta})]] - R(\theta^*) = \frac{1}{2K} \text{tr}(\bar{H}^{-1}\Sigma_g) + \frac{\sigma^2}{2} \text{tr}(\bar{H}) + O(K^{-3/2}) + O(\sigma^3), \quad (31)$$

779 which is Eq. (12).  $\square$   $\square$

### 780 B.3 Proof of Counterexample 1: Decoupling of flatness and gradient alignment

781 *Proof of Counterexample 1.* We restate the counterexample included in the main text for complete-  
782 ness.

**Counterexample 1** (Decoupling of flatness and gradient alignment). Consider the family of multi-distribution learning problems with per-distribution losses

$$\mathcal{L}_e(\theta) = \frac{1}{2}\theta^T A\theta + b_e^T \theta, \quad \mathbb{E}_e[b_e] = 0, \quad (32)$$

parameterized by a symmetric PD matrix  $A \in \mathbb{R}^{d \times d}$  and a collection of linear coefficients  $\{b_e\}$  satisfying  $\mathbb{E}_e[b_e] = 0$ . For every  $M > 0$ , there exist instances of Eq. (32) such that:

- (i) **Flat but misaligned:**  $\text{tr}(\bar{H}) \leq M^{-1}$  and  $\text{tr}(\bar{H}^{-1}\Sigma_g) \geq M$ .
- (ii) **Aligned but sharp:**  $\text{tr}(\bar{H}^{-1}\Sigma_g) \leq M^{-1}$  and  $\text{tr}(\bar{H}) \geq M$ .

In particular, neither  $\text{tr}(\bar{H})$  nor  $\text{tr}(\bar{H}^{-1}\Sigma_g)$  can be bounded above by any function of the other alone.

783

784 For the family in the counterexample we have  $\nabla \mathcal{L}_e(\theta) = A\theta + b_e$  and  $\nabla^2 \mathcal{L}_e(\theta) = A$ . Therefore  
785  $\bar{H} = A$  is fixed by the choice of  $A$  alone, and the population minimizer is  $\theta^* = -A^{-1}\mathbb{E}_e[b_e] = 0$ .  
786 At  $\theta^*$ , the per-environment gradients are  $g_e = b_e$ , so

$$\Sigma_g = \mathbb{E}_e[b_e b_e^T], \quad (33)$$

787 which is fixed by the choice of  $\{b_e\}$  alone. The two quantities are thus independently controllable:  $A$   
788 determines  $\bar{H}$  without affecting  $\Sigma_g$ , and  $\{b_e\}$  determines  $\Sigma_g$  without affecting  $\bar{H}$ .

789 **Construction of (i): flat but misaligned.** Take  $d = 2$ ,  $A = \lambda I$  with  $\lambda = (2M)^{-1}$ , and let  $b_e$  be  
790 supported on  $e_2$  with variance  $v = 1$ . Then  $\bar{H} = \lambda I$ , so

$$\text{tr}(\bar{H}) = 2\lambda = M^{-1}, \quad (34)$$

791 satisfying the first inequality. Also  $\Sigma_g = \text{diag}(0, v)$ , so

$$\text{tr}(\bar{H}^{-1}\Sigma_g) = \text{tr}(\lambda^{-1}I \cdot \text{diag}(0, v)) = v/\lambda = 2Mv. \quad (35)$$

792 Choosing  $v = 1$  gives  $\text{tr}(\bar{H}^{-1}\Sigma_g) = 2M \geq M$ . Concretely, two environments with  $b_1 = (0, +1)^T$   
793 and  $b_2 = (0, -1)^T$  realize this  $\Sigma_g$ .

794 **Construction of (ii): aligned but sharp.** Take  $d = 2$ ,  $A = \lambda I$  with  $\lambda = M/2$ , and  $b_e$  supported  
795 on  $e_1$  with variance  $v = 1/2$ . Then

$$\text{tr}(\bar{H}) = 2\lambda = M, \quad \text{tr}(\bar{H}^{-1}\Sigma_g) = v/\lambda = M^{-1}. \quad (36)$$

796 Two environments with  $b_1 = (+1/\sqrt{2}, 0)^T$  and  $b_2 = (-1/\sqrt{2}, 0)^T$  realize the required  $\Sigma_g$ .

797 **Conclusion.** In both constructions,  $\bar{H}$  and  $\Sigma_g$  are chosen independently to drive one quantity  
798 arbitrarily small and the other arbitrarily large. If there existed a function  $\phi$  such that  $\text{tr}(\bar{H}^{-1}\Sigma_g) \leq$   
799  $\phi(\text{tr}(\bar{H}))$  for all instances in the family, then construction (i) with  $M \rightarrow \infty$  would contradict  
800  $\phi(M^{-1}) < \infty$ , symmetrically for the reverse direction.  $\square$   $\square$

801 **B.4 Limitations**

802 **Local, not global.** Theorem 1 is an asymptotic expansion valid in a neighborhood of  $\theta^*$  where the  
 803 quadratic approximation holds. In non-convex deep learning landscapes, there may be many  $\bar{g} = 0$   
 804 points, and the theorem does not address which one the optimizer reaches. This is the “hypothesis  
 805 selection” gap: among the manifold of stationary points, flatness and alignment act as refinement  
 806 criteria, but the decomposition alone does not prove that optimization converges to the jointly-optimal  
 807 refinement.

