

GIST: Targeted Data Selection for Instruction Tuning via Coupled Optimization Geometry

Guanghui Min

Department of Computer Science,
University of Virginia

Outline

1

Problem

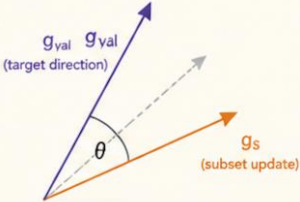
Targeted data selection



2

Baseline

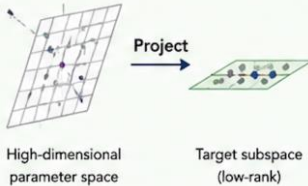
LESS and its geometry mismatch



3

Method

GIST via target subspace projection



4

Results

Better accuracy with lower cost



Outline

1

Problem

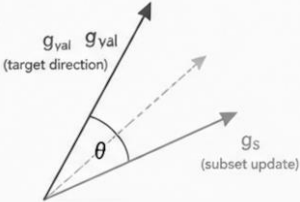
Targeted data selection



2

Baseline

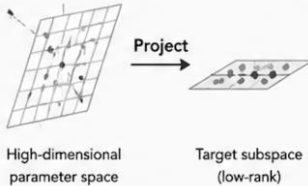
LESS and its geometry mismatch



3

Method

GIST via target subspace projection



4

Results

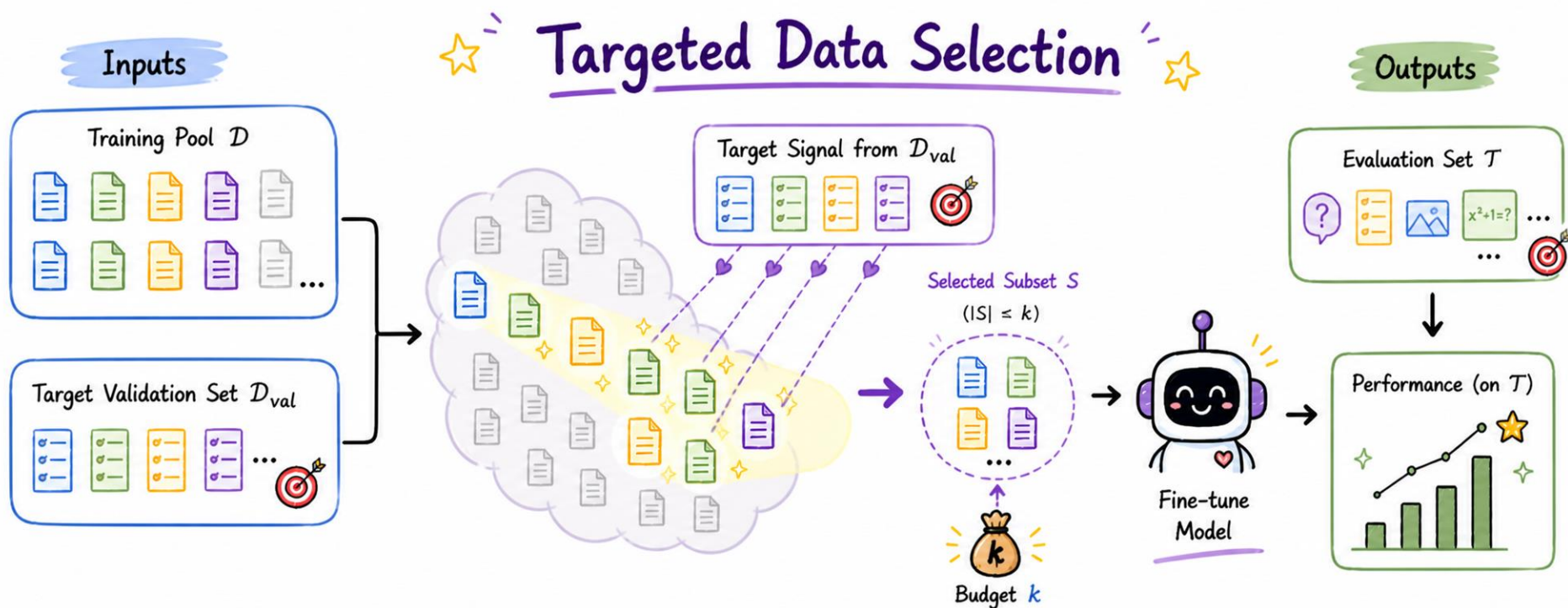
Better accuracy with lower cost



Problem Setup

Targeted Data Selection for Instruction Tuning

- **Given:** a large training pool \mathcal{D} , a small target validation set \mathcal{D}_{val} and a budget k .
- **Goal:** select $S \subset \mathcal{D}$, $|S| = k$ such that fine-tuning on S achieves the best performance on the target evaluation set \mathcal{T} .



Outline

1

Problem

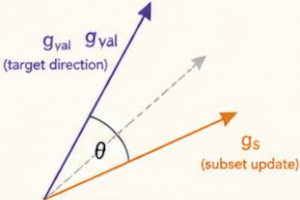
Targeted data selection



2

Baseline

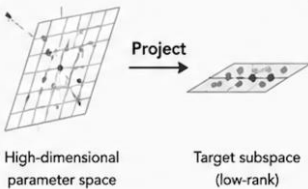
LESS and its geometry mismatch



3

Method

GIST via target subspace projection



4

Results

Better accuracy with lower cost

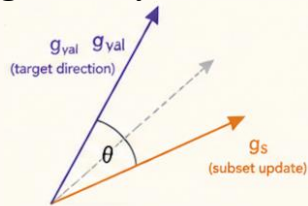


Outline

2

Baseline

LESS and its
geometry mismatch



Introduction to LESS and related works

Underline Assumption of LESS

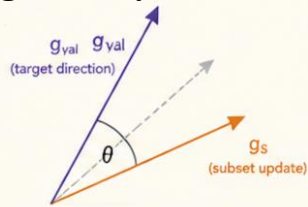
Geometry Mismatch of LESS

Outline

2

Baseline

LESS and its
geometry mismatch



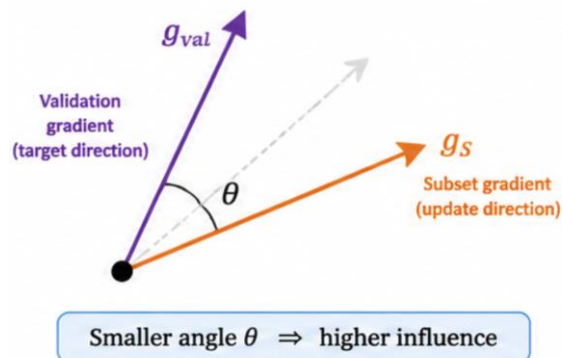
Introduction to LESS and related works

Underline Assumption of LESS

Geometry Mismatch of LESS

LESS: Gradient Alignment

Core idea: Useful data should move the model toward the target task.



