

SUPPLEMENTARY MATERIAL FOR “EXPLORING MODE CONNECTIVITY IN KRYLOV SUB- SPACE FOR DOMAIN GENERALIZATION”

This supplementary file is structured as follows: It begins by presenting the Billiard Optimization Algorithm (BOA) framework in Section A, detailing its two core algorithms for line search and reflection phases within Krylov subspaces to enhance domain generalization. The theoretical foundation is then established through a comprehensive subspace analysis (see Section B), which compares the approximation capabilities of Krylov and random subspaces for test gradients. This includes rigorous mathematical proofs (Theorems B.1, B.8, and Corollary B.9) and supporting lemmas on Gamma/Beta functions in Section B.1. The subsequent Section B.2 provides an empirical investigation into Krylov subspace performance across diverse datasets and architectures, with results visualized in figures. The document concludes with reproducibility details (in Section C), a note on LLM Utilization (in Section D) for language refinement, and a discussion of broader impacts (in Section E).

A SKETCH OF BILLIARD OPTIMIZATION ALGORITHM

This section presents a detailed sketch of the Billiard Optimization Algorithm (BOA), a novel optimization technique inspired by the dynamics of billiard motion. It consists of two main algorithms: Algorithm 1 outlines the core BOA framework, which starts with initialization using pre-trained parameters and a loss threshold, then iterates through a line search phase (leveraging Algorithm 2 to find solution intervals and golden-section search for parameter updates) and a reflection phase (computing normal vectors, projecting them into a Krylov subspace, and deriving reflected directions to navigate loss landscapes). Algorithm 2 supports this by efficiently bracketing intervals containing solutions to loss threshold equations. Overall, BOA aims to traverse high-dimensional parameter spaces via physics-inspired operations, enhancing domain generalization by maintaining low loss while exploring connected modes.

Algorithm 1: Billiard Optimization Algorithm

Input: Training dataset \mathcal{D}_{train} , reflection count N , hyperparameters (K, ϵ, Δ_ℓ , step size h).

Output: Parameter trajectory $\{\theta_i\}_{i=1}^N$.

1: **Initialization:**

2: Load pretrained parameters θ_0 .

3: Compute loss threshold $\ell_{th} \leftarrow \ell(\theta_0) + \Delta_\ell$.

4: Compute initial gradient $\mathbf{g}_0 \leftarrow \nabla \ell(\theta_0)$ on \mathcal{D}_{train} .

5: Initialize incident direction via the following equation:

$$\mathbf{p}_0 \leftarrow -\text{normalize}\left(\sum_{k=0}^{K-1} \mathbf{H}^k \mathbf{g}_0\right).$$

6: **for** $i = 1$ **to** N **do**

7: **Line Search Phase:**

8: Determine solution interval $[h_L, h_R]$ via Algorithm 2.

9: Get the approximate solution α^* of $\ell(\theta_{i-1} + \alpha \mathbf{p}_0) = \ell_{th}$ via golden-section search using the above interval.

10: Update the parameters: $\theta_i \leftarrow \theta_{i-1} + \alpha^* \mathbf{p}_{i-1}$.

11: **Reflection Phase:**

12: Compute normal vector: $\mathbf{n}_i \leftarrow -\nabla \ell(\theta_i) / \|\nabla \ell(\theta_i)\|$.

13: Project normal vector into the Krylov subspace: $\tilde{\mathbf{n}}_i \leftarrow \text{proj}_{\mathcal{K}} \mathbf{n}_i$.

14: Compute reflected direction: $\mathbf{p}_i \leftarrow (\mathbf{I} - 2\tilde{\mathbf{n}}_i \tilde{\mathbf{n}}_i^\top) \mathbf{p}_{i-1}$.

15: **end for**

16: **return** Parameter trajectory $\{\theta_i\}_{i=1}^N$.

Algorithm 2: Interval Search for Solution Bracket.

Input: Training dataset \mathcal{D}_{train} , current model parameter θ , the incident direction \mathbf{p} , loss threshold ℓ_{th} , step size h .

Output: Interval $[h_L, h_R]$ containing solution to $\ell(\theta + \alpha\mathbf{p}) = \ell_{th}$.

```

1: Initialize  $\mathbf{x} = \theta + h\mathbf{p}$ .
2: if  $\ell(\mathbf{x}) < \ell_{th}$  then
3:   Initialize  $h_L \leftarrow 0$ .
4:   while  $\ell(\mathbf{x}) < \ell_{th}$  do
5:      $h_L \leftarrow h_L + h; h \leftarrow 2h; \mathbf{x} \leftarrow \mathbf{x} + h\mathbf{p}$ .
6:   end while
7:    $h_R \leftarrow h_L + h$ .
8: else
9:   while  $\ell(\mathbf{x}) > \ell_{th}$  do
10:     $h \leftarrow h/2; \mathbf{x} \leftarrow \theta + h\mathbf{p}$ .
11:   end while
12:    $h_L \leftarrow h; h_R \leftarrow 2h$ .
13: end if
14: return The interval  $[h_L, h_R]$ .

```

B SUBSPACE ANALYSIS

As demonstrated in Table 1 presented in the main manuscript, prevalent optimization algorithms, such as SGD (SGLD Welling & Teh (2011)) and SAM Foret et al. (2020), can uniformly interpret gradient updates through either Krylov subspaces or random subspaces. In this section, we demonstrate the superiority of Krylov subspaces over random subspaces for approximating test gradients. Specifically, we examine the best possible approximation of test gradients achievable when using a set of basis vectors, which mathematically corresponds to projecting the test gradients onto the subspace spanned by those bases.

To analyze random subspaces, we consider the case where the subspace is constrained to include the training gradient, under the rationale that the training gradient may itself be a reasonable approximation to the test gradients. Under this assumption, we derive the expected length of the projection of a unit directional vector onto all such random subspaces that contain the training gradient. For the Krylov subspace, we empirically investigate how the projection length of the test gradients onto the Krylov subspace varies as the dimension of the subspace increases. This dynamic is illustrated through a plotted relationship, highlighting the characteristics of Krylov-based approximation. We also compare this projection length from the Krylov subspace with the expected projection length obtained from random subspaces of the same dimension, clearly showcasing the comparative advantage of Krylov subspaces in capturing the directional information of test gradients more effectively.

B.1 RANDOM SUBSPACE ANALYSIS

In this section, we focus on analyzing random subspaces that explicitly contain the training gradient vector \mathbf{a} . We begin by introducing Theorem B.1, which establishes that all such subspaces can be mathematically represented as a direct sum structure: specifically, $\text{span}(\mathbf{a}) \oplus W$, where W denotes an arbitrary subspace lying within the orthogonal complement of \mathbf{a} , i.e., $\text{span}(\mathbf{a})^\perp$. This decomposition provides a foundational framework for the subsequent understanding of these constrained random subspaces.

Subsequently, we present Theorem B.8, which offers a detailed statistical characterization (such as the expectation, variance, and concentration bounds) about the projection length of a fixed vector onto random subspaces. To ensure a comprehensive and rigorous derivation of Theorem B.8, we first introduce several essential definitions and preliminary results related to Gamma and Beta functions, along with supporting lemmas. These mathematical tools are critical for facilitating the proofs and subsequent analyses.

By synthesizing the structural insights from Theorem B.1 and the probabilistic findings from Theorem B.8, we derive a significant corollary (Corollary B.9). This corollary delivers an expectation bound for the projection of a fixed vector \mathbf{b} (interpreted as the test gradient) onto random subspaces that contain \mathbf{a} (interpreted as the training gradient). According to this corollary, we can obtain theo-

retical cosine increases between test gradients and their n -dimensional Krylov Subspace projections, which are present in Table 1. Moreover, Corollary B.9 demonstrates a sharp concentration phenomenon, indicating that the probability of the projection magnitude deviating significantly from its expected value is exceedingly low.