808 **Not a PAC-Bayes bound.** Although Step 3 of the proof of Theorem 1 echoes the structure of  
 809 Gaussian PAC-Bayes analyses the result is an *expectation-form* excess-risk decomposition, not a  
 810 high-probability uniform bound. A full PAC-Bayes statement would add a KL-divergence term,  
 811 bound the failure probability explicitly, and require a different argument for the alignment term. The  
 812 result stated here is sufficient for its purpose, i.e, identifying which quantities appear in the excess  
 813 risk and in what functional form, but should not be read as a finite-sample generalization guarantee.

814 **Quadratic family is toy.** The decoupling construction in Counterexample 1 uses a purely quadratic  
 815 family, in which  $\bar{H}$  is globally constant and independent of  $\{b_e\}$ . A skeptical reader might object  
 816 that real deep networks are highly non-convex and that the independence of  $\bar{H}$  and  $\Sigma_g$  observed in  
 817 the quadratic family may not hold globally. The correct response is that Counterexample 1 provides  
 818 a strong indication, but no guarantees. If even in a simple setting (exact quadratic, infinite data,  
 819 closed-form) neither criterion alone suffices, then a fortiori neither is expected to suffice in the harder,  
 820 non-convex setting.

821 **C Motivating Example: Why Both Properties Are Necessary**

822 Theorem 1 decomposes the excess risk into an alignment term and a curvature term, while Theorem 1  
 823 shows that neither term can be bounded as a function of the other. The example below illustrates  
 824 how this is reflected in a concrete setting where each failure mode appears along an orthogonal  
 825 eigendirection of  $\bar{H}$ . The example uses a linear model on a Gaussian data-generating process with  
 826 two domains, chosen so that the loss landscape admits closed-form analysis.

827 **Data-generating process.** Consider a binary classification task with label  $y \in \{-1, +1\}$  drawn  
 828 uniformly, and input  $x = (x_{\text{inv}}, x_{\text{spur}}) \in \mathbb{R}^2$  consisting of an invariant feature and a spurious feature.  
 829 For domain  $k \in \{1, 2\}$ , the features are generated conditionally on  $y$  as:

$$x_{\text{inv}} | y \sim \mathcal{N}(y \cdot \mu_{\text{inv}}, \sigma_{\text{inv}}^2), \tag{37}$$

$$x_{\text{spur}}^{(k)} | y \sim \mathcal{N}(y \cdot c_k, \sigma_{\text{spur}}^2), \tag{38}$$

830 where  $c_1 = +\mu_{\text{spur}}$  and  $c_2 = -\mu_{\text{spur}}$ . The invariant feature is causally linked to  $y$  identically in both  
 831 domains, while the spurious feature’s correlation with  $y$  reverses across domains. Crucially, we set  
 832  $\sigma_{\text{inv}}^2 \gg \sigma_{\text{spur}}^2$ : the invariant feature is noisy (high variance) while the spurious feature is clean (low  
 833 variance). We adopt the concrete values  $\mu_{\text{inv}} = 1$ ,  $\sigma_{\text{inv}}^2 = 9$ ,  $\mu_{\text{spur}} = 2$ ,  $\sigma_{\text{spur}}^2 = 0.01$  throughout the  
 834 example.

835 **Loss landscape.** We train a linear model  $f(x; \theta) = \theta_{\text{inv}} x_{\text{inv}} + \theta_{\text{spur}} x_{\text{spur}}$  with mean squared error.  
 836 The expected loss for domain  $k$  expands to:

$$\mathcal{L}_k(\theta) = \frac{1}{2} - \theta_{\text{inv}} \mu_{\text{inv}} - \theta_{\text{spur}} c_k + \frac{1}{2} \theta_{\text{inv}}^2 (\sigma_{\text{inv}}^2 + \mu_{\text{inv}}^2) + \frac{1}{2} \theta_{\text{spur}}^2 (\sigma_{\text{spur}}^2 + \mu_{\text{spur}}^2) + \theta_{\text{inv}} \theta_{\text{spur}} \mu_{\text{inv}} c_k, \tag{39}$$

837 using the moments  $\mathbb{E}[y^2] = 1$ ,  $\mathbb{E}[y x_{\text{inv}}] = \mu_{\text{inv}}$ ,  $\mathbb{E}[x_{\text{inv}}^2] = \sigma_{\text{inv}}^2 + \mu_{\text{inv}}^2$ , and so on. Because each  
 838 sub-objective loss is quadratic, the Hessian of  $\mathcal{L}_k$  is constant:

$$H_k = \begin{pmatrix} \sigma_{\text{inv}}^2 + \mu_{\text{inv}}^2 & \mu_{\text{inv}} c_k \\ \mu_{\text{inv}} c_k & \sigma_{\text{spur}}^2 + \mu_{\text{spur}}^2 \end{pmatrix}. \tag{40}$$

839 Substituting our numerical values gives  $H_{11} = 10$ ,  $H_{22} = 4.01$ , and  $H_{12}^{(1)} = +2$ ,  $H_{12}^{(2)} = -2$ . Two  
840 features of this landscape are critical. First, the curvature along  $\theta_{\text{inv}}$  is  $2.5\times$  larger than along  $\theta_{\text{spur}}$ ,  
841 meaning the invariant direction forms a narrow valley while the spurious direction is comparatively  
842 flat. Second, the off-diagonal term  $\mu_{\text{inv}C_k}$  flips sign across domains, coupling the two parameters  
843 differently in each domain’s loss surface.