Gradient alignment score:

$$\max_{S \subset \mathcal{D}, |S|=k} \nabla_{\theta_t} \mathcal{L}(\mathcal{D}_{\text{val}}, \theta_t)^\top \nabla_{\theta_t} \mathcal{L}(S, \theta_t),$$

LESS implementation:

$$\text{Influence}_{x \in \mathcal{D}, z \in \mathcal{D}_{\text{val}}}(x, z) = \sum_{i=1}^N \bar{\eta}_i \cos(\nabla \mathcal{L}(z, \theta_i), \Gamma(x, \theta_i)),$$

N : number of warmup checkpoints, $\bar{\eta}_i$: average learning rate, $\Gamma(x, \theta_i)$: Adam-preconditioned gradient.

LESS further improves the basic score with:



Adam-preconditioned gradients
Capture optimizer geometry.

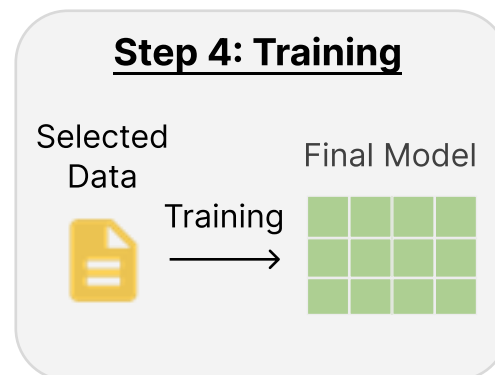
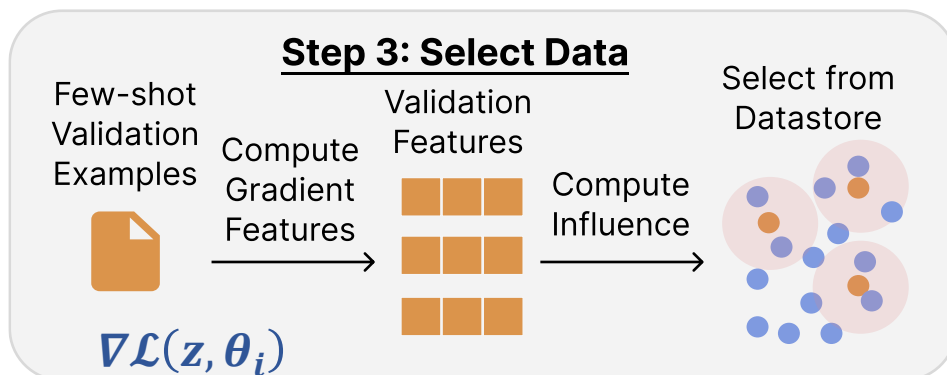
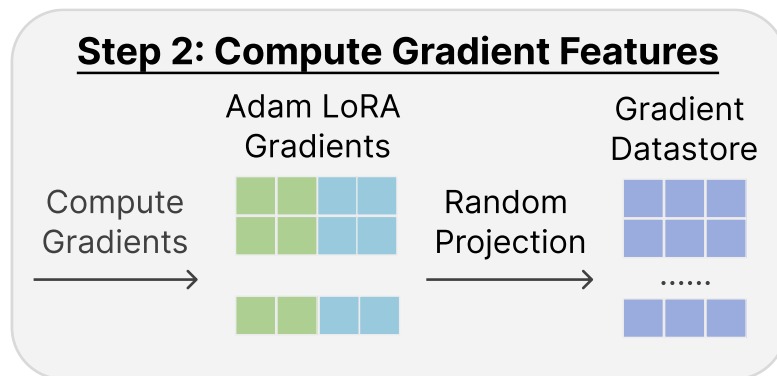
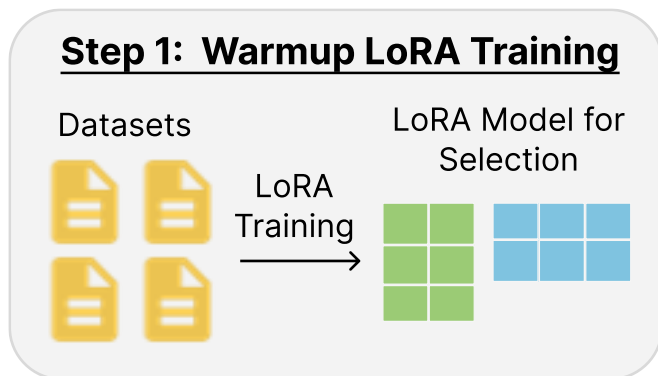


Multi-checkpoint aggregation
More stable and robust influence scores.



Length normalization
Mitigate bias from variable-length instructions.

LESS Pipeline



$$\text{Influence}(x, z) = \sum_{i=1}^N \bar{\eta}_i \cos(\nabla \mathcal{L}(z, \theta_i), \Gamma(x, \theta_i))$$

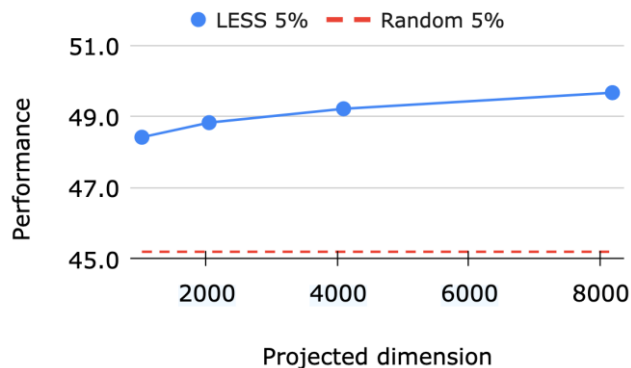
Key idea: precompute gradient features once, then use target validation examples as queries for data selection.

LESS is Expensive

Main bottleneck: computing and storing gradient features

Table 4: Asymptotic complexity, wall-clock runtime (measured as **single** A100 GPU hours) and storage cost associated with each step in LESS. Gradient computation is the most costly step, followed by the warmup LoRA training stage, but this expense is incurred only once. The actual data selection process requires minimal computation.

	Warmup LoRA Training		Gradient Features Computation		Data Selection	
	Complexity	Actual	Complexity	Actual	Complexity	Actual
Compute	$\mathcal{O}(\mathcal{D}_{\text{warmup}} \cdot N)$	6 Hours	$\mathcal{O}(\mathcal{D} \cdot N)$	48 Hours	$\mathcal{O}(\mathcal{D} \cdot \mathcal{D}_{\text{val}} \cdot d)$	< 1 Min
Storage	-	-	$\mathcal{O}(\mathcal{D} \cdot N \cdot d)$	17.7 GB	-	-



Key idea: LESS is effective, but its selection performance relies on expensive multi-checkpoint gradient features.