Theorem B.1. *Let V be a finite-dimensional vector space over a field \mathbb{F} endowed with an inner product (\cdot, \cdot) , and let $\mathbf{a} \in V$ be a fixed nonzero vector. Define the collection \mathcal{A} as the set of all subspaces of V that contain \mathbf{a} , that is, $\mathcal{A} = \{U \subseteq V \mid U \text{ is a subspace and } \mathbf{a} \in U\}$. Define the collection \mathcal{B} as the set of all subspaces expressible as the direct sum: $\text{span}(\mathbf{a}) \oplus W$, where W is any subspace of the orthogonal complement $\text{span}(\mathbf{a})^\perp$ (denoted by \mathbf{a}^\perp hereafter). Then $\mathcal{A} = \mathcal{B}$.*

Proof. To prove $\mathcal{A} = \mathcal{B}$, we demonstrate mutual set inclusion.

First, every subspace U in \mathcal{B} is of the form $\text{span}(\mathbf{a}) \oplus W$ for some $W \subseteq \mathbf{a}^\perp$. Since \mathbf{a} lies in $\text{span}(\mathbf{a})$, it is contained in U , confirming that U belongs to \mathcal{A} . Thus, \mathcal{B} is a subset of \mathcal{A} .

Conversely, take any subspace U in \mathcal{A} , meaning $\mathbf{a} \in U$. Define $W = U \cap \mathbf{a}^\perp$, which is a subspace of \mathbf{a}^\perp because it is the intersection of two subspaces. We claim that U decomposes orthogonally as $\text{span}(\mathbf{a}) \oplus W$. To verify this, observe that for any vector $\mathbf{u} \in U$, it can be orthogonally decomposed as $\mathbf{u} = \frac{(\mathbf{u}, \mathbf{a})}{(\mathbf{a}, \mathbf{a})} \mathbf{a} + \left(\mathbf{u} - \frac{(\mathbf{u}, \mathbf{a})}{(\mathbf{a}, \mathbf{a})} \mathbf{a} \right)$. The first term is a scalar multiple of \mathbf{a} and hence belongs to $\text{span}(\mathbf{a})$. The second term is orthogonal to \mathbf{a} by construction, and since both \mathbf{u} and \mathbf{a} are in U , the second term also lies in U ; therefore, it belongs to W . Moreover, the intersection of $\text{span}(\mathbf{a})$ and W is trivial: if a vector $c\mathbf{a}$ (for some scalar c) is in W , then it must be orthogonal to \mathbf{a} , implying $c = 0$ since $\mathbf{a} \neq \mathbf{0}$. This shows $U \subseteq \text{span}(\mathbf{a}) \oplus W$. The reverse inclusion is immediate because both $\text{span}(\mathbf{a})$ and W are contained in U . Hence, U is exactly the direct sum $\text{span}(\mathbf{a}) \oplus W$, and so U is an element of \mathcal{B} . This proves $\mathcal{A} \subseteq \mathcal{B}$.

By mutual inclusion, we conclude that $\mathcal{A} = \mathcal{B}$. \square

Definition B.2 (Gamma Function). *The Gamma function, denoted as $\Gamma(z)$, generalizes the factorial function to complex and real numbers (excluding non-positive integers). For $\text{Re}(z) > 0$, it is defined by the improper integral:*

$$\Gamma(z) = \int_0^\infty t^{z-1} e^{-t} dt. \quad (1)$$

This is also called Euler's integral of the second kind. Key properties of the Gamma Function include:

- (1) *Recurrence relation:* $\Gamma(z + 1) = z\Gamma(z)$.
- (2) *For positive integers n ,* $\Gamma(n) = (n - 1)!$.

Lemma B.3 (Stirling's Approximation for the Gamma Function). *Let $\Gamma(z)$ denote the gamma function defined for $\Re(z) > 0$ by the integral $\Gamma(z) = \int_0^\infty t^{z-1} e^{-t} dt$. As $|z| \rightarrow \infty$ in the complex plane with $|\arg z| < \pi$ (i.e., excluding the negative real axis), $\Gamma(z)$ satisfies the asymptotic expansion:*

$$\Gamma(z) \sim \sqrt{\frac{2\pi}{z}} \left(\frac{z}{e}\right)^z \left(1 + \frac{1}{12z} + \frac{1}{288z^2} + \dots\right). \quad (2)$$

where the series arises from higher-order corrections in Laplace's method applied to the integral representation of $\Gamma(z)$. The leading term $\sqrt{\frac{2\pi}{z}} \left(\frac{z}{e}\right)^z$ constitutes the basic Stirling approximation, while the subsequent series forms a divergent but asymptotic expansion that converges factorially fast for sufficiently large $|z|$.

Definition B.4 (Beta Function). *The Beta function, denoted as $B(x, y)$, is a special function defined for positive real numbers $x > 0$ and $y > 0$ by the integral:*

$$B(x, y) = \int_0^1 t^{x-1} (1-t)^{y-1} dt. \quad (3)$$

This integral converges for $x > 0$ and $y > 0$, and it is also known as Euler's integral of the first kind. It relates to the Gamma function via the following equation:

$$B(x, y) = \frac{\Gamma(x)\Gamma(y)}{\Gamma(x+y)}. \quad (4)$$

It exhibits symmetry (i.e., $B(x, y) = B(y, x)$), and is fundamental in probability theory, particularly for normalizing the Beta distribution.

Definition B.5 (Beta Distribution). *The probability density function (PDF) of the Beta distribution is defined for $x \in [0, 1]$ as:*

$$f_X(x) = \frac{1}{B(\alpha, \beta)} x^{\alpha-1} (1-x)^{\beta-1}, \quad (5)$$

where $B(\alpha, \beta)$ represents the Beta function.

Lemma B.6. *Let X be a random variable following a Beta distribution with parameters $\alpha > 0$ and $\beta > 0$, denoted as $X \sim \text{Beta}(\alpha, \beta)$. The expectation and variance of \sqrt{X} are given by:*

$$\mathbb{E}[\sqrt{X}] = \frac{\Gamma(\alpha + \frac{1}{2})\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\alpha + \beta + \frac{1}{2})}, \quad (6)$$

$$\mathbb{V}[\sqrt{X}] = \frac{\alpha}{\alpha + \beta} - \left(\frac{\Gamma(\alpha + \frac{1}{2})\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\alpha + \beta + \frac{1}{2})} \right)^2. \quad (7)$$

Proof. We derive each component step by step in the following.

1. Compute the expectation.

The expectation $\mathbb{E}[\sqrt{X}]$ is computed by integrating \sqrt{x} against the probability density function of the Beta distribution:

$$\mathbb{E}[\sqrt{X}] = \int_0^1 \sqrt{x} f_X(x) dx = \int_0^1 x^{1/2} \cdot \frac{1}{B(\alpha, \beta)} x^{\alpha-1} (1-x)^{\beta-1} dx. \quad (8)$$

We simplify the integrand to obtain:

$$\mathbb{E}[\sqrt{X}] = \frac{1}{B(\alpha, \beta)} \int_0^1 x^{\alpha+\frac{1}{2}-1} (1-x)^{\beta-1} dx. \quad (9)$$

The integral $\int_0^1 x^{\alpha+\frac{1}{2}-1} (1-x)^{\beta-1} dx$ corresponds to the Beta function $B(\alpha + \frac{1}{2}, \beta)$. Therefore, we derive:

$$\mathbb{E}[\sqrt{X}] = \frac{B(\alpha + \frac{1}{2}, \beta)}{B(\alpha, \beta)}. \quad (10)$$

We express both Beta functions in terms of Gamma functions as follows:

$$B\left(\alpha + \frac{1}{2}, \beta\right) = \frac{\Gamma(\alpha + \frac{1}{2})\Gamma(\beta)}{\Gamma(\alpha + \frac{1}{2} + \beta)}, \quad B(\alpha, \beta) = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha + \beta)}. \quad (11)$$

Substituting these expressions into the equation yields:

$$\mathbb{E}[\sqrt{X}] = \frac{\frac{\Gamma(\alpha + \frac{1}{2})\Gamma(\beta)}{\Gamma(\alpha + \frac{1}{2} + \beta)}}{\frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha + \beta)}} = \frac{\Gamma(\alpha + \frac{1}{2})\Gamma(\beta)}{\Gamma(\alpha + \frac{1}{2} + \beta)} \cdot \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)}. \quad (12)$$

The $\Gamma(\beta)$ terms cancel, which leads to the final result:

$$\mathbb{E}[\sqrt{X}] = \frac{\Gamma(\alpha + \frac{1}{2})\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\alpha + \beta + \frac{1}{2})}. \quad (13)$$