844 **Aggregate loss and its minima.** The aggregate (ERM) loss  $\bar{\mathcal{L}} = \frac{1}{2}(\mathcal{L}_1 + \mathcal{L}_2)$  has the Hessian:

$$\bar{H} = \frac{1}{2}(H_1 + H_2) = \begin{pmatrix} 10 & 0 \\ 0 & 4.01 \end{pmatrix}, \quad (41)$$

845 where the cross-domain averaging cancels the off-diagonal terms. The aggregate minimum is  
846  $\theta^* = \bar{H}^{-1}\nabla_{\theta}\bar{\mathcal{L}}(0) = (0.1, 0)^T$ , confirming that the ERM optimum correctly assigns zero weight to  
847 the spurious feature.

848 However, the individual domain minima tell a different story. Setting  $\nabla_{\theta}\mathcal{L}_k = 0$  and solving yields  
849 domain-specific optima that place substantial weight on  $\theta_{\text{spur}}$  with opposite signs across domains. In  
850 each domain individually, the spurious feature is a low-noise, high-signal predictor and can lead to  
851 training models that will not generalize.

852 **Failure mode 1: Flatness without alignment.** A method that searches for low-curvature solutions  
853 within a single domain finds the per-domain optima  $\theta_k^* = H_k^{-1}b_k \approx (0.0003, \pm 0.499)^T$ , which  
854 lie almost entirely along the spurious axis with reversed signs across domains. This is a structural  
855 consequence of the within-domain feature statistics, as the spurious feature has signal-to-noise ratio  
856  $\mu_{\text{spur}}^2/\sigma_{\text{spur}}^2 = 400$ , far exceeding the invariant feature’s  $\mu_{\text{inv}}^2/\sigma_{\text{inv}}^2 \approx 0.11$ . Within either domain  
857 individually, the spurious feature is the better predictor. A flatness-only criterion cannot distinguish a  
858 flat minimum that uses spurious features from one that uses invariant ones, and the alignment term of  
859 Theorem 1 is what penalizes the difference.

860 **Failure mode 2: Alignment without flatness.** A method that enforces gradient agreement across  
861 domains correctly identifies the shared direction. The per-domain gradients at the origin are  
862  $\nabla\mathcal{L}_1(0) = (-1, -2)^T$  and  $\nabla\mathcal{L}_2(0) = (-1, +2)^T$ , which agree on the invariant component and  
863 conflict on the spurious one. The optimizer converges to  $\theta^* = (0.1, 0)^T$ , i.e. the aggregate minimum,  
864 which uses exclusively the invariant feature. However,  $\theta^*$  sits on the *sharp* eigendirection of  $\bar{H}$   
865 (eigenvalue 10), so a parameter displacement  $\delta$  along  $\theta_{\text{inv}}$  leads to a loss increase of  $\frac{1}{2} \cdot 10 \cdot \delta^2 = 5\delta^2$ .  
866 Therefore, under distribution shift the optimal coefficient moves and the steep curvature amplifies the  
867 resulting error. The curvature term of Theorem 1 is what penalizes this.

868 **Connection to Counterexample 1.** The two failure modes occur along orthogonal eigendirec-  
869 tions of  $\bar{H}$ , which is what makes this two-dimensional construction minimal. The cross-domain  
870 gradient covariance at the aggregate minimum is  $\Sigma_g = \text{diag}(0, 3.24)$ , supported entirely on the flat  
871 eigendirection of  $\bar{H}$ . The alignment term evaluates to  $\text{tr}(\bar{H}^{-1}\Sigma_g) = 3.24/4.01 \approx 0.81$ , which is  
872 dominated by the ratio of gradient variance to the small eigenvalue. This is, structurally, an instance  
873 of construction (i) of Counterexample 1, as the support of  $\Sigma_g$  coincides with the flat eigendirection  
874 of  $\bar{H}$ , so the alignment term is amplified by the inverse curvature in that direction. The two failure  
875 modes the example exhibits are not just numerical coincidences but the structural examples of the  
876 connection between  $\Sigma_g$  and  $\bar{H}$ .

## 877 D Scale Invariance and Scale-Dependent Flatness in SAM

878 A well-documented limitation that SAGE addresses is with regards to the interaction between SAM’s  
879 fixed perturbation radius and the scale invariance of modern architectures [34]. Neural networks  
880 employing ReLU activations and normalization layers exhibit positive scale invariance, meaning that  
881 scaling a weight matrix  $W$  by  $\alpha > 0$  to obtain  $W' = \alpha W$  (with the subsequent layer scaled by  $1/\alpha$ )  
882 leaves the network’s output unchanged. The loss landscape is therefore functionally identical, yet the  
883 gradient scales inversely:

$$\nabla_{W'}\mathcal{L} = \frac{1}{\alpha}\nabla_W\mathcal{L}. \quad (42)$$

884 Because SAM uses a fixed  $L_2$ -radius  $\rho$ , the perturbation’s size relative to the parameter norm  
 885 changes drastically with  $\alpha$ . Shrinking the weights ( $\alpha < 1$ ) makes the fixed-size perturbation  
 886 push parameters proportionally further from their current values, inflating the apparent sharpness;  
 887 conversely, scaling weights up ( $\alpha \gg 1$ ) makes the same perturbation negligibly small relative to the  
 888 parameter norm, producing artificially flat sharpness estimates. This phenomenon, termed “scale-  
 889 dependent flatness” [8], means that SAM’s sharpness measure conflates the true geometry of the loss  
 890 surface with the arbitrary scale of the parameterization.