Table 6: Number of checkpoints (N) used for select data with LESS. Using fewer checkpoints still outperforms random selection but is less effective.

	MMLU	TYDIQA	BBH	Avg.
Random	46.5 (0.5)	52.7 (0.4)	38.9 (0.5)	46.0
$N = 1$	48.2 (0.4)	54.9 (0.4)	40.2 (0.2)	47.8
$N = 4$ (default)	50.2 (0.5)	56.2 (0.7)	41.5 (0.6)	49.3

Follow-ups Keep LESS Geometry

Recent works mainly improve the LESS framework in two directions:

➤ **Faster gradient computation:**

1. **InfDist** (Nikdan et al., NeurIPS'25): Uses JVP-based embeddings to approximate compressed Adam-gradient features, so it avoids explicitly computing all gradient features.
2. **IProX** (Chen et al., ICLR'26): Uses a smaller proxy model to compute influence features, reducing gradient-computation cost while preserving selection quality.

➤ **Better retrieval over gradient features:**

1. **G2IS** (Zhao et al., ACL'25): Builds a graph over gradient features and uses PCA on validation gradients to identify core target knowledge for retrieval.

These methods improve efficiency or retrieval, but largely keep the same Adam-gradient influence formulation.

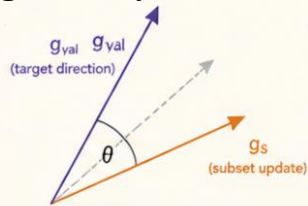
Our question: Is the Adam-preconditioned influence score the right geometry for LoRA fine-tuning?

Outline

2

Baseline

LESS and its
geometry mismatch



Introduction to LESS and related works

Underline Assumption of LESS

Geometry Mismatch of LESS

Bi-level View

Given the training set \mathcal{D} , the validation set \mathcal{D}_{val} and the selection budget k , we try to find the optimal subset $S \subset \mathcal{D}$ whose fine-tuned model minimizes target validation loss.

This naturally gives a bi-level optimization problem:

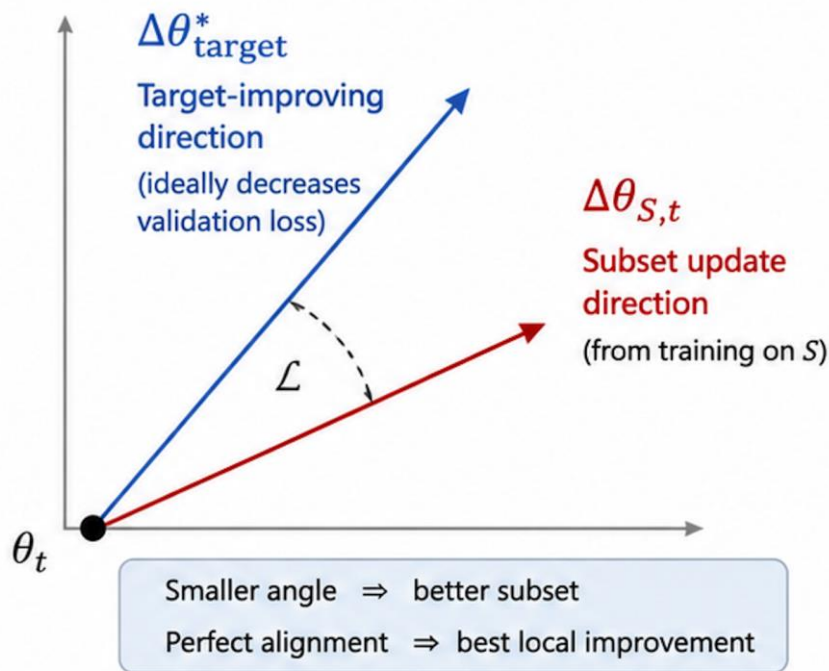
Problem 1 (Targeted Data Selection). *Find the optimal subset S^* that solves the following bi-level optimization:*

$$\begin{aligned} \min_{S \subseteq \mathcal{D}} \quad & \mathcal{L}(\mathcal{D}_{\text{val}}, \boldsymbol{\theta}_S^*) \\ \text{s.t.} \quad & |S| = k, \\ & \boldsymbol{\theta}_S^* = \arg \min_{\boldsymbol{\theta}} \mathcal{L}(S, \boldsymbol{\theta}). \end{aligned} \tag{1}$$

Fine-tune the model on the selected subset

Local Direction Matching

Intuition: choose data whose training update points toward the direction that most reduces the validation loss.



Target-improving (Newton) direction

$$\Delta\theta_{\text{target}}^* = -\mathbf{H}_{\text{val},t}^\dagger \nabla_{\theta_t} \mathcal{L}(\mathcal{D}_{\text{val}}, \theta_t)$$

$\mathbf{H}_{\text{val},t}^\dagger$: validation loss Hessian at θ_t
 \dagger : (pseudo-)inverse

Actual update from training on subset S (first-order)

$$\Delta\theta_{S,t} = -\eta \nabla_{\theta_t} \mathcal{L}(S, \theta_t).$$

Objective
(local view)

Solve **Problem 1** is locally equivalent to **maximizing the alignment** between the two directions

Ideal Influence

Bi-Level Objective

$$\begin{aligned} \min_{S \subseteq \mathcal{D}} \quad & \mathcal{L}(\mathcal{D}_{val}, \theta_S^*) \\ \text{s.t.} \quad & |S| = k, \\ & \theta_S^* = \arg \min_{\theta} \mathcal{L}(S, \theta). \end{aligned}$$



Assume Local Convexity

$$\begin{aligned} & \mathcal{L}(\mathcal{D}_{val}, \theta_t) + \nabla_{\theta_t} \mathcal{L}(\mathcal{D}_{val}, \theta_t)^\top \Delta \theta_t \\ & + \frac{1}{2} \Delta \theta_t^\top \mathbf{H}_{val,t} \Delta \theta_t \end{aligned}$$

Plugging the target-improving direction and the subset update direction into the **local alignment objective** gives:

$$\begin{aligned} \max_{\substack{S \subseteq \mathcal{D} \\ |S|=k}} \quad & \underbrace{\nabla_{\theta} L(\mathcal{D}_{val}, \theta_t)^\top}_{\text{validation gradient target direction}} \underbrace{H_{val,t}^\dagger}_{\text{target Hessian geometry}} \underbrace{\nabla_{\theta} L(S, \theta_t)}_{\text{subset gradient training update}} \end{aligned}$$

Newton Step for target set

up to constants independent of S and higher-order terms.

**Influence is not just gradient similarity.
It is gradient similarity under the target geometry.**



Adam = Diagonal Scaling

LESS replaces the raw gradient with the **Adam-preconditioned gradient**.