2. Compute the variance.

Define $Y = \sqrt{X}$. The variance of Y is:

$$\mathbb{V}[Y] = \mathbb{E}[Y^2] - (\mathbb{E}[Y])^2 = \mathbb{E}[X] - \left(\mathbb{E}[\sqrt{X}]\right)^2. \quad (14)$$

The expectation of X is:

$$\mathbb{E}[X] = \int_0^1 x \cdot f_X(x) dx = \frac{1}{B(\alpha, \beta)} \int_0^1 x^\alpha (1-x)^{\beta-1} dx. \quad (15)$$

The integral is the Beta function $B(\alpha + 1, \beta)$:

$$\mathbb{E}[X] = \frac{B(\alpha + 1, \beta)}{B(\alpha, \beta)} = \frac{\Gamma(\alpha + 1)\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\alpha + \beta + 1)}. \quad (16)$$

Using $\Gamma(z + 1) = z\Gamma(z)$, this simplifies to:

$$\mathbb{E}[X] = \frac{\alpha}{\alpha + \beta}. \quad (17)$$

Substitute equation (13) and (17) into the variance formula (14), we obtain:

$$\mathbb{V}[\sqrt{X}] = \frac{\alpha}{\alpha + \beta} - \left(\frac{\Gamma(\alpha + \frac{1}{2})\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\alpha + \beta + \frac{1}{2})} \right)^2. \quad (18)$$

This is the exact closed-form expression. \square

Lemma B.7 (Concentration of Lipschitz functions on the sphere Vershynin (2018)). *Let X be a random vector uniformly distributed on the Euclidean sphere of radius \sqrt{n} , i.e., $X \sim \text{Unif}(\sqrt{n}S^{n-1})$. Then for any Lipschitz function¹ $f : \sqrt{n}S^{n-1} \rightarrow \mathbb{R}$ and any $t \geq 0$, we have*

$$\mathbb{P}\{|f(X) - \mathbb{E}f(X)| \geq t\} \leq 2 \exp\left(-\frac{ct^2}{\|f\|_{\text{Lip}}^2}\right). \quad (19)$$

Theorem B.8. *Let \mathbf{a} be a fixed unit vector in \mathbb{R}^N , and let V be an n -dimensional subspace of \mathbb{R}^N sampled uniformly from the Grassmannian manifold $G(n, N)$ of all n -dimensional subspaces. Let $\mathbf{P}_V(\mathbf{a})$ denote the orthogonal projection of \mathbf{a} onto V . Then, we have the following three conclusions: (1) The cosine of the angle θ between the original vector \mathbf{a} and its projection $\mathbf{P}_V(\mathbf{a})$ satisfies the relation $\cos \theta = \|\mathbf{P}_V(\mathbf{a})\|$, and the expected value of this cosine scales with the subspace dimension n and the ambient dimension N according to:*

$$\mathbb{E}[\cos \theta] = \mathbb{E}[\|\mathbf{P}_V(\mathbf{a})\|] = \frac{\Gamma(\frac{n}{2} + \frac{1}{2})\Gamma(\frac{N}{2})}{\Gamma(\frac{n}{2})\Gamma(\frac{N}{2} + \frac{1}{2})}, \quad (20)$$

where Γ denotes the gamma function. This expectation increases monotonically with n , approaching 1 as $n \rightarrow N$. For large N and n , $\mathbb{E}[\cos \theta] \approx \sqrt{\frac{n}{N}}$.

(2) The variance of the cosine is given by:

$$\mathbb{V}[\cos \theta] = \mathbb{V}[\|\mathbf{P}_V(\mathbf{a})\|] = \frac{n}{N} - \left(\frac{\Gamma(\frac{n+1}{2})\Gamma(\frac{N}{2})}{\Gamma(\frac{n}{2})\Gamma(\frac{N+1}{2})} \right)^2, \quad (21)$$

which exhibits a decay rate of $\mathcal{O}(1/N)$ as N increases, for any fixed dimension n .

(3) $\cos^2 \theta$ is $\mathcal{O}(1/N)$ -subgaussian. That is, for all $t > 0$, the tail probability is bounded by:

$$\mathbb{P}\left(\left|\cos^2 \theta - \frac{n}{N}\right| \geq t\right) \leq 2 \exp\left(-\frac{t^2}{\mathcal{O}(1/N)}\right). \quad (22)$$

Proof. The orthogonal projection $\mathbf{P}_V(\mathbf{a})$ resides within the subspace V , and by the properties of orthogonal projection in inner product spaces, the cosine of the angle θ between \mathbf{a} and its projection is given by $\cos \theta = \frac{(\mathbf{a}, \mathbf{P}_V(\mathbf{a}))}{\|\mathbf{a}\|\|\mathbf{P}_V(\mathbf{a})\|}$. Since \mathbf{a} is a unit vector and because $(\mathbf{a}, \mathbf{P}_V(\mathbf{a})) = \|\mathbf{P}_V(\mathbf{a})\|^2$ (a fundamental property of orthogonal projections), this expression simplifies to $\cos \theta = \|\mathbf{P}_V(\mathbf{a})\|$. The squared norm $\|\mathbf{P}_V(\mathbf{a})\|^2$ represents the fraction of the vector's energy preserved in the projection.

¹A function $f : \mathcal{X} \rightarrow \mathbb{R}$ is Lipschitz continuous if there exists a constant $\|f\|_{\text{Lip}} \geq 0$ (called the Lipschitz constant) such that for all $x, y \in \mathcal{X}$, $|f(x) - f(y)| \leq \|f\|_{\text{Lip}} \cdot \rho(x, y)$, where $\rho(x, y)$ is a distance metric on the domain \mathcal{X} (e.g., the Euclidean distance in \mathbb{R}^n).

Due to the uniform random sampling of V from $G(n, N)$, the quantity $\|\mathbf{P}_V(\mathbf{a})\|^2$ is a random variable that follows a Beta distribution: $\|\mathbf{P}_V(\mathbf{a})\|^2 \sim \text{Beta}\left(\frac{n}{2}, \frac{N-n}{2}\right)$.

To determine the distribution of $\|\mathbf{P}_V(\mathbf{a})\|^2$, we exploit the rotational symmetry of the uniform distribution on $G(n, N)$. Specifically, the uniform distribution on $G(n, N)$ is invariant under the action of the orthogonal group $O(N)$. Thus, without loss of generality, we may fix $\mathbf{a} = \mathbf{e}_1 = (1, 0, \dots, 0)^\top$, as the distribution of $\|\mathbf{P}_V(\mathbf{a})\|^2$ depends only on the geometry of the subspace and not on the specific direction of \mathbf{a} . Note that a random subspace $V \in G(n, N)$ can be generated by selecting a random orthogonal matrix $Q \in O(N)$ uniformly and defining V as the span of the first n columns of Q , denoted $\mathbf{q}_1, \dots, \mathbf{q}_n$. The squared projection length is then:

$$\|\mathbf{P}_V(\mathbf{a})\|^2 = \sum_{i=1}^n (\mathbf{a}^\top \mathbf{q}_i)^2 = \sum_{i=1}^n \mathbf{e}_1^\top \mathbf{q}_i \mathbf{q}_i^\top \mathbf{e}_1 = \sum_{i=1}^n q_{i1}^2, \quad (23)$$

where q_{i1} is the first component of \mathbf{q}_i . That is, $\|\mathbf{P}_V(\mathbf{a})\|^2 = \sum_{i=1}^n w_i^2$, where $\mathbf{w} = (w_1, \dots, w_N)^\top = Q^\top \mathbf{e}_1$ is a random vector uniformly distributed on the unit sphere S^{N-1} in \mathbb{R}^N . It is well-known that a random vector \mathbf{w} uniformly distributed on S^{N-1} can be generated by normalizing a standard Gaussian vector. Let X_1, \dots, X_N be independent and identically distributed (i.i.d.) standard normal random variables, $X_i \sim \mathcal{N}(0, 1)$. Then:

$$\mathbf{w} = \left(\frac{X_1}{R}, \frac{X_2}{R}, \dots, \frac{X_N}{R} \right)^\top, \quad \text{where } R = \sqrt{\sum_{j=1}^N X_j^2}. \quad (24)$$