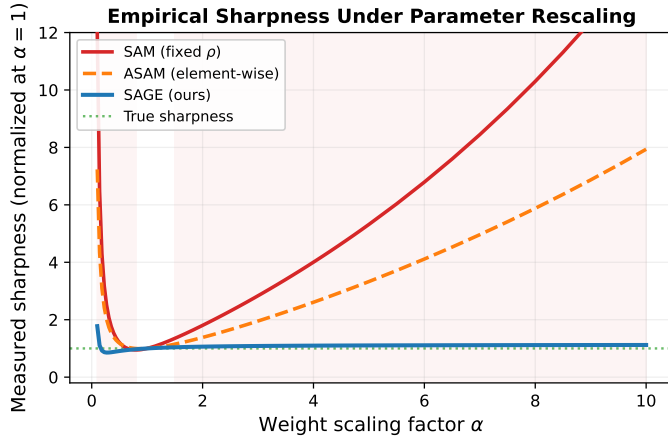


Figure 4: **Empirical scale invariance.** Measured sharpness as a function of weight rescaling factor  $\alpha$  for a two-layer MLP (no normalization layers). The network computes the same function at all  $\alpha$ , so the true sharpness (green, dotted) is constant. Both SAM (red) and ASAM (orange, dashed) produce scale-dependent sharpness estimates that diverge away from  $\alpha = 1$ , despite the underlying function being identical. SAGE (blue) is approximately invariant to the rescaling, correctly reflecting the true geometry.

## 891 D.1 Illustrative Example

892 To empirically demonstrate the scale-invariance failure, we train a two-layer MLP without normal-  
 893 ization layers ( $\text{Linear}(2, 64) \rightarrow \text{ReLU} \rightarrow \text{Linear}(64, 2)$ ) on a synthetic concentric-circles task  
 894 until convergence. At the converged parameters  $\theta^*$ , we apply the rescaling  $W_1 \mapsto \alpha W_1$ ,  $b_1 \mapsto \alpha b_1$ ,  
 895  $W_2 \mapsto W_2/\alpha$  for  $\alpha \in [0.1, 10.0]$ . As previously mentioned, the network computes the same function  
 896 for all  $\alpha$ . For this example, we measure the sharpness  $\mathcal{L}(\theta + \epsilon) - \mathcal{L}(\theta)$  at each  $\alpha$  using SAM (fixed  
 897  $\rho$ ), ASAM (element-wise adaptive), and SAGE (spectral, scale-adaptive). As shown in Figure 4, both  
 898 SAM and ASAM produce sharpness estimates that diverge away from  $\alpha = 1$ , while SAGE remains  
 899 approximately constant, correctly reflecting the invariance of the underlying loss geometry<sup>5</sup>.

## 900 E Experimental details & Hyperparameters

### 901 E.1 Details & Hyperparameters

902 **Domain Generalization** For DG, we follow the protocol of the widely adopted and challenging Do-  
 903 mainBed [41] benchmark, and conduct experiments on five image classification datasets. Specifically,  
 904 we follow the leave-one-domain-out evaluation protocol for PACS [42] (9,991 images, 4 domains and  
 905 7 classes), VLCS [43] (10,729 images, 4 domains and 5 classes), OfficeHome [44] (15,588 images,  
 906 4 domains and 65 classes), TerraIncognita [45] (24,788 images, 4 domains and 10 classes), and  
 907 DomainNet [46] (586,575 images, 6 domains and 345 classes), and report the average top-1 accuracy  
 908 over 3 runs based on training-domain split validation. Following previous implementations [51, 10],  
 909 we select a batch size of 12 where images are uniformly sampled across source domains, set an  
 910 initial learning rate of  $1e - 6$  and employ the Adam optimizer for finetuning a CLIP [52] pretrained

<sup>5</sup>The effect on SAM’s fixed- $\rho$  perturbation depends on which layer dominates. In this architecture,  $W_2$ ’s gradient grows and  $\alpha$  causes the perturbation to push harder on  $W_2$ , thus increasing sharpness for large  $\alpha$ .

911 ViT-B/16 [53]. Regarding SAGE-specific hyperparameters, we select  $\rho = 1e - 6$  and  $\gamma = 1e - 6$   
912 across all datasets.

913 **Multi-Task Learning** For the MTL setting, we follow the standard protocol in recent literature [11,  
914 47, 48] and evaluate on Cityscapes [49], a dataset of urban street scenes containing 5,000 images with  
915 pixel-level annotations, and it supports two tasks: 7-class semantic segmentation and depth estimation.  
916 In contrast, NYU-v2 [50] focuses on indoor environments, providing 1,449 densely annotated images  
917 and enabling three tasks: 13-class semantic segmentation, depth estimation, and surface normal  
918 prediction. For training, we employ MTAN [55] as the shared backbone, with task-specific attention  
919 modules built on top of a SegNet [56]. Following previous implementations [11, 47] the model is  
920 trained for 200 epochs with a batch size of 8 for Cityscapes and 2 for NYU-v2. The learning rate  
921 is set to  $1e - 4$  for the first 100 epochs and is then halved for the remainder. Regarding SAGE  
922 hyperparameters, we set  $\rho = 1e - 4$  and  $\gamma = 1e - 5$  for both datasets.

## 923 E.2 Infrastructure

924 All experiments were conducted on a cluster containing  $4 \times 40$  GB NVIDIA A100 GPU cards, split  
925 into 8 20GB virtual MIG devices and  $1 \times 24$ GB NVIDIA RTX A5000 GPU card, via a SLURM  
926 workload manager.