LESS update (per sample x)

$$\underbrace{\Gamma(x, \theta_t)}_{\text{Adam-preconditioned gradient (used by LESS)}} \approx \underbrace{D_{\text{Adam},t}}_{\text{Adam-preconditioned gradient (used by LESS)}} \underbrace{\nabla \mathcal{L}(x, \theta_t)}_{\text{raw gradient of sample } x}$$

Diagonal scaling matrix

$$D_{\text{Adam},t} = \text{diag}(G_{\text{train}} G_{\text{train}}^\top)^{-\frac{1}{2}} \quad (\text{stacked training gradient matrix})$$
$$G_{\text{train},t} = [\nabla \mathcal{L}(x, \theta_t)]_{x \in \mathcal{D}} \in \mathbb{R}^{d \times |\mathcal{D}|}$$

Geometric view
Adam multiplies the raw gradient by a diagonal matrix.

$$\Gamma(x, \theta_t) \approx \begin{bmatrix} 1 & \dots & 0 \\ g_1^\top g_1 & \dots & \vdots \\ \vdots & \ddots & \vdots \\ 0 & \dots & \frac{1}{g_{|\mathcal{D}|}^\top g_{|\mathcal{D}|}} \end{bmatrix} \cdot \nabla \mathcal{L}(x, \theta_t)$$

LESS Uses Diagonal Geometry

Ignoring EMA terms and ϵ in Adam, the influence function in LESS (single epoch) becomes an alignment under **diagonal optimizer geometry**.

Ideal Influence
(target geometry)

$$\nabla_{\theta_t} \mathcal{L}(\mathcal{D}_{\text{val}}, \theta_t)^\top \mathbf{H}_{\text{val},t}^\dagger \nabla_{\theta_t} \mathcal{L}(S, \theta_t)$$

validation gradient target Hessian geometry subset gradient

LESS Influence
(diagonal geometry)

$$\nabla_{\theta_t} \mathcal{L}(\mathcal{D}_{\text{val}}, \theta_t)^\top \left[\text{diag}(\mathbf{G}_{\text{train}} \mathbf{G}_{\text{train}}^\top) \right]^{-\frac{1}{2}} \nabla_{\theta_t} \mathcal{L}(S, \theta_t)$$

validation gradient Adam diagonal preconditioner subset gradient



Key idea

LESS measures influence under a **diagonal (axis-aligned) geometry**. It can **rescale** each coordinate independently, but cannot model **rotation** or **coupling**.

Key implication

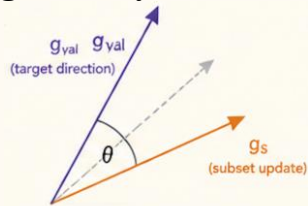
- Parameters are treated as coordinate-wise independent.
- Off-diagonal structure in the target geometry is **discarded**.

Outline

2

Baseline

LESS and its
geometry mismatch



Introduction to LESS and related works

Underline Assumption of LESS

Geometry Mismatch of LESS

Target Geometry is Low-rank

The target (validation) gradients span a low-dimensional subspace.

Stacked validation gradient matrix

$$G_{\text{val},t} = [\nabla \mathcal{L}(x, \theta_t)]_{x \in \mathcal{D}_{\text{val}}} \in \mathbb{R}^{d \times |\mathcal{D}_{\text{val}}|}$$

Target curvature proxy

$$\hat{H}_{\text{val},t} \approx G_{\text{val},t} G_{\text{val},t}^{\top} \quad \text{empirical Fisher / Gauss-Newton approximation}$$

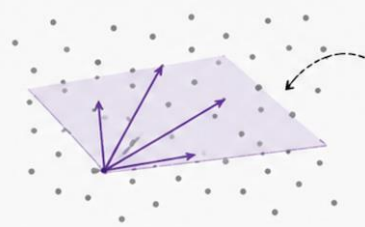
Rank upper bound

$$\text{rank}(G_{\text{val},t} G_{\text{val},t}^{\top}) = \text{rank}(\hat{H}_{\text{val},t}) \leq |\mathcal{D}_{\text{val}}| \ll d$$

What this means

The target geometry lives in a low-dimensional subspace whose dimension is at most $|\mathcal{D}_{\text{val}}|$, which is typically much smaller than d .

Parameter space (\mathbb{R}^d)



Low-dimensional subspace spanned by validation gradients (dimension $\leq |\mathcal{D}_{\text{val}}| \ll d$)

Key insight: Low-rank target geometry = a few coupled update directions. Diagonal Adam scaling cannot represent this **coupled subspace**.

Target Geometry is Low-rank

A small principal subspace captures most target-gradient variance.