Thus, the squared projection becomes:

$$\|\mathbf{P}_V(\mathbf{a})\|^2 = \sum_{i=1}^n w_i^2 = \frac{\sum_{i=1}^n X_i^2}{\sum_{j=1}^N X_j^2}. \quad (25)$$

Define $S_n = \sum_{i=1}^n X_i^2$ and $S_{N-n} = \sum_{k=n+1}^N X_k^2$. Since the X_i are i.i.d. $\mathcal{N}(0, 1)$, $S_n \sim \chi^2(n)$ and $S_{N-n} \sim \chi^2(N-n)$ are independent chi-square random variables with n and $N-n$ degrees of freedom, respectively. Consequently:

$$\|\mathbf{P}_V(\mathbf{a})\|^2 = \frac{S_n}{S_n + S_{N-n}}. \quad (26)$$

The ratio $\frac{S_n}{S_n + S_{N-n}}$ is a Beta-distributed random variable. Specifically, if $A \sim \chi^2(a)$ and $B \sim \chi^2(b)$ are independent, then $\frac{A}{A+B} \sim \text{Beta}\left(\frac{a}{2}, \frac{b}{2}\right)$. Here, $a = n$ and $b = N-n$, so:

$$\frac{S_n}{S_n + S_{N-n}} \sim \text{Beta}\left(\frac{n}{2}, \frac{N-n}{2}\right). \quad (27)$$

Therefore, $\cos^2 \theta = \|\mathbf{P}_V(\mathbf{a})\|^2$ follows a Beta distribution with parameters $\alpha = \frac{n}{2}$ and $\beta = \frac{N-n}{2}$.

Proof of (1). The expected value of the cosine, $\mathbb{E}[\|\mathbf{P}_V(\mathbf{a})\|]$, is the expectation of the square root of this Beta-distributed variable. For a random variable $X \sim \text{Beta}(\alpha, \beta)$, the expectation $\mathbb{E}[\sqrt{X}]$ is given by $\frac{\Gamma(\alpha + \frac{1}{2})\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\alpha + \beta + \frac{1}{2})}$, as depicted in Lemma B.6. Substituting $\alpha = \frac{n}{2}$ and $\beta = \frac{N-n}{2}$ yields the stated result:

$$\mathbb{E}[\|\mathbf{P}_V(\mathbf{a})\|] = \frac{\Gamma\left(\frac{n}{2} + \frac{1}{2}\right)\Gamma\left(\frac{N}{2}\right)}{\Gamma\left(\frac{n}{2}\right)\Gamma\left(\frac{N}{2} + \frac{1}{2}\right)}. \quad (28)$$

The monotonic increase of $\mathbb{E}[\cos \theta]$ with n follows intuitively because a higher-dimensional random subspace has a greater probability of capturing a significant component of any fixed vector. In the limit as n approaches N , the subspace V encompasses the entire ambient space, the projection becomes the identity operation ($\|\mathbf{P}_V(\mathbf{a})\| = 1$), and thus $\mathbb{E}[\cos \theta] \rightarrow 1$. To analyze how $\mathbb{E}[\cos \theta]$ varies with n and N , we next apply Stirling's approximation. Stirling's approximation states that for large z , $\Gamma(z) \sim \sqrt{\frac{2\pi}{z}} \left(\frac{z}{e}\right)^z$. Applying this to each Gamma function in equation (28) yields:

$$\mathbb{E}[\|\mathbf{P}_V(\mathbf{a})\|] \approx \frac{\sqrt{\frac{2\pi}{\frac{n}{2} + \frac{1}{2}}} \left(\frac{\frac{n}{2} + \frac{1}{2}}{e}\right)^{\frac{n}{2} + \frac{1}{2}} \cdot \sqrt{\frac{2\pi}{\frac{N}{2}}} \left(\frac{N}{2}\right)^{\frac{N}{2}}}{\sqrt{\frac{2\pi}{\frac{n}{2}}} \left(\frac{n}{2}\right)^{\frac{n}{2}} \cdot \sqrt{\frac{2\pi}{\frac{N}{2} + \frac{1}{2}}} \left(\frac{N}{2} + \frac{1}{2}\right)^{\frac{N}{2} + \frac{1}{2}}}. \quad (29)$$

Simplifying radicals and exponents, this reduces to:

$$\mathbb{E}[\|\mathbf{P}_V(\mathbf{a})\|] \approx \frac{\frac{1}{\sqrt{(n+1)N}}(n+1)^{(n+1)/2}N^{N/2}}{\frac{1}{\sqrt{n(N+1)}}n^{n/2}(N+1)^{(N+1)/2}}. \quad (30)$$

The expression is then algebraically rearranged to isolate dominant terms:

$$\mathbb{E}[\|\mathbf{P}_V(\mathbf{a})\|] \approx \frac{\sqrt{n(N+1)}}{\sqrt{(n+1)N}} \cdot \frac{\sqrt{(n+1) \cdot (1 + \frac{1}{n})^n}}{\sqrt{(N+1) \cdot (1 + \frac{1}{N})^N}}. \quad (31)$$

This rearrangement is exact and leverages properties of exponents and radicals to highlight the asymptotic behavior of the ratio. As $n, N \rightarrow \infty$, the terms $(1 + \frac{1}{n})^n \rightarrow e$, $(1 + \frac{1}{N})^N \rightarrow e$. Thus, we have:

$$\mathbb{E}[\|\mathbf{P}_V(\mathbf{a})\|] \approx \sqrt{\frac{n}{N}}. \quad (32)$$

The final approximation $\sqrt{n/N}$ holds for large n and N .

Proof of (2). To derive the asymptotic behavior of the variance for the norm of the projection $\|\mathbf{P}_V(\mathbf{a})\|$, we begin with the exact variance expression obtained from the Beta distribution properties:

$$\mathbb{V}[\|\mathbf{P}_V(\mathbf{a})\|] = \frac{n}{N} - \left(\frac{\Gamma(\frac{n+1}{2})\Gamma(\frac{N}{2})}{\Gamma(\frac{n}{2})\Gamma(\frac{N+1}{2})} \right)^2. \quad (33)$$

This expression arises because $\|\mathbf{P}_V(\mathbf{a})\|^2 \sim \text{Beta}(\frac{n}{2}, \frac{N-n}{2})$, and the variance of $\|\mathbf{P}_V(\mathbf{a})\|$ is derived from the moments of this distribution. The term $\frac{n}{N}$ corresponds to $\mathbb{E}[\|\mathbf{P}_V(\mathbf{a})\|^2]$, while the subtracted term is the square of $\mathbb{E}[\|\mathbf{P}_V(\mathbf{a})\|]$.

To analyze the behavior for large N , we apply Stirling's approximation to the Gamma functions. Using Stirling's approximation:

$$\Gamma\left(\frac{N}{2}\right) \approx \sqrt{\frac{4\pi}{N}} \left(\frac{N}{2e}\right)^{N/2}, \quad (34)$$

$$\Gamma\left(\frac{N+1}{2}\right) \approx \sqrt{\frac{4\pi}{N+1}} \left(\frac{N+1}{2e}\right)^{(N+1)/2}. \quad (35)$$

Thus, their ratio simplifies to:

$$\frac{\Gamma\left(\frac{N}{2}\right)}{\Gamma\left(\frac{N+1}{2}\right)} \approx \frac{\sqrt{\frac{4\pi}{N}} \left(\frac{N}{2e}\right)^{N/2}}{\sqrt{\frac{4\pi}{N+1}} \left(\frac{N+1}{2e}\right)^{(N+1)/2}} = \sqrt{\frac{N+1}{N}} \cdot \frac{\left(\frac{N}{2e}\right)^{N/2}}{\left(\frac{N+1}{2e}\right)^{(N+1)/2}}. \quad (36)$$

Simplifying the exponential terms:

$$\frac{\left(\frac{N}{2e}\right)^{N/2}}{\left(\frac{N+1}{2e}\right)^{(N+1)/2}} = \left(\frac{N}{N+1}\right)^{N/2} \cdot \frac{1}{\sqrt{\frac{N+1}{2e}}}. \quad (37)$$

As $N \rightarrow \infty$, $(1 + \frac{1}{N})^{-N/2} \rightarrow e^{-1/2}$. Combining all factors:

$$\frac{\Gamma\left(\frac{N}{2}\right)}{\Gamma\left(\frac{N+1}{2}\right)} \approx \sqrt{\frac{N+1}{N}} \cdot \frac{e^{-1/2}}{\sqrt{\frac{N+1}{2e}}} = \sqrt{\frac{2}{N}}. \quad (38)$$