## 927 F Additional Ablations

928 In Figure 5 we provide plots of the cosine similarity between the average training-domain gradient  
929 and the gradient on the unseen target during training for the DomainBed datasets. The similarity is  
930 calculated at every training step, while the lines are smoothed for better visibility.

## 931 G Computational Analysis

932 Each SAGE training iteration requires two forward-backward passes through the network, identical  
933 in structure to SAM and all SAM-like methods. The first pass computes the aggregate and per-  
934 environment gradients  $g_{e_k}$  from a single forward call over the concatenated mini-batches from all  
935  $K$  distributions (domains in DG or tasks in MTL), incurring no additional forward passes beyond  
936 what ERM already performs with the same batch. From these per-distribution gradients, the pairwise  
937 cosine similarity  $S(\theta)$  and the noise scale  $\beta$  are computed in  $\mathcal{O}(K^2d)$  time, which is negligible  
938 relative to the cost of backpropagation for any practical number of environments ( $K \leq 6$  in all of  
939 our experiments). The spectral perturbation replaces SAM’s  $L_2$ -normalization of the gradient with  
940  $T = 5$  Newton-Schulz iterations per weight matrix, each consisting of a single matrix multiplication  
941 and a matrix subtraction. These operations are fully parallelized on modern GPU hardware and add  
942 negligible overhead. The second forward-backward pass at the perturbed point  $\theta + \epsilon$  is identical  
943 to that of SAM. Finally, the noise injection  $\beta \xi$  at the descent step amounts to sampling a Gaussian  
944 vector and a single multiply-add operation, which is  $\mathcal{O}(d)$ . In total, SAGE’s per-iteration cost is  
945 approximately  $2 \times 2$  that of ERM, the same factor as SAM, with the additional operations (cosine  
946 similarities, Newton-Schulz iterations, and noise sampling) contributing an overhead that is negligible  
947 in practice compared to the forward-backward passes. Importantly, inference is entirely unaffected  
948 by SAGE, as all modifications occur exclusively during the training optimization loop.

## 949 H Pseudo code of the proposed SAGE algorithm

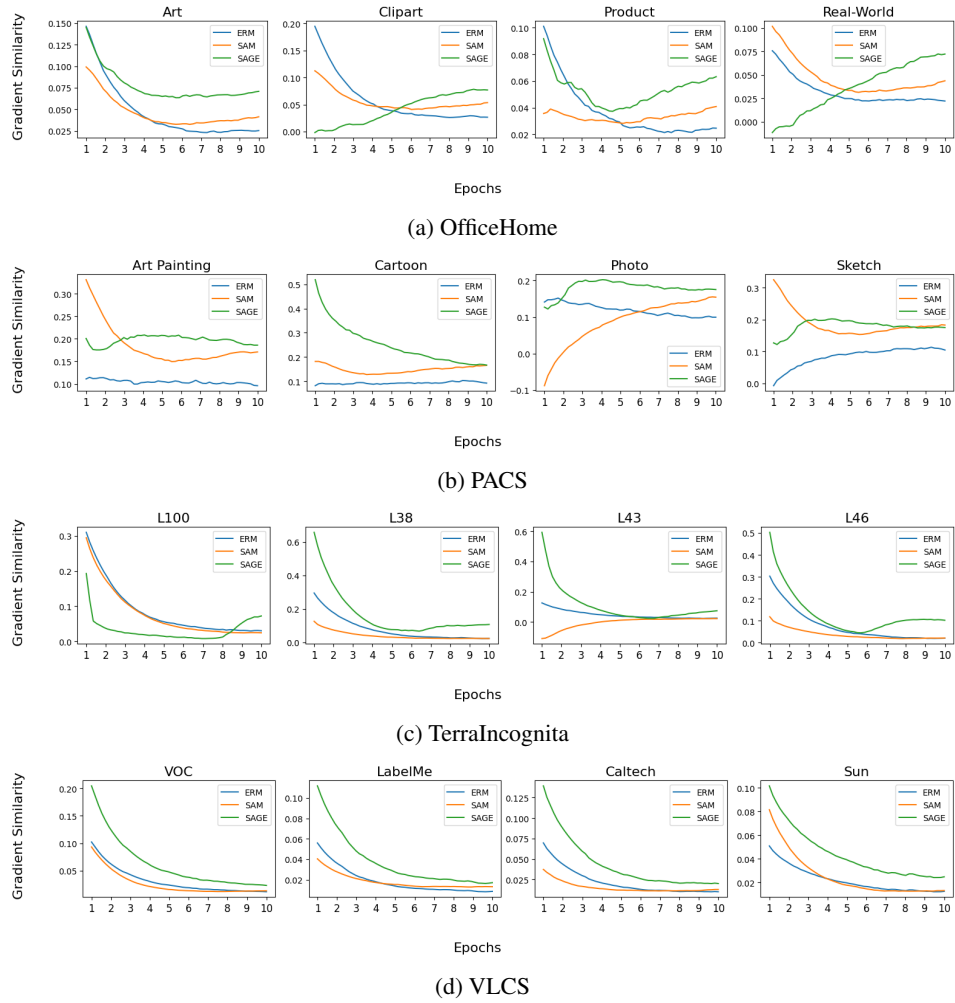


Figure 5: Gradient alignment for unseen domains during model training on the datasets of DomainBed.