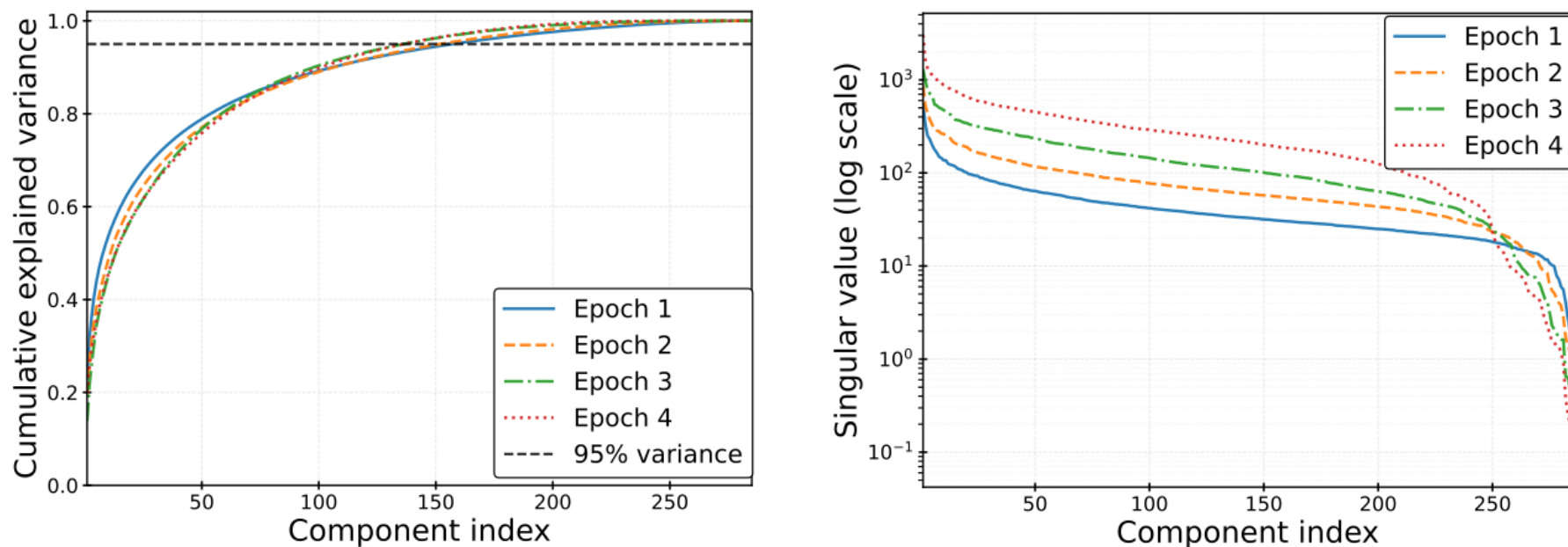


Figure 1. **Spectral analysis of the MMLU validation gradient G_{val} on Llama2-7B.** We decompose the gradient matrix via SVD to obtain singular values σ_i . (a) Cumulative explained variance. A steeper curve indicates that a **smaller principal subspace dimension is sufficient to capture the majority of the variance (e.g., Rank 150 captures 95%)**, confirming high directional information density. (b) The singular values (σ_k) exhibit **precipitous decay, further verifying the intrinsic low-rank structure**.

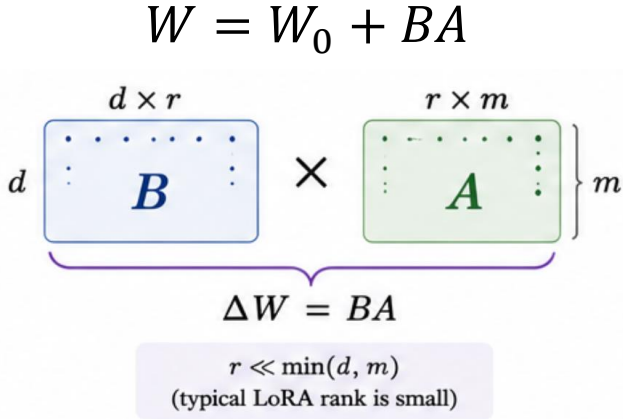
Key insight: The Observation shows that the effective rank can even be smaller than $|\mathcal{D}_{\text{val}}|$.



LoRA Induces Coupling

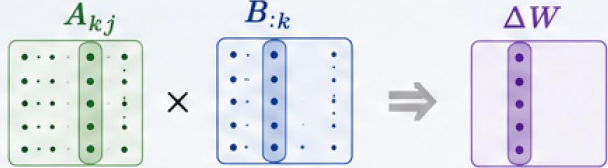
LoRA uses a bilinear parameterization. Parameters in A and B jointly determine the update ΔW , which leads to off-diagonal curvature in LoRA parameter space.

LoRA parameterization

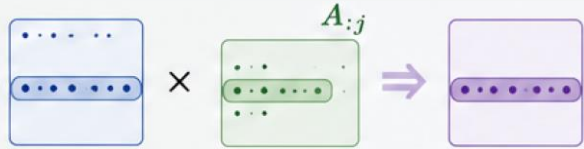


Why are A and B coupled ?

Change one element in A affects ΔW through the entire column of B .

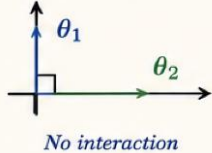


Change one element in B affects ΔW through B_i the entire row of A .

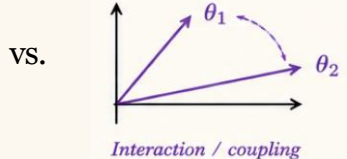


Key insight: LoRA parameters are **not** coordinate-wise independent. The bilinear BA **structure structurally induces coupling** and off-diagonal curvature.

Coordinate-wise independent (LESS's assumption)