Then, the variance formula contains the squared ratio:

$$\left(\frac{\Gamma\left(\frac{n+1}{2}\right)\Gamma\left(\frac{N}{2}\right)}{\Gamma\left(\frac{n}{2}\right)\Gamma\left(\frac{N+1}{2}\right)} \right)^2 \approx \frac{2}{N} \left(\frac{\Gamma\left(\frac{n+1}{2}\right)}{\Gamma\left(\frac{n}{2}\right)} \right)^2. \quad (39)$$

Thus, the variance becomes:

$$\mathbb{V}[\|\mathbf{P}_V(\mathbf{a})\|] \approx \frac{n}{N} - \frac{2}{N} \left(\frac{\Gamma(\frac{n+1}{2})}{\Gamma(\frac{n}{2})} \right)^2. \quad (40)$$

The term $\left(\frac{\Gamma(\frac{n+1}{2})}{\Gamma(\frac{n}{2})} \right)^2$ depends only on the fixed dimension n , not N . Regardless of n , this term is bounded by a constant C_n . Thus,

$$\mathbb{V}[\|\mathbf{P}_V(\mathbf{a})\|] \approx \frac{n}{N} - \frac{C_n}{N} = \mathcal{O}\left(\frac{1}{N}\right). \quad (41)$$

Proof of (3). We will prove this conclusion by using Lemma B.7. Now we start with a Lipschitz function.

Consider the function $f : \mathbb{R}^N \rightarrow \mathbb{R}$ defined by

$$f(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^n x_i^2, \quad (42)$$

where $\mathbf{x} = (x_1, \dots, x_N)$ and $n \leq N$. The domain of f is the sphere of radius \sqrt{N} in \mathbb{R}^N , denoted $\sqrt{N}S^{N-1}$. This sphere consists of all vectors $\mathbf{x} \in \mathbb{R}^N$ with Euclidean norm $\|\mathbf{x}\| = \sqrt{N}$, i.e.,

$$\sqrt{N}S^{N-1} = \{\mathbf{x} \in \mathbb{R}^N \mid \|\mathbf{x}\| = \sqrt{N}\}. \quad (43)$$

This sphere is a Riemannian manifold with intrinsic dimension $N - 1$, and its geometry plays a crucial role in the analysis of concentration phenomena. The function f arises naturally in the context of orthogonal projections: if $\mathbf{w} = \mathbf{x}/\sqrt{N}$ is a unit vector on S^{N-1} , then $f(\mathbf{x}) = \sum_{i=1}^n w_i^2 = \|\mathbf{P}_V(\mathbf{a})\|^2$, where $\mathbf{P}_V(\mathbf{a})$ is the orthogonal projection of a fixed unit vector \mathbf{a} onto a random n -dimensional subspace V of \mathbb{R}^N . To establish the Lipschitz property of the function f , let $\mathbf{a}, \mathbf{b} \in \sqrt{N}S^{N-1}$ be arbitrary points on the sphere. The difference $|f(\mathbf{a}) - f(\mathbf{b})|$ is given by:

$$|f(\mathbf{a}) - f(\mathbf{b})| = \left| \frac{1}{N} \sum_{i=1}^n (a_i^2 - b_i^2) \right|. \quad (44)$$

Using the identity $a_i^2 - b_i^2 = (a_i - b_i)(a_i + b_i)$ and the triangle inequality, we have:

$$\left| \sum_{i=1}^n (a_i^2 - b_i^2) \right| \leq \sum_{i=1}^n |a_i - b_i| \cdot |a_i + b_i|. \quad (45)$$

By the Cauchy-Schwarz inequality, we obtain:

$$\sum_{i=1}^n |a_i - b_i| \cdot |a_i + b_i| \leq \sqrt{\sum_{i=1}^n (a_i - b_i)^2} \cdot \sqrt{\sum_{i=1}^n (a_i + b_i)^2}. \quad (46)$$

Since $\|\mathbf{a}\|^2 = \|\mathbf{b}\|^2 = N$, we have $\sum_{i=1}^N a_i^2 = N$ and $\sum_{i=1}^N b_i^2 = N$. Thus,

$$\sum_{i=1}^n (a_i + b_i)^2 \leq \sum_{i=1}^N (a_i + b_i)^2 = \|\mathbf{a}\|^2 + \|\mathbf{b}\|^2 + 2\mathbf{a} \cdot \mathbf{b} \leq 2N + 2|\mathbf{a} \cdot \mathbf{b}| \leq 4N, \quad (47)$$

where the last step follows from $|\mathbf{a} \cdot \mathbf{b}| \leq \|\mathbf{a}\|\|\mathbf{b}\| = N$. Combining these, we can get:

$$|f(\mathbf{a}) - f(\mathbf{b})| \leq \frac{1}{N} \cdot \sqrt{\sum_{i=1}^n (a_i - b_i)^2} \cdot \sqrt{4N} = \frac{2}{\sqrt{N}} \cdot \sqrt{\sum_{i=1}^n (a_i - b_i)^2}. \quad (48)$$

Finally, $\sqrt{\sum_{i=1}^n (a_i - b_i)^2} \leq \|\mathbf{a} - \mathbf{b}\|$, so

$$|f(\mathbf{a}) - f(\mathbf{b})| \leq \frac{2}{\sqrt{N}} \|\mathbf{a} - \mathbf{b}\|. \quad (49)$$

This shows that f is Lipschitz continuous on $\sqrt{N}S^{N-1}$ with Lipschitz constant $L = \frac{2}{\sqrt{N}}$.

Then, according to Lemma B.7, any Lipschitz function concentrates sharply around its mean. Specifically, if f is L -Lipschitz, then for a random vector \mathbf{x} uniformly distributed on $\sqrt{N}S^{N-1}$, we have

$$\mathbb{P}(|f(\mathbf{x}) - \mathbb{E}[f(\mathbf{x})]| \geq t) \leq 2 \exp\left(-\frac{ct^2}{L^2}\right), \quad \forall t > 0, \quad (50)$$

where $c > 0$ is an absolute constant. Substituting $L = \frac{2}{\sqrt{N}}$ yields

$$\mathbb{P}(|f(\mathbf{x}) - \mathbb{E}[f(\mathbf{x})]| \geq t) \leq 2 \exp\left(-\frac{ct^2}{(2/\sqrt{N})^2}\right) = 2 \exp\left(-\frac{cNt^2}{4}\right). \quad (51)$$

This bound is valid for all $t \geq 0$, demonstrating that $f(\mathbf{x})$ is subgaussian with variance proxy $\sigma^2 = \mathcal{O}(1/N)$. □

Corollary B.9. *Let V be an N -dimensional inner product space over \mathbb{R} , and let $\mathbf{a}, \mathbf{b} \in V$ be fixed unit vectors. Define the set of admissible subspaces:*

$$\mathcal{A} := \{U \subseteq V \mid \dim(U) = n, \mathbf{a} \in U\}, \quad 1 < n < N. \quad (52)$$

Let U be a random subspace sampled uniformly from \mathcal{A} . Denote by $\mathbf{P}_U(\mathbf{b})$ the orthogonal projection of \mathbf{b} onto U , and let θ be the angle between \mathbf{b} and $\mathbf{P}_U(\mathbf{b})$. Then, we have:

(1) *Subgaussian concentration:*

$$\mathbb{P}(|\cos^2 \theta - \mathbb{E}[\cos^2 \theta]| \geq t) \leq 2 \exp\left(-\frac{t^2}{\mathcal{O}(1/N)}\right), \quad \forall t > 0. \quad (53)$$

(2) *Expectation bound:*

$$\mathbb{E}[\cos \theta] = \mathbb{E}[\|\mathbf{P}_U(\mathbf{b})\|] \leq |(\mathbf{a}, \mathbf{b})| + \frac{\Gamma\left(\frac{n}{2}\right) \Gamma\left(\frac{N-1}{2}\right)}{\Gamma\left(\frac{n-1}{2}\right) \Gamma\left(\frac{N}{2}\right)}. \quad (54)$$

