

486  
487  
488  
489  
490  
491  
492  
493  
494  
495  
496  
497  
498  
499  
500  
501  
502  
503  
504  
505  
506  
507  
508  
509  
510  
511  
512  
513  
514  
515  
516  
517  
518  
519  
520  
521  
522  
523  
524  
525  
526  
527  
528  
529  
530  
531  
532  
533  
534  
535  
536  
537  
538  
539

---

## Supplementary Material: TaskMixPGM: Task Mixtures via Probabilistic Graphical Modelling for Language Model Finetuning

---

### CONTENTS

<b>Appendices</b>	<b>10</b>
<b>A Organization of the Appendix</b>	<b>11</b>
<b>B Broader Impact</b>	<b>11</b>
<b>C Main Theoretical Results</b>	<b>11</b>
C.1 Closed-Form Solution of Quadratic Minimization over the Simplex . . . . .	11
C.2 Lagrangian and First-Order Conditions . . . . .	12
C.3 Solution under Interior Assumption . . . . .	12
C.4 Discussion . . . . .	13
C.5 Monotonicity and Submodular Properties of Energy Potential . . . . .	13
C.6 Weak Submodularity of Set function $f$ . . . . .	14
<b>D Comparative Analysis across various Notions of Task Similarities</b>	<b>16</b>
D.1 Similarity across Task Vectors via Linearized finetuning . . . . .	16
D.2 Algorithms for computing PMI and JSD . . . . .	18
<b>E Experimental Details</b>	<b>20</b>
<b>F Additional Results</b>	<b>21</b>
<b>G Code</b>	<b>25</b>

540  
541  
542  
543  
544  
545  
546  
547  
548  
549  
550  
551  
552  
553  
554  
555  
556  
557  
558  
559  
560  
561  
562  
563  
564  
565  
566  
567  
568  
569  
570  
571  
572  
573  
574  
575  
576  
577  
578  
579  
580  
581  
582  
583  
584  
585  
586  
587  
588  
589  
590  
591  
592  
593

---

## Supplementary Material: TaskMixPGM: Task Mixtures via Probabilistic Graphical Modelling for Language Model Finetuning

---

### A ORGANIZATION OF THE APPENDIX

This appendix provides supporting material for the main text, organized into the following sections. Section B presents the overall broader impact of our work. Section C presents the theoretical foundations underpinning our approach, including monotonicity and submodularity results relevant to energy-based models. Section D provides a comparative analysis of task similarity measures, starting with linearized fine-tuning vectors and extending to distributional metrics such as Pointwise Mutual Information (PMI) and Jensen-Shannon Divergence (JSD), along with algorithms for their computation. Section E details the experimental setup, datasets, and model configurations used in our evaluations. Section F includes extended results, such as tabular comparisons, that complement those in the main paper. Finally, Section G outlines the structure of our codebase and provides guidance for reproducing the experiments.

### B BROADER IMPACT

Our proposed work on TASKMIXPGM has significant broader impact across multiple domains of machine learning research and real-world applications.

- In **natural language understanding and multilingual benchmarks**, the selection of fine-tuning data mixtures is critical to model generalization. By explicitly optimizing for both representativeness and diversity, TASKMIXPGM enhances performance on complex, multi-domain evaluations such as MMLU and BIG-Bench-Hard. This enables more robust LLMs capable of reasoning across languages, topics, and task formats.
- In **AI deployment for low-resource and specialized domains**, TASKMIXPGM provides a scalable and principled solution to constructing effective mixtures from limited or domain-specific task collections. Applications include legal document analysis, medical QA, and scientific literature synthesis—areas where manually tuning mixtures is costly and error-prone.
- In **AI safety and interpretability** research, our framework offers interpretable insights into task interactions and data influence. The use of functional similarity via output divergences, rather than opaque semantic features, facilitates transparency in fine-tuning decisions. This can assist auditing pipelines and mitigate risks associated with over-representation of narrow task distributions.
- In **efficient model training and green AI initiatives**, TASKMIXPGM can reduce unnecessary computation and data usage by guiding mixture construction toward high-impact tasks. This aligns with ongoing efforts to lower the carbon footprint of large-scale model development while maintaining or improving downstream performance.