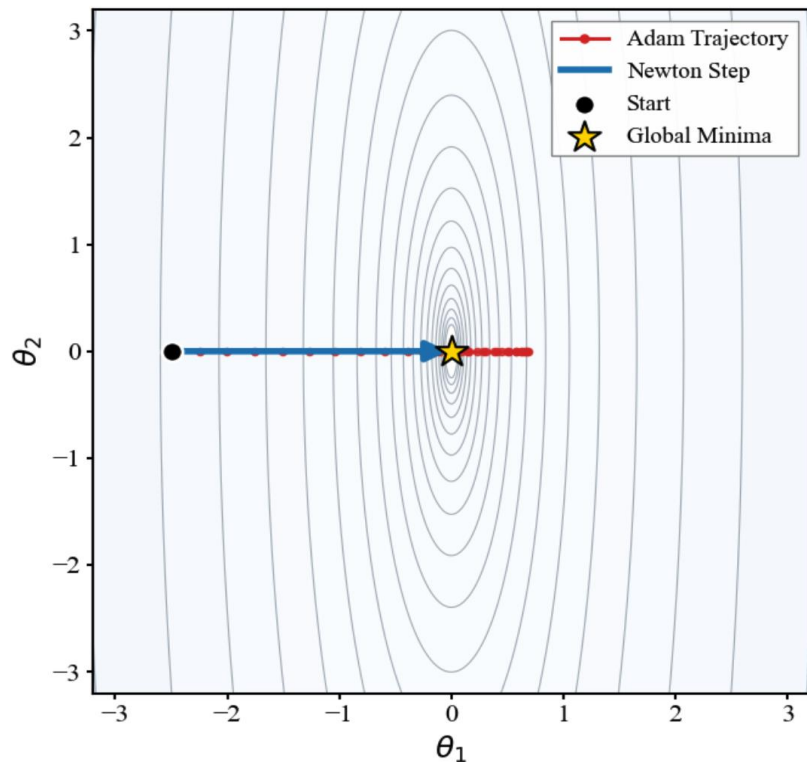


LoRA coupled geometry

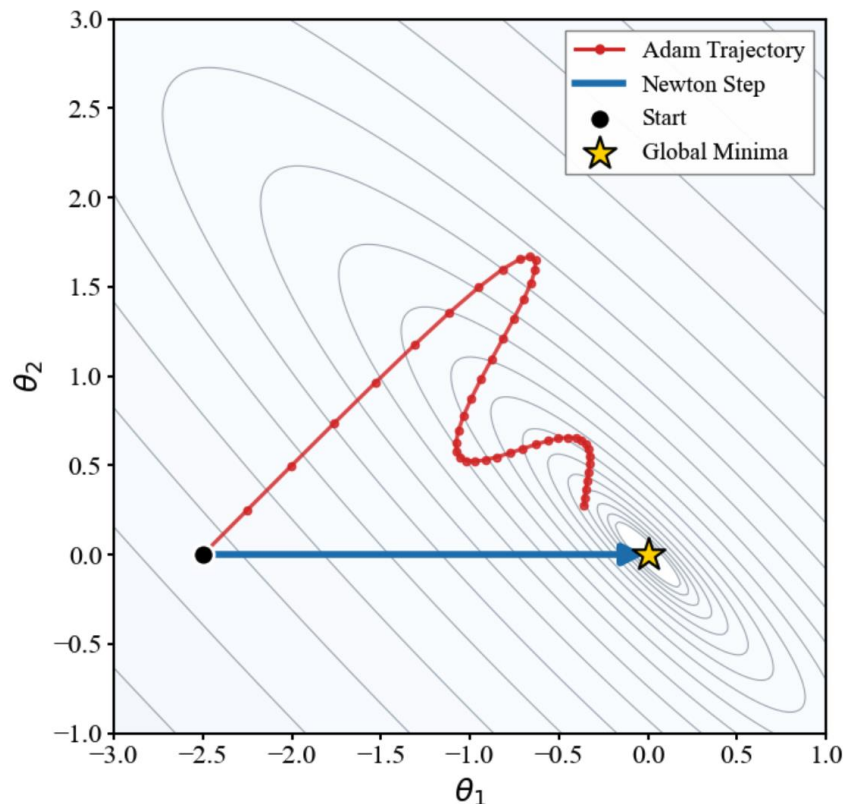


Adam Zig-zags in Coupled Geometry

Toy example: when curvature directions are coupled, diagonal Adam scaling can zig-zag instead of following the target Newton direction.



Axis-aligned geometry: Adam follows the descent direction.

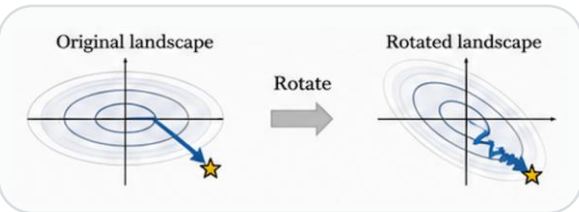


Coupled geometry: Adam zig-zags across rotated curvature directions.

Key insight: diagonal scaling works when the geometry is **axis-aligned**, but becomes **unstable** when curvature is coupled.

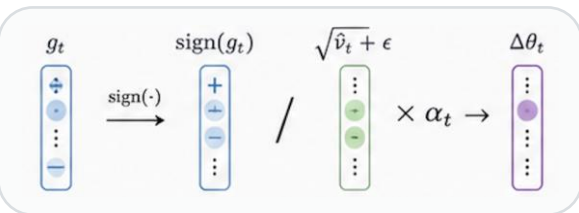
Adam is Coupling-sensitive

Several optimization studies show that Adam's diagonal adaptive scaling is not rotation-invariant and can be **unstable on coupled geometries**.



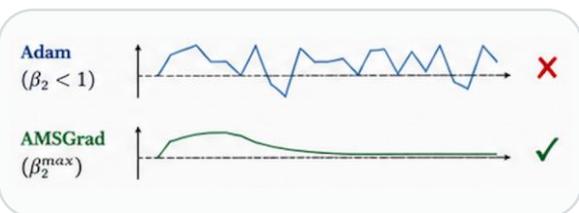
1 Adam is rotation dependent.
SGD is rotation-equivariant, while Adam can produce different trajectories under coupled loss landscapes.

Understanding Adam Requires Better Rotation Dependent Assumptions
Zhang et al., 2025
(arXiv: 2410.19964)



2 Adam is coordinate-wise.
Adam adapts the update using the sign and (second-moment) magnitude of each coordinate independently.

Dissecting Adam: The Sign, Magnitude and Variance of Stochastic Gradients
Balles & Hennig, 2018
(ICML 2018)



3 Adam can fail to converge.
For some convex objectives, Adam does not converge without modifications (e.g., AMSGrad).

On the Convergence of Adam and Beyond
Reddi et al., 2018
(ICLR 2018)

Key insight: Adam is effective in practice, but its diagonal adaptive geometry is **not designed to capture coupled curvature**.

Outline

1

Problem

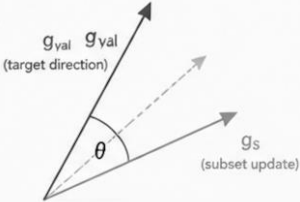
Targeted data selection



2

Baseline

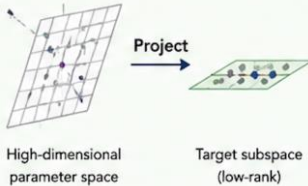
LESS and its geometry mismatch



3

Method

GIST via target subspace projection



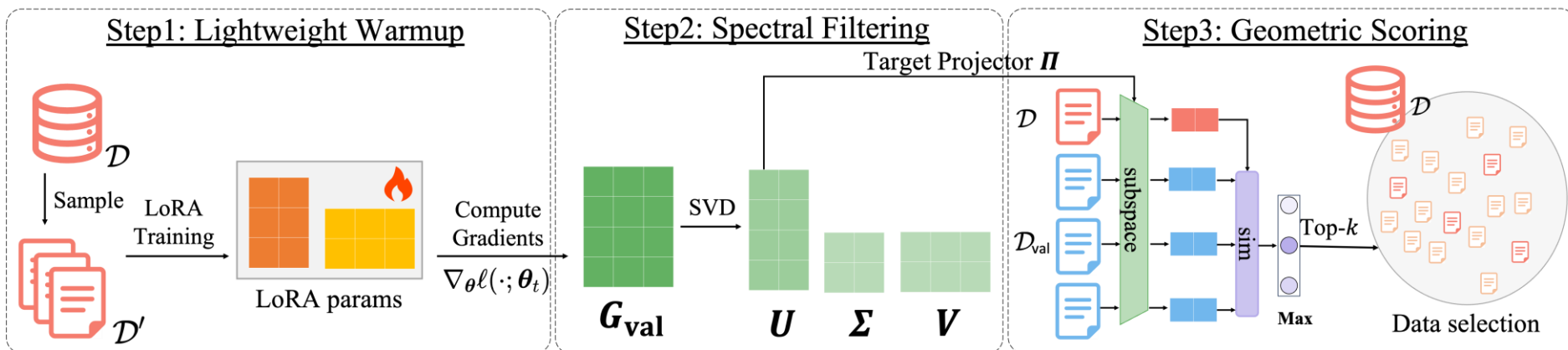
4

Results

Better accuracy with lower cost



GIST Pipeline



- **Step 1: Lightweight Warmup**

Train a LoRA model briefly and compute gradients.

- **Step 2: Spectral Filtering**

Compute SVD of validation gradients and keep top- r target directions.