Proof. Since $\mathbf{a} \in U$, we decompose \mathbf{b} orthogonally relative to \mathbf{a} :

$$\mathbf{b} = (\mathbf{a}, \mathbf{b})\mathbf{a} + \mathbf{b}_\perp, \quad \mathbf{b}_\perp \in \mathbf{a}^\perp. \quad (55)$$

According to Theorem B.1, the projection onto U splits into components parallel and orthogonal to \mathbf{a} :

$$\|\mathbf{P}_U(\mathbf{b})\|^2 = (\mathbf{a}, \mathbf{b})^2 + \|\mathbf{P}_{U \cap \mathbf{a}^\perp}(\mathbf{b}_\perp)\|^2, \quad (56)$$

where $U \cap \mathbf{a}^\perp$ is a uniformly random $(n-1)$ -dimensional subspace in the $(N-1)$ -dimensional space \mathbf{a}^\perp . The term $\|\mathbf{P}_{U \cap \mathbf{a}^\perp}(\mathbf{b}_\perp)\|^2$ is the squared norm of the projection of \mathbf{b}_\perp onto a random $(n-1)$ -dimensional subspace of \mathbf{a}^\perp . By the rotational invariance of the uniform distribution on the Grassmannian $G(n-1, N-1)$, this follows a Beta distribution:

$$\|\mathbf{P}_{U \cap \mathbf{a}^\perp}(\mathbf{b}_\perp)\|^2 / \|\mathbf{b}_\perp\|^2 \sim \text{Beta}\left(\frac{n-1}{2}, \frac{N-n}{2}\right). \quad (57)$$

According to the Theorem B.8, we have:

$$\mathbb{P}(|\|\mathbf{P}_{U \cap \mathbf{a}^\perp}(\mathbf{b}_\perp)\|^2 - \mathbb{E}[\|\mathbf{P}_{U \cap \mathbf{a}^\perp}(\mathbf{b}_\perp)\|^2]| \geq t) \leq 2 \exp(-c(N-1)t^2). \quad (58)$$

Naturally,

$$\mathbb{P}(|\|\mathbf{P}_U(\mathbf{b})\|^2 - \mathbb{E}[\|\mathbf{P}_U(\mathbf{b})\|^2]| \geq t) \leq 2 \exp(-c(N-1)t^2). \quad (59)$$

The expectation $\mathbb{E}[\|\mathbf{P}_U(\mathbf{b})\|]$ is bounded using Jensen's inequality and the moments of the Beta distribution:

$$\mathbb{E}[\|\mathbf{P}_U(\mathbf{b})\|] \leq \sqrt{\mathbb{E}[\|\mathbf{P}_U(\mathbf{b})\|^2]} = \sqrt{(\mathbf{a}, \mathbf{b})^2 + \mathbb{E}[\|\mathbf{P}_{U \cap \mathbf{a}^\perp}(\mathbf{b}_\perp)\|^2]}. \quad (60)$$

Table 1: **Theoretical Cosine Increase Between Test Gradients and Their n -dimensional Krylov Subspace Projections.**

Dataset	ViT-B/32 & ViT-B/16			ViT-L/14		
	n	N	$\frac{\Gamma(\frac{n}{2})\Gamma(\frac{N-1}{2})}{\Gamma(\frac{n-1}{2})\Gamma(\frac{N}{2})}$	n	N	$\frac{\Gamma(\frac{n}{2})\Gamma(\frac{N-1}{2})}{\Gamma(\frac{n-1}{2})\Gamma(\frac{N}{2})}$
VLCS	10	94725	9.48e-3	10	249605	5.84e-3
PACS	10	95751	9.43e-3	10	251143	5.82e-3
OfficeHome	10	176805	6.94e-3	10	372645	4.78e-3
TerraIncognita	10	97290	9.36e-3	10	253450	5.80e-3
DomainNet	10	269145	5.62e-3	10	511065	4.08e-3

The second term relates to the mean of Beta $(\frac{n-1}{2}, \frac{N-n}{2})$:

$$\mathbb{E} [\|\mathbf{P}_{U\mathbf{a}^\perp}(\mathbf{b}_\perp)\|^2] = \frac{n-1}{N-1} \|\mathbf{b}_\perp\|^2 \leq \frac{n-1}{N-1}. \quad (61)$$

However, a tighter bound arises from the expectation of the square root:

$$\mathbb{E} \left[\frac{\|\mathbf{P}_{U\mathbf{a}^\perp}(\mathbf{b}_\perp)\|}{\|\mathbf{b}_\perp\|} \right] = \mathbb{E} [\sqrt{Y}], \quad (62)$$

$$Y \sim \text{Beta} \left(\frac{n-1}{2}, \frac{N-n}{2} \right). \quad (63)$$

For $Y \sim \text{Beta}(\alpha, \beta)$, $\mathbb{E}[\sqrt{Y}] = \frac{\Gamma(\alpha+\frac{1}{2})\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\alpha+\beta+\frac{1}{2})}$. Substituting $\alpha = \frac{n-1}{2}$, $\beta = \frac{N-n}{2}$:

$$\mathbb{E} [\|\mathbf{P}_{U\mathbf{a}^\perp}(\mathbf{b}_\perp)\|] = \|\mathbf{b}_\perp\| \frac{\Gamma(\frac{n}{2})\Gamma(\frac{N-1}{2})}{\Gamma(\frac{n-1}{2})\Gamma(\frac{N}{2})}. \quad (64)$$

Since $\|\mathbf{P}_U(\mathbf{b})\| \leq |(\mathbf{a}, \mathbf{b})| + \|\mathbf{P}_{U\mathbf{a}^\perp}(\mathbf{b}_\perp)\|$ by the triangle inequality and $\|\mathbf{b}_\perp\| < \|\mathbf{b}\| = 1$, we obtain:

$$\mathbb{E} [\|\mathbf{P}_U(\mathbf{b})\|] \leq |(\mathbf{a}, \mathbf{b})| + \frac{\Gamma(\frac{n}{2})\Gamma(\frac{N-1}{2})}{\Gamma(\frac{n-1}{2})\Gamma(\frac{N}{2})}. \quad (65)$$

□

B.2 KRYLOV SUBSPACE ANALYSIS

In this section, we conduct an empirical investigation into the behavior of the projection length of test gradients onto Krylov subspaces as the subspace dimension increases. This relationship is systematically examined across five diverse datasets and three distinct neural architectures, with the comprehensive experimental results visualized in Figures 1 through 5. As quantitatively summarized in Table 1, the improvement in cosine similarity remains remarkably modest (≤ 0.01) across all tested datasets and architectures when approximations are made via a 10-dimensional random subspace. This minimal average gain, however, contrasts sharply with the dynamic patterns revealed in the detailed subplots of Figures 1-5. For each individual experimental run and configuration, the increase in cosine similarity far surpasses the seemingly negligible average of 0.01. A particularly striking example that underscores this discrepancy is observed on the ‘‘LabelMe’’ domain within the VLCS dataset when utilizing the ViT-B/16 architecture, where the cosine similarity exhibits a substantial increase, exceeding 0.6. This pronounced enhancement strongly suggests that Krylov subspaces possess an inherent capability to align more closely with the true direction of test gradients, thereby offering a more effective and powerful framework for gradient approximation tasks.

C REPRODUCIBILITY

Our research is grounded in DomainBed², a comprehensive and widely-adopted benchmark suite designed for domain generalization studies, which is openly available under the permissive MIT license. DomainBed provides a standardized framework for evaluating the robustness of machine learning models across diverse data distributions, ensuring fair comparisons and reproducibility in domain shift scenarios. All experiments were conducted on a single NVIDIA Tesla V100 GPU.

D UTILIZATION OF LLMs

In the process of academic writing, Large Language Models (LLMs) are employed solely for language refinement purposes, including modifying grammatical errors, selecting more appropriate vocabulary, and polishing sentences to enhance clarity and flow, without being involved in any core scientific research activities such as idea generation or data analysis. This utilization leverages LLMs' strengths in text generation and syntactic understanding to improve the readability and overall quality of the manuscript while ensuring the intellectual contributions remain entirely human-driven.