---

**Algorithm 1** SAGE: Spectral-Aware and Gradient-Aligned Exploration
 

---

**Require:** Learning rate  $\eta$ , perturbation radius  $\rho$ , noise scale  $\gamma$ , sub-objectives  $\{D_1, \dots, D_K\}$ , feature extractor  $f_\theta$ , cosine similarity function (*cos*), Newton-Schulz iterations  $T$ , loss function  $\ell$ , base optimizer  $\text{OPT}()$ , total number of training operations  $n$ .

```

1: for  $t \leftarrow 1$  to  $n$  do
2:   Sample mini-batches  $\{(x_k, y_k)\}_{k=1}^K$  from each sub-objective
3:   // Phase 1: Sub-objective and total gradients
4:   Compute logits  $\hat{y} = f_\theta(\text{concat}(x_1, \dots, x_K))$ 
5:   for  $k = 1, \dots, K$  do
6:      $\mathcal{L}_k \leftarrow \ell(\hat{y}_k, y_k)$ 
7:      $g_k \leftarrow \nabla_\theta \mathcal{L}_k$ 
8:   end for
9:    $\bar{g} \leftarrow \sum_k w_k g_k$  { sample-weighted total gradient }
10:  // Phase 2: Calculate gradient agreement & noise scale
11:   $S \leftarrow \frac{2}{K(K-1)} \sum_{i < j} \cos(g_i, g_j)$ 
12:   $\beta \leftarrow \gamma (1 - S)$ 
13:  // Phase 3: Compute Spectral Perturbations (Muon Logic)
14:  for each weight matrix  $W \in \theta$  and its corresponding gradient  $G \in g$  do
15:    if  $G$  is a matrix or higher-order tensor then
16:       $G_{\text{ortho}} \leftarrow \text{NewtonSchulz5}(G)$  { Orthogonalize via Newton-Schulz }
17:       $\epsilon_W \leftarrow \rho \|W\|_F G_{\text{ortho}}$  { Scaled spectral perturbation }
18:    else
19:       $\epsilon_W \leftarrow \rho \frac{G}{\|G\|_2}$ 
20:    end if
21:  end for
22:  // Phase 4: Apply Perturbation (Ascent Step)
23:   $\theta \leftarrow \theta + \epsilon$ 
24:  // Phase 5: Adversarial Forward & Backward Pass
25:  Compute perturbed loss  $\mathcal{L}_{\text{pert}} \leftarrow \ell(f_\theta(\text{concat}(x_1, \dots, x_K)))$ 
26:  Compute perturbed gradient  $g_{\text{pert}} \leftarrow \nabla_\theta \mathcal{L}_{\text{pert}}$ 
27:  // Phase 6: Restore Weights & Inject Gradient Noise (Descent Step)
28:   $\theta \leftarrow \theta - \epsilon$  { Restore clean weights }
29:   $g_{\text{final}} \leftarrow g_{\text{pert}} + \beta \mathcal{N}(0, \mathbf{I})$  { Inject domain-conflict noise }
30:  // Phase 7: Optimizer Step
31:   $\theta \leftarrow \text{OPT}(\theta, \eta, g_{\text{final}})$ 
32: end for

```

---

## 950 **NeurIPS Paper Checklist**

951 The checklist is designed to encourage best practices for responsible machine learning research,  
952 addressing issues of reproducibility, transparency, research ethics, and societal impact. Do not remove  
953 the checklist: **The papers not including the checklist will be desk rejected.** The checklist should  
954 follow the references and follow the (optional) supplemental material. The checklist does NOT count  
955 towards the page limit.

956 Please read the checklist guidelines carefully for information on how to answer these questions. For  
957 each question in the checklist:

- 958 • You should answer [Yes], [No], or [N/A].
- 959 • [N/A] means either that the question is Not Applicable for that particular paper or the  
960 relevant information is Not Available.
- 961 • Please provide a short (1–2 sentence) justification right after your answer (even for [N/A]).

962 **The checklist answers are an integral part of your paper submission.** They are visible to the  
963 reviewers, area chairs, senior area chairs, and ethics reviewers. You will also be asked to include it  
964 (after eventual revisions) with the final version of your paper, and its final version will be published  
965 with the paper.

966 The reviewers of your paper will be asked to use the checklist as one of the factors in their evaluation.  
967 While [Yes] is generally preferable to [No], it is perfectly acceptable to answer [No] provided a  
968 proper justification is given (e.g., error bars are not reported because it would be too computationally  
969 expensive” or “we were unable to find the license for the dataset we used”). In general, answering  
970 [No] or [N/A] is not grounds for rejection. While the questions are phrased in a binary way, we  
971 acknowledge that the true answer is often more nuanced, so please just use your best judgment and  
972 write a justification to elaborate. All supporting evidence can appear either in the main paper or the  
973 supplemental material, provided in appendix. If you answer [Yes] to a question, in the justification  
974 please point to the section(s) where related material for the question can be found.

975 **IMPORTANT, please:**

- 976 • **Delete this instruction block, but keep the section heading “NeurIPS Paper Checklist”,**
- 977 • **Keep the checklist subsection headings, questions/answers and guidelines below.**
- 978 • **Do not modify the questions and only use the provided macros for your answers.**

### 979 **1. Claims**

980 Question: Do the main claims made in the abstract and introduction accurately reflect the  
981 paper’s contributions and scope?