- **Step 3: Geometric Scoring**

Project gradients by $\mathbf{\Pi} = \mathbf{U}_r^{\top}$, compute cosine scores, select top- k .

$$\text{Inf}_t(\mathbf{z}_i, \mathbf{z}_{\text{val}}^{(j)}) \triangleq \frac{(\mathbf{\Pi} \mathbf{g}_{i,t})^{\top} (\mathbf{\Pi} \mathbf{g}_{\text{val},t}^{(j)})}{\|\mathbf{\Pi} \mathbf{g}_{i,t}\|_2 \|\mathbf{\Pi} \mathbf{g}_{\text{val},t}^{(j)}\|_2}$$

Key idea: GIST uses the validation-gradient subspace as the target geometry, then selects training examples by projected alignment.

GIST: Spectral Projection

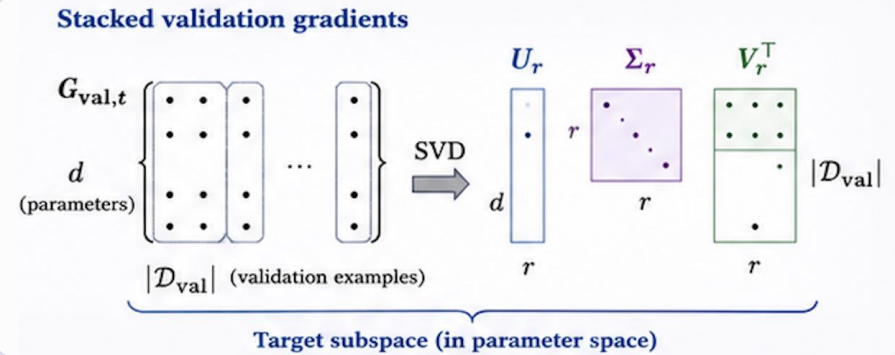
GIST replaces diagonal Adam scaling with **spectral projection onto the target subspace** identified by validation gradients.

1 Identify target subspace from validation gradients

Stack validation gradients at checkpoint θ_t into $G_{\text{val},t}$ and compute compact SVD:

$$G_{\text{val},t} = U_r \Sigma_r V_r^\top$$

$U_r \in \mathbb{R}^{d \times r}$ (parameter subspace)
 $\Sigma_r \in \mathbb{R}^{r \times r}$ (singular values)
 $V_r \in \mathbb{R}^{|\mathcal{D}_{\text{val}}| \times r}$ (example space)

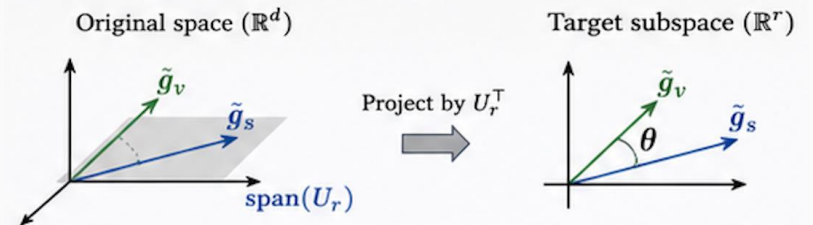


2 Use rotation (projection) to decouple

Take the task projector $\Pi = U_r^\top$.
It diagonalizes the surrogate in the subspace:

$$U_r^\top G_{\text{val},t} G_{\text{val},t}^\top U_r = \Sigma_r V_r^\top V_r \Sigma_r = \Sigma_r^2.$$

Projection geometry



3 Project gradients and compare in the subspace

project the per-point gradient by U_r^\top , and measure the alignment of the projected training gradient descent direction $U_r^\top \nabla_{\theta_t} \mathcal{L}(S, \theta_t)$ and validation gradient descent direction $U_r^\top \nabla_{\theta_t} \mathcal{L}(\mathcal{D}_{\text{val}}, \theta_t)$.

GIST Approximates Ideal Influence

GIST approximates the ideal target geometry by replacing the Hessian pseudo-inverse with a **low-rank projector from validation-gradient SVD**.

Ideal Influence
(target geometry)

$$\nabla_{\theta_t} \mathcal{L}(\mathcal{D}_{\text{val}}, \theta_t)^\top \mathbf{H}_{\text{val},t}^\dagger \nabla_{\theta_t} \mathcal{L}(S, \theta_t)$$

validation gradient target Hessian geometry subset gradient

GIST Influence
(Coupled geometry)

$$\nabla_{\theta_t} \mathcal{L}(\mathcal{D}_{\text{val}}, \theta_t)^\top [\mathbf{U}_r \mathbf{U}_r^\top] \nabla_{\theta_t} \mathcal{L}(S, \theta_t)$$

validation gradient target subspace projector subset gradient



Key idea: GIST rotates gradients into the validation-gradient subspace (target geometry) and compares alignments there, effectively **decoupling the coupled geometry** that diagonal Adam scaling cannot capture.

Outline

1

Problem

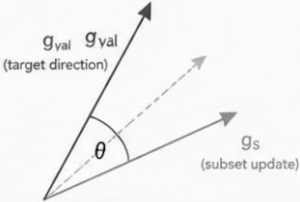
Targeted data selection



2

Baseline

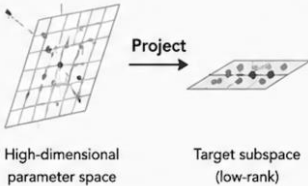
LESS and its geometry mismatch



3

Method

GIST via target subspace projection



4

Results

Better accuracy with lower cost



Evaluation Settings

- **Base Models (Big vs. Small, Modern vs. Old)**

- Llama-2-7B, Llama-3.2-3B, Qwen-2.5-1.5B

- **Training Sets (~270k samples)**

- Flan v2 (Longpre et al., 2023), Dolly (Conover et al., 2023), COT (Wei et al., 2022), Open Assistant 1 (Kopf et al., 2023).

- **Target Task**

- MMLU (Hendrycks et al., 2020)
- TydiQA (Clark et al., 2020)
- BBH (Suzgun et al., 2023)

Table 1. Statistics of evaluation datasets. We select tasks covering diverse output formats, including multiple-choice, extractive spans, and generative reasoning.

Dataset	Shots	Tasks	Data Size		Output Format
			Val.	Test	
MMLU	5	57	285	18,721	Multiple Choice
TYDIQA	1	9	9	1,713	Extractive Span
BBH	3	23	81	920	Generation (CoT)

Main Results

GIST Consistently Improves Targeted Selection

Table 2. Accuracy across datasets and models. *Base* denotes 0% (no selection) and *Full* denotes 100% (full dataset). All other methods select 5% of the data under the same finetuning budget. Gray \pm values report standard deviations. **Bold** numbers denote the best selected subset in each row, and underlined numbers denote the second best selected subset. Avg. Δ reports the average absolute improvement over *Base* over MMLU, TYDIQA and BBH.