### C MAIN THEORETICAL RESULTS

#### C.1 CLOSED-FORM SOLUTION OF QUADRATIC MINIMIZATION OVER THE SIMPLEX

We consider the problem of minimizing a quadratic energy function over the probability simplex  $\Delta_n = \{\mathbf{p} \in \mathbb{R}^n : \mathbf{p}^\top \mathbf{1}_n = 1, \mathbf{p} \geq 0\}$ :

$$\min_{\mathbf{p} \in \Delta_n} E(\mathbf{p}) := -\Psi_{\text{un}}^\top \mathbf{p} + \frac{1}{2} \mathbf{p}^\top \Psi_{\text{pair}} \mathbf{p} \quad (7)$$

where  $\Psi_{\text{un}} \in \mathbb{R}^n$  denotes a unary potential vector and  $\Psi_{\text{pair}} \in \mathbb{R}^{n \times n}$  is a symmetric positive semi-definite (PSD) matrix encoding pairwise interactions.

## C.2 LAGRANGIAN AND FIRST-ORDER CONDITIONS

To enforce the affine constraint  $\mathbf{p}^\top \mathbf{1}_n = 1$ , and inequality constraints  $\mathbf{p} \geq 0$ , we consider the KKT conditions for optimality. Define the Lagrangian:

$$L(\mathbf{p}, \nu, \boldsymbol{\mu}) = -\Psi_{\text{un}}^\top \mathbf{p} + \frac{1}{2} \mathbf{p}^\top \Psi_{\text{pair}} \mathbf{p} + \nu (\mathbf{p}^\top \mathbf{1}_n - 1) - \boldsymbol{\mu}^\top \mathbf{p} \quad (8)$$

with dual variables  $\nu \in \mathbb{R}$  (equality) and  $\boldsymbol{\mu} \in \mathbb{R}_+^n$  (inequality).

The **KKT optimality conditions** are:

Stationarity Condition

$$\begin{aligned} \nabla_{\mathbf{p}} L &= -\Psi_{\text{un}} + \Psi_{\text{pair}} \mathbf{p} + \nu \mathbf{1}_n - \boldsymbol{\mu} = 0 \\ \frac{\partial L}{\partial \mathbf{p}_{[i]}} &= -\Psi_{\text{un}[i]} + \sum_{j=1}^n \Psi_{\text{pair}[ij]} \mathbf{p}_{[j]} + \nu - \boldsymbol{\mu}_{[i]} = 0. \end{aligned}$$

This implies:

$$\nu = \Psi_{\text{un}[i]} - \sum_{j=1}^n \Psi_{\text{pair}[ij]} \mathbf{p}_{[j]} + \boldsymbol{\mu}_{[i]}.$$

Primal Feasibility

$$0 \leq \mathbf{p}_{[i]} \leq 1, \quad \text{for } i = 1, \dots, n, \quad \sum_{i=1}^n \mathbf{p}_{[i]} = 1$$

Dual Feasibility

$$\nu \in \mathbb{R}^n, \quad \boldsymbol{\mu}_{[i]} \geq 0$$

Complementary Slackness

$$\boldsymbol{\mu}_{[i]} \mathbf{p}_{[i]} = 0 \quad \forall i \in [n]$$

Coordinate wise analysis for each edge cases

- **(Interior Points)**  $0 < \mathbf{p}_{[i]} < 1$  Due to complementary slackness, we have  $\boldsymbol{\mu}_{[i]} = 0 \quad \forall i \in [n]$ . Hence  $\nu = \Psi_{\text{un}[i]} - \sum_{j=1}^n \Psi_{\text{pair}[ij]} \mathbf{p}_{[j]}$  and therefore  $\mathbf{p}_{[i]} = \sum_{j=1}^n \Psi_{\text{pair}[ij]}^{-1} (\Psi_{\text{un}} - \nu \mathbf{1}_n)$
- **(Boundary Point)**  $\mathbf{p}_{[i]} = 0$  From complementary slackness, we know  $\boldsymbol{\mu}_{[i]} \geq 0$

Let  $k$  points lie in the interior and  $n - k$  points lie on the boundary

$$\sum_{i \in k} e^{\frac{s_{\alpha}[i] - \beta}{\lambda} - 1} + \sum_{i \in (n-k)} \mathbf{0}_{[i]} = 1$$

## C.3 SOLUTION UNDER INTERIOR ASSUMPTION

We first consider the case where the solution lies in the relative interior of the simplex; that is,  $\mathbf{p}^* > 0$  and hence  $\boldsymbol{\mu} = \mathbf{0}$ . Substituting into (??), we obtain:

$$\Psi_{\text{pair}} \mathbf{p} = \Psi_{\text{un}} - \nu \mathbf{1}_n \quad (9)$$