E BROADER IMPACTS

This paper is primarily dedicated to the development of an effective domain generalization method aimed at addressing the pervasive challenge of domain shifts in machine learning systems. Given the prevalence of domain shifts and their detrimental impact on model performance with unseen data, our work focuses on improving model robustness and adaptability. By seeking models with better generalization, the proposed method reduces model bias toward spurious correlations, thereby promoting fairness and ethical alignment in automated decision-making. Consequently, we anticipate that this research will have a positive societal impact by contributing to more reliable and equitable AI systems. After careful consideration, we do not foresee any significant negative social implications arising from this work.

REFERENCES

- Pierre Foret, Ariel Kleiner, Hossein Mobahi, and Behnam Neyshabur. Sharpness-aware minimization for efficiently improving generalization. *arXiv preprint arXiv:2010.01412*, 2020.
- Roman Vershynin. *High-dimensional probability: An introduction with applications in data science*, volume 47. Cambridge university press, 2018.
- Max Welling and Yee W Teh. Bayesian learning via stochastic gradient langevin dynamics. In *Proceedings of the 28th international conference on machine learning (ICML-11)*, pp. 681–688, 2011.

²<https://github.com/facebookresearch/DomainBed>

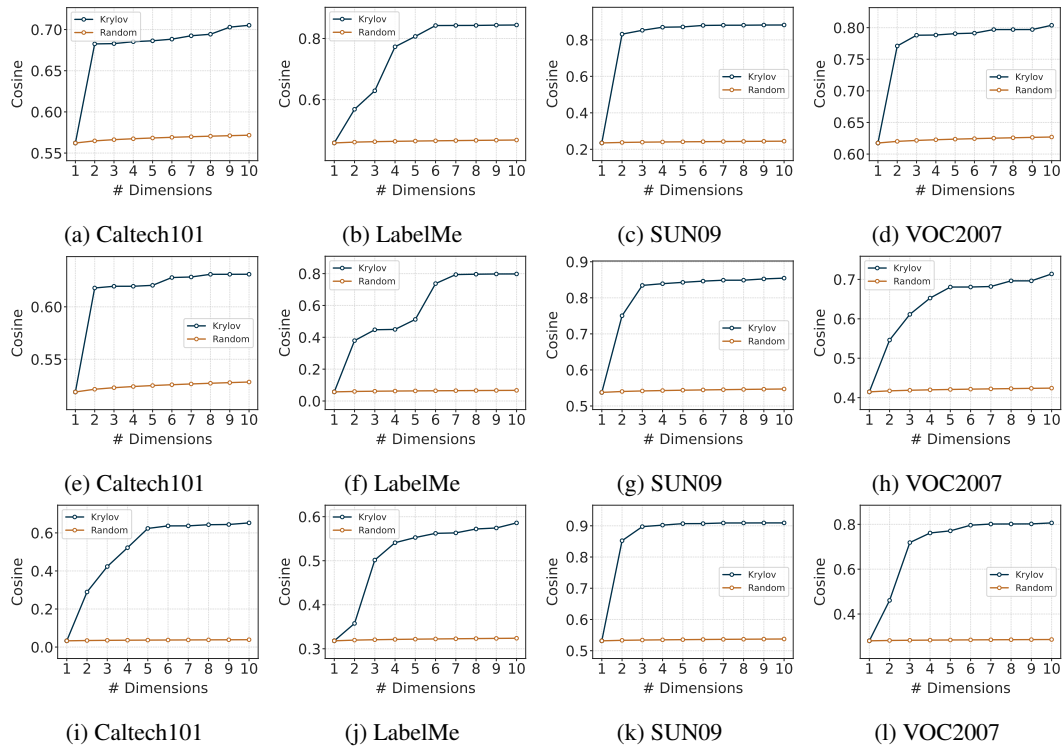


Figure 1: **Cosine Similarity Between Test Gradients and Their Krylov Subspace Projections Across Dimensions on VLCS.** The subplots in the first, second, and third rows display the performance for the ViT-B/32, ViT-B/16, and ViT-L/14 models, respectively.

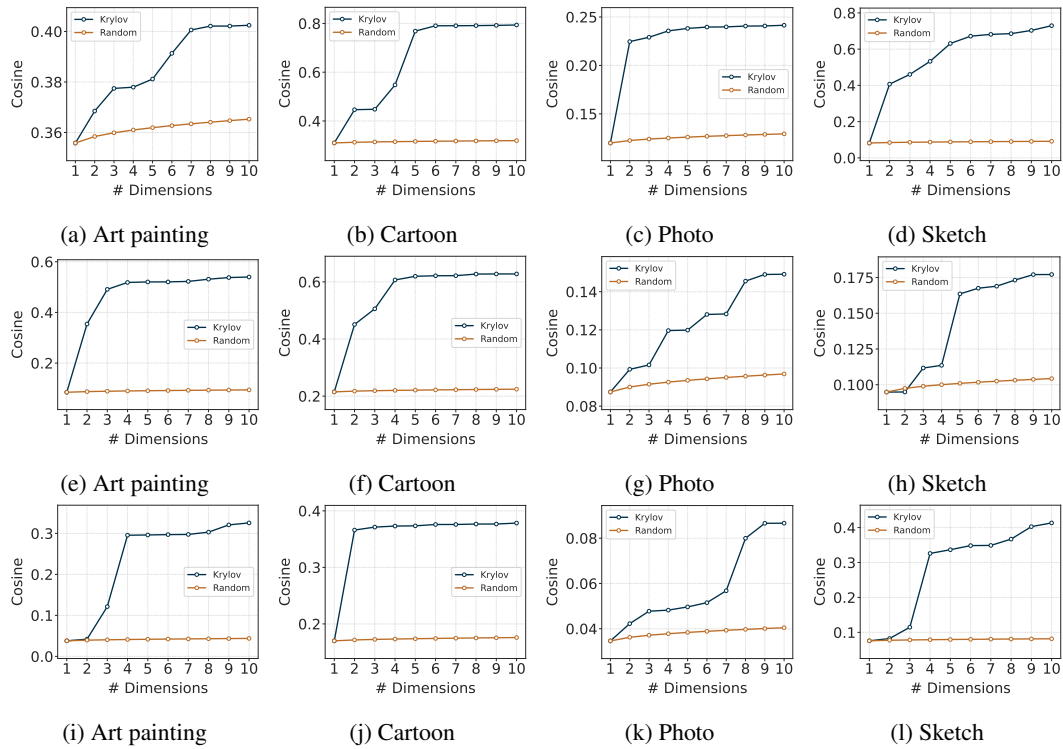


Figure 2: **Cosine Similarity Between Test Gradients and Their Krylov Subspace Projections Across Dimensions on PACS.** The subplots in the first, second, and third rows display the performance for the ViT-B/32, ViT-B/16, and ViT-L/14 models, respectively.

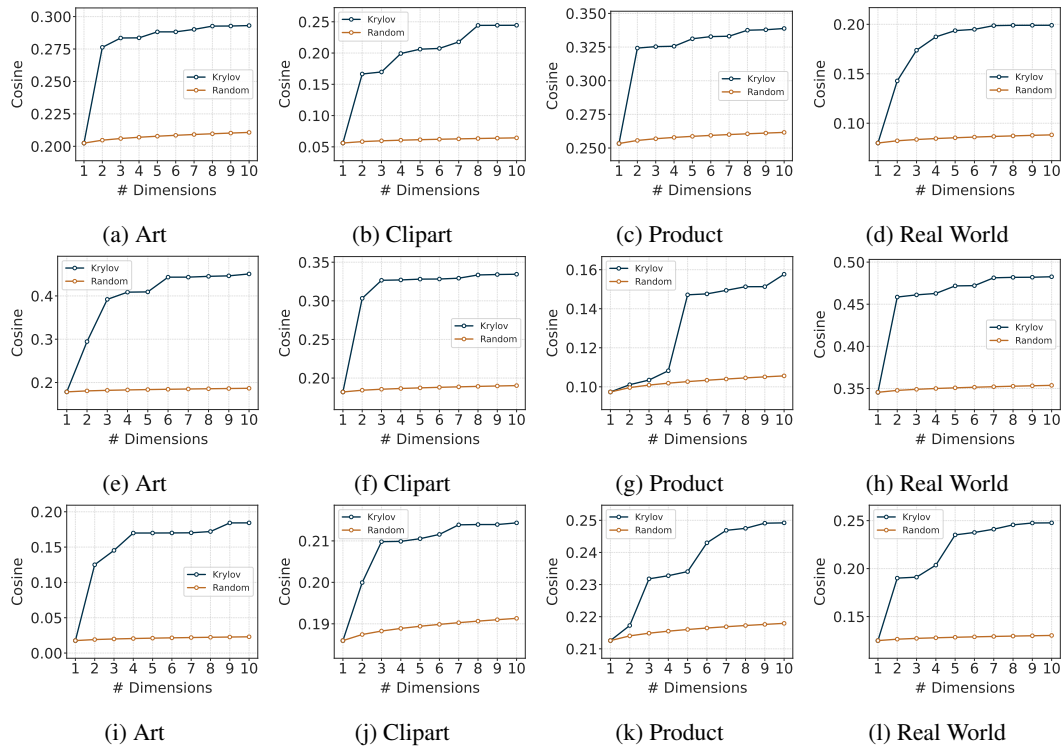


Figure 3: **Cosine Similarity Between Test Gradients and Their Krylov Subspace Projections Across Dimensions on OfficeHome.** The subplots in the first, second, and third rows display the performance for the ViT-B/32, ViT-B/16, and ViT-L/14 models, respectively.