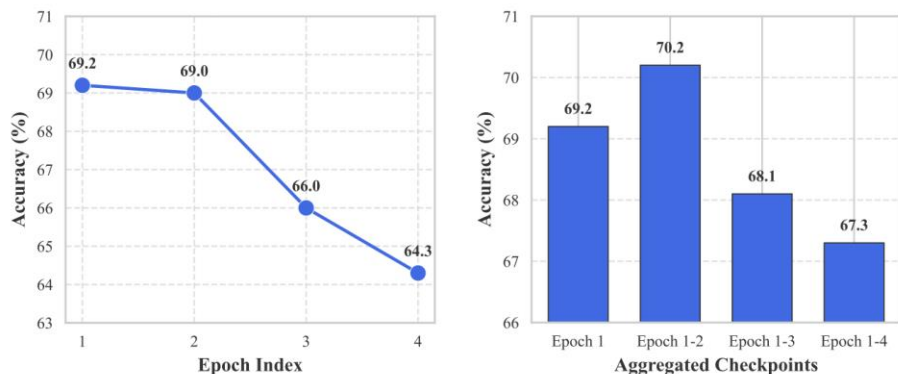
Model	Dataset	Base (0%)	Full (100%)	Rand.	Length	PPL.	Embed.	RDS+	LESS	GIST
Llama2-7B	MMLU	45.6	51.6	46.5 \pm 0.5	<u>50.5</u> \pm 0.3	38.9 \pm 0.9	47.3 \pm 0.3	46.1 \pm 0.6	50.2 \pm 0.5	51.2 \pm 0.7
	TYDIQA	46.4	54	52.7 \pm 0.4	51.1 \pm 0.3	32.9 \pm 0.4	49.8 \pm 0.8	51.0 \pm 0.3	56.2 \pm 0.7	<u>55.8</u> \pm 0.6
	BBH	38.3	43.2	38.9 \pm 0.5	39.8 \pm 0.7	34.5 \pm 0.2	38.6 \pm 0.6	39.0 \pm 0.7	<u>41.5</u> \pm 0.6	41.8 \pm 0.8
	Avg. Δ	–	+6.2	+2.6	+3.7	-8.0	+1.8	+1.9	<u>+5.9</u>	+6.2
Llama3.2-3B	MMLU	53.9	55.2	53.2 \pm 0.6	57.6 \pm 0.9	51.8 \pm 0.3	53.9 \pm 0.6	51.9 \pm 0.2	<u>56.5</u> \pm 0.6	56.1 \pm 0.4
	TYDIQA	60.4	66.6	64.1 \pm 0.4	63.9 \pm 1.3	57.1 \pm 0.7	62.8 \pm 0.4	62.6 \pm 0.5	<u>67.1</u> \pm 0.8	69.2 \pm 0.3
	BBH	45.5	47.8	45.1 \pm 0.2	45.3 \pm 0.3	43.0 \pm 0.4	44.6 \pm 0.5	45.1 \pm 0.4	<u>46.1</u> \pm 0.7	48.0 \pm 0.5
	Avg. Δ	–	+3.3	+0.9	+2.3	-2.6	+0.5	-0.1	<u>+3.3</u>	+4.5
Qwen2.5-1.5B	MMLU	62.0	62.3	61.5 \pm 0.3	<u>62.8</u> \pm 0.1	61.7 \pm 0.2	62.3 \pm 0.2	62.4 \pm 0.2	62.2 \pm 0.2	62.9 \pm 0.6
	TYDIQA	57.8	59.9	55.6 \pm 0.3	<u>56.7</u> \pm 0.3	58.4 \pm 0.4	54.1 \pm 0.3	56.2 \pm 0.3	61.4 \pm 0.8	<u>61.2</u> \pm 0.5
	BBH	43.8	44.0	43.0 \pm 0.4	42.8 \pm 0.2	43.1 \pm 0.2	<u>43.6</u> \pm 0.2	39.0 \pm 0.6	43.2 \pm 0.5	44.1 \pm 0.4
	Avg. Δ	–	+0.9	-1.2	-0.4	-0.1	-1.2	-2	<u>+1.1</u>	+1.5

Key insights:

- **Best Avg. Δ** across all three backbones.
- 5% selected data can **match or outperform** 100% full fine-tuning.
- **Compared with LESS:** +0.3, +1.2, +0.4 Avg. Δ across the three backbones.

Is One Checkpoint Enough?

Yes. Early single-checkpoint geometry is sufficient, and later/more checkpoints do not consistently help.



(a) Single Epoch Performance (b) Cumulative Performance

Figure 3. Impact of Checkpoint Selection. (a) Using single-epoch gradients shows a clear performance drop in later epochs. (b) Aggregating multiple checkpoints (weighted) does not outperform the early-stop strategy, confirming that early gradients contain the essential task optimization directions.

Table 3. Number of checkpoints (N) used for select data with GIST for Llama3.2-3B.

	MMLU	TyDIQA	BBH	Avg.
LESS	56.5 \pm 0.6	67.1 \pm 0.8	46.1 \pm 0.7	56.6
$N = 4$	56.3 \pm 0.3	67.3 \pm 0.1	46.8 \pm 0.4	56.8
$N = 1$ (default)	56.1 \pm 0.4	69.2 \pm 0.3	48.0 \pm 0.2	57.8

Table. Warmup ablation with GIST for Llama3.2-3B.

Task	Base	GIST w/o warmup	GIST
MMLU	53.9	53.0 \pm 0.2	56.1\pm0.4
TyDIQA	60.4	64.5 \pm 1.7	69.2\pm0.3
BBH	45.5	46.2 \pm 0.4	48.0\pm0.5

Key insights: GIST only needs a short warmup and one early checkpoint to recover useful target geometry. Later or aggregated checkpoints do not consistently improve selection.

GIST is 4× faster

Table 4. Efficiency comparison between GIST (Ours) and LESS. Runtime is measured in **single A100 GPU hours**. GIST significantly reduces the overhead in Warmup and Feature Extraction stages (approx. 4× speedup) and requires negligible time for SVD.

Method	Metric	Warmup	Target SVD	Grad. Feats.
LESS	Time	6.0 h	–	48.0 h
	Compl.	$\mathcal{O}(\mathcal{D}' \cdot N)$	–	$\mathcal{O}(\mathcal{D} \cdot N)$
GIST	Time	1.5 h	< 1 m	12.0 h
	Compl.	$\mathcal{O}(\mathcal{D}')$	$\mathcal{O}(\mathcal{D}_{\text{val}} ^2 d)$	$\mathcal{O}(\mathcal{D})$

Key insights:

- GIST reduces the dominant cost of targeted selection by **avoiding multi-checkpoint** feature extraction.
- For each target task, runtime drops from 54.0h to 13.5h, with warmup reduced from 6.0h to 1.5h and gradient-feature extraction reduced from 48.0h to 12.0h.
- The extra target **SVD cost is negligible**, taking less than 1 minute.
- Storage is also reduced from 75GB to 217MB, about a 350× reduction.

Thanks for listening!