982 Answer: [Yes]

983 Justification: All claims made in the Abstract and Introduction accurately reflect the paper’s  
984 contribution and scope. Specifically, for the 4 contributions presented in the Introduction  
985 are respectively reflected in Sections 2, 3 and 5.

986 Guidelines:

- 987 • The answer [N/A] means that the abstract and introduction do not include the claims  
988 made in the paper.
- 989 • The abstract and/or introduction should clearly state the claims made, including the  
990 contributions made in the paper and important assumptions and limitations. A [No] or  
991 [N/A] answer to this question will not be perceived well by the reviewers.
- 992 • The claims made should match theoretical and experimental results, and reflect how  
993 much the results can be expected to generalize to other settings.
- 994 • It is fine to include aspirational goals as motivation as long as it is clear that these goals  
995 are not attained by the paper.

### 996 **2. Limitations**

997 Question: Does the paper discuss the limitations of the work performed by the authors?

998 Answer: [Yes]

999  
1000  
1001  
1002  
1003  
1004  
1005  
1006  
1007  
1008  
1009  
1010  
1011  
1012  
1013  
1014  
1015  
1016  
1017  
1018  
1019  
1020  
1021  
1022  
1023  
1024  
1025  
1026  
1027  
1028  
1029  
1030  
1031  
1032  
1033  
1034  
1035  
1036  
1037  
1038  
1039  
1040  
1041  
1042  
1043  
1044  
1045  
1046  
1047  
1048  
1049  
1050  
1051

Justification: The limitations are mentioned both in the main text and in the Appendix.

Guidelines:

- The answer [N/A] means that the paper has no limitation while the answer [No] means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate “Limitations” section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren’t acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

**3. Theory assumptions and proofs**

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [Yes]

Justification: For Theorem 1 and Counterexample1, we provide the full set of assumptions in Section 2.2, 2.3 and in Appendices B.2 and B.3.

Guidelines:

- The answer [N/A] means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

**4. Experimental result reproducibility**

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: The manuscript adequately describes all experimental details needed to reproduce the results. Code is submitted as supplementary material and will be made publicly available upon acceptance.

1052  
1053  
1054  
1055  
1056  
1057  
1058  
1059  
1060  
1061  
1062  
1063  
1064  
1065  
1066  
1067  
1068  
1069  
1070  
1071  
1072  
1073  
1074  
1075  
1076  
1077  
1078  
1079  
1080  
1081  
1082  
1083  
1084  
1085  
1086  
1087  
1088  
1089  
1090  
1091  
1092  
1093  
1094  
1095  
1096  
1097  
1098  
1099  
1100  
1101  
1102  
1103  
1104  
1105  
1106

Guidelines:

- The answer [N/A] means that the paper does not include experiments.
- If the paper includes experiments, a [No] answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
  - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
  - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
  - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
  - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

**5. Open access to data and code**

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: Code will be made publicly available upon acceptance. All datasets are open-source.

Guidelines:

- The answer [N/A] means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://neurips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so [No] is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://neurips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.

- 1107 • At submission time, to preserve anonymity, the authors should release anonymized  
1108 versions (if applicable).  
1109 • Providing as much information as possible in supplemental material (appended to the  
1110 paper) is recommended, but including URLs to data and code is permitted.

## 1111 6. Experimental setting/details

1112 Question: Does the paper specify all the training and test details (e.g., data splits, hyperpa-  
1113 rameters, how they were chosen, type of optimizer) necessary to understand the results?

1114 Answer: [Yes]

1115 Justification: All experimental details are mentioned and included in the main text (Section 5)  
1116 and in Appendix E.

1117 Guidelines:

- 1118 • The answer [N/A] means that the paper does not include experiments.
- 1119 • The experimental setting should be presented in the core of the paper to a level of detail  
1120 that is necessary to appreciate the results and make sense of them.
- 1121 • The full details can be provided either with the code, in appendix, or as supplemental  
1122 material.

## 1123 7. Experiment statistical significance

1124 Question: Does the paper report error bars suitably and correctly defined or other appropriate  
1125 information about the statistical significance of the experiments?

1126 Answer: [No]

1127 Justification: We have followed the experimental protocol and visual presentation of most  
1128 recent DG [10, 51] and MTL [11, 70] literature which do not present standard deviation.

1129 Guidelines:

- 1130 • The answer [N/A] means that the paper does not include experiments.
- 1131 • The authors should answer [Yes] if the results are accompanied by error bars, confidence  
1132 intervals, or statistical significance tests, at least for the experiments that support the  
1133 main claims of the paper.
- 1134 • The factors of variability that the error bars are capturing should be clearly stated (for  
1135 example, train/test split, initialization, random drawing of some parameter, or overall  
1136 run with given experimental conditions).
- 1137 • The method for calculating the error bars should be explained (closed form formula,  
1138 call to a library function, bootstrap, etc.)
- 1139 • The assumptions made should be given (e.g., Normally distributed errors).
- 1140 • It should be clear whether the error bar is the standard deviation or the standard error  
1141 of the mean.
- 1142 • It is OK to report 1-sigma error bars, but one should state it. The authors should  
1143 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis  
1144 of Normality of errors is not verified.
- 1145 • For asymmetric distributions, the authors should be careful not to show in tables or  
1146 figures symmetric error bars that would yield results that are out of range (e.g., negative  
1147 error rates).
- 1148 • If error bars are reported in tables or plots, the authors should explain in the text how  
1149 they were calculated and reference the corresponding figures or tables in the text.