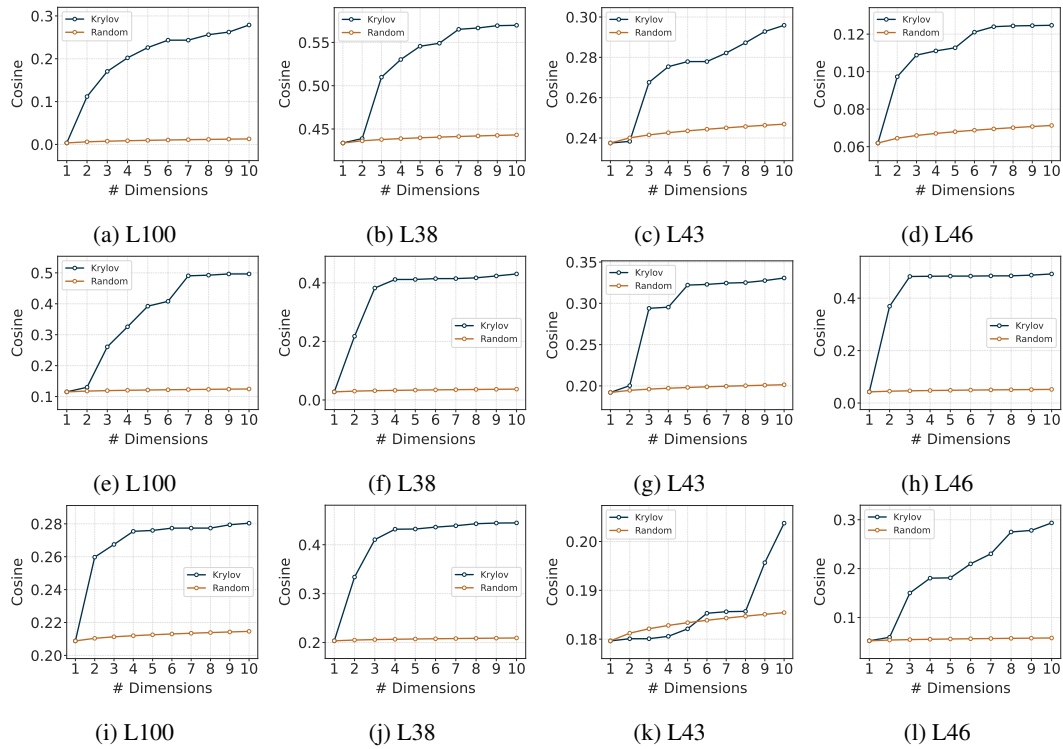


Figure 4: **Cosine Similarity Between Test Gradients and Their Krylov Subspace Projections Across Dimensions on TerraIncognita.** The subplots in the first, second, and third rows display the performance for the ViT-B/32, ViT-B/16, and ViT-L/14 models, respectively.

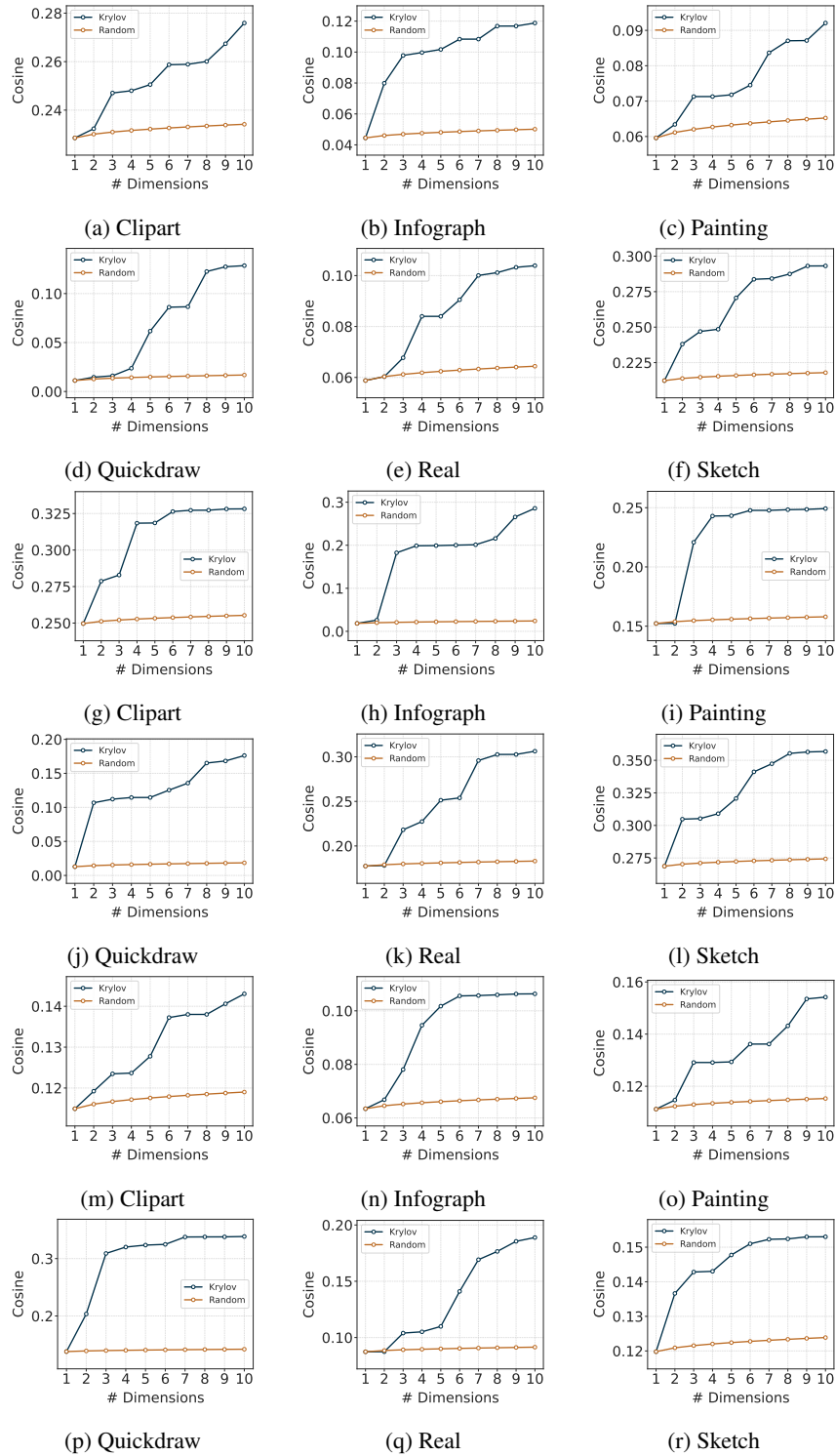


Figure 5: **Cosine Similarity Between Test Gradients and Their Krylov Subspace Projections Across Dimensions on DomainNet.** The subplots in the first two rows ((a)-(f)) correspond to the ViT-B/32 model; The subplots in the next two rows ((g)-(l)) correspond to the ViT-B/16 model; The subplots in the final two rows ((m)-(r)) correspond to the ViT-L/14 model.