## 1150 8. Experiments compute resources

1151 Question: For each experiment, does the paper provide sufficient information on the com-  
1152 puter resources (type of compute workers, memory, time of execution) needed to reproduce  
1153 the experiments?

1154 Answer: [Yes]

1155 Justification: All experimental details are mentioned and included in the main text (Section 5)  
1156 and in Appendix E.

1157 Guidelines:

- 1158 • The answer [N/A] means that the paper does not include experiments.
- 1159 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
- 1160 or cloud provider, including relevant memory and storage.
- 1161 • The paper should provide the amount of compute required for each of the individual
- 1162 experimental runs as well as estimate the total compute.
- 1163 • The paper should disclose whether the full research project required more compute
- 1164 than the experiments reported in the paper (e.g., preliminary or failed experiments that
- 1165 didn't make it into the paper).

## 1166 9. Code of ethics

1167 Question: Does the research conducted in the paper conform, in every respect, with the  
1168 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

1169 Answer: [Yes]

1170 Justification: The research conducted in the paper conforms, in every respect, with the  
1171 NeurIPS Code of Ethics.

1172 Guidelines:

- 1173 • The answer [N/A] means that the authors have not reviewed the NeurIPS Code of
- 1174 Ethics.
- 1175 • If the authors answer [No], they should explain the special circumstances that require a
- 1176 deviation from the Code of Ethics.
- 1177 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-
- 1178 eration due to laws or regulations in their jurisdiction).

## 1179 10. Broader impacts

1180 Question: Does the paper discuss both potential positive societal impacts and negative  
1181 societal impacts of the work performed?

1182 Answer: [N/A]

1183 Justification: [N/A]

1184 Guidelines:

- 1185 • The answer [N/A] means that there is no societal impact of the work performed.
- 1186 • If the authors answer [N/A] or [No], they should explain why their work has no societal
- 1187 impact or why the paper does not address societal impact.
- 1188 • Examples of negative societal impacts include potential malicious or unintended uses
- 1189 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations
- 1190 (e.g., deployment of technologies that could make decisions that unfairly impact specific
- 1191 groups), privacy considerations, and security considerations.
- 1192 • The conference expects that many papers will be foundational research and not tied
- 1193 to particular applications, let alone deployments. However, if there is a direct path to
- 1194 any negative applications, the authors should point it out. For example, it is legitimate
- 1195 to point out that an improvement in the quality of generative models could be used to
- 1196 generate Deepfakes for disinformation. On the other hand, it is not needed to point out
- 1197 that a generic algorithm for optimizing neural networks could enable people to train
- 1198 models that generate Deepfakes faster.
- 1199 • The authors should consider possible harms that could arise when the technology is
- 1200 being used as intended and functioning correctly, harms that could arise when the
- 1201 technology is being used as intended but gives incorrect results, and harms following
- 1202 from (intentional or unintentional) misuse of the technology.
- 1203 • If there are negative societal impacts, the authors could also discuss possible mitigation
- 1204 strategies (e.g., gated release of models, providing defenses in addition to attacks,
- 1205 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from
- 1206 feedback over time, improving the efficiency and accessibility of ML).

## 1207 11. Safeguards

1208 Question: Does the paper describe safeguards that have been put in place for responsible  
1209 release of data or models that have a high risk for misuse (e.g., pre-trained language models,  
1210 image generators, or scraped datasets)?

1211  
1212  
1213  
1214  
1215  
1216  
1217  
1218  
1219  
1220  
1221  
1222  
1223  
1224  
1225  
1226  
1227  
1228  
1229  
1230  
1231  
1232  
1233  
1234  
1235  
1236  
1237  
1238  
1239  
1240  
1241  
1242  
1243  
1244  
1245  
1246  
1247  
1248  
1249  
1250  
1251  
1252  
1253  
1254  
1255  
1256  
1257  
1258  
1259  
1260  
1261  
1262

Answer: [N/A]

Justification: [N/A]

Guidelines:

- The answer [N/A] means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

## 12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: All creators or original owners of assets (e.g., code, data, models), used in the paper, are properly credited. The license and terms of use are explicitly mentioned and properly respected.

Guidelines:

- The answer [N/A] means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, [paperswithcode.com/datasets](https://paperswithcode.com/datasets) has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

## 13. New assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [N/A]

Justification: Code will be released upon acceptance which will include appropriate documentation.

Guidelines:

- The answer [N/A] means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

1263  
1264  
1265  
1266  
1267  
1268  
1269  
1270  
1271  
1272  
1273  
1274  
1275  
1276  
1277  
1278  
1279  
1280  
1281  
1282  
1283  
1284  
1285  
1286  
1287  
1288  
1289  
1290  
1291  
1292  
1293  
1294  
1295  
1296  
1297  
1298  
1299  
1300  
1301  
1302  
1303  
1304  
1305  
1306  
1307  
1308  
1309

#### 14. Crowdsourcing and research with human subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [N/A]

Justification: [N/A]

Guidelines:

- The answer [N/A] means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

#### 15. Institutional review board (IRB) approvals or equivalent for research with human subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [N/A]

Justification: [N/A]

Guidelines:

- The answer [N/A] means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

#### 16. Declaration of LLM usage

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does *not* impact the core methodology, scientific rigor, or originality of the research, declaration is not required.

Answer: [N/A]

Justification: [N/A]

Guidelines:

- The answer [N/A] means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy in the NeurIPS handbook for what should or should not be described.