Assuming  $\Psi_{\text{pair}}$  is invertible (i.e., strictly positive definite), we may solve:

$$\mathbf{p} = \Psi_{\text{pair}}^{-1} \Psi_{\text{un}} - \nu \Psi_{\text{pair}}^{-1} \mathbf{1}_n \quad (10)$$

Imposing the constraint  $\mathbf{p}^\top \mathbf{1}_n = 1$ , we find:

$$\mathbf{1}_n^\top \mathbf{p} = \mathbf{1}_n^\top \Psi_{\text{pair}}^{-1} \Psi_{\text{un}} - \nu \mathbf{1}_n^\top \Psi_{\text{pair}}^{-1} \mathbf{1}_n = 1 \quad (11)$$

Letting

$$a := \mathbf{1}_n^\top \Psi_{\text{pair}}^{-1} \Psi_{\text{un}}, \quad b := \mathbf{1}_n^\top \Psi_{\text{pair}}^{-1} \mathbf{1}_n,$$

we obtain  $\nu = \frac{a-1}{b}$ .

Substituting back into the expression for  $\mathbf{p}$ , we conclude:

$$\mathbf{p}^* = \Psi_{\text{pair}}^{-1} \Psi_{\text{un}} - \frac{\mathbf{1}_n^\top \Psi_{\text{pair}}^{-1} \Psi_{\text{un}} - 1}{\mathbf{1}_n^\top \Psi_{\text{pair}}^{-1} \mathbf{1}_n} \cdot \Psi_{\text{pair}}^{-1} \mathbf{1}_n \quad (12)$$

#### C.4 DISCUSSION

The closed-form expression (12) satisfies the affine constraint by construction. If  $\mathbf{p}^* \geq 0$  componentwise, it is the unique global minimizer. Otherwise, if any coordinate is negative, the interior assumption fails, and active-set refinement or projection onto the simplex is required. In practice, one may use projection-based algorithms (e.g., conditional gradient, projected gradient descent) or iteratively restrict to the support set of nonnegative entries and resolve (12) over that face of the simplex.

#### C.5 MONOTONICITY AND SUBMODULAR PROPERTIES OF ENERGY POTENTIAL

**Lemma 1** (Monotonicity). *Let  $f$  be the set function defined in Eq (4). Then  $f$  is monotonic: for any sets  $\tilde{A} \subseteq \tilde{B}$ ,  $f(\tilde{A}) \leq f(\tilde{B})$ .*

*Proof.* Let  $|\tilde{A}| = n_1$  and  $|\tilde{B}| = n_2$  and since  $\tilde{A} \subseteq \tilde{B}$  we have  $n_1 < n_2$ . We index the elements in  $\tilde{B}$  such that the first  $n_1$  elements are contained in  $\tilde{A}$ .

$$f(\tilde{B}) = \max_{\mathbf{p} \in \Delta_{n_2}^{\mathbb{R}}; \text{supp}(\mathbf{p}) \subseteq \tilde{B}} \bar{\mathbb{E}}(\mathbf{p}) \geq \max_{\mathbf{p} \in \Delta_{n_1}^{\mathbb{R}}; \text{supp}(\mathbf{p}) \subseteq \tilde{A}} \bar{\mathbb{E}}(\mathbf{p}) = f(\tilde{A})$$

□

This indicates the function under consideration is monotonically increasing under task mixture.

**Lemma 2** (Finite RSC and RSM of Quadratic Term). *Let  $\mathbf{S} \in \mathbb{R}^{n \times n}$  be a symmetric positive definite similarity matrix. Then the quadratic function  $\mathbb{E}(\mathbf{p}) = \mathbf{p}^\top \mathbf{S} \mathbf{p}$  satisfies Restricted Strong Convexity (RSC) and Restricted Smoothness (RSM) over the probability simplex  $\Delta_n = \{\mathbf{p} \in \mathbb{R}^n : \mathbf{p} \geq 0, \|\mathbf{p}\|_1 = 1\}$  with finite constants  $c_\Omega > 0$  and  $C_\Omega > 0$ , respectively. That is, for all  $\mathbf{p}, \mathbf{q} \in \Delta_n$ ,*

$$\frac{c_\Omega}{2} \|\mathbf{p} - \mathbf{q}\|_2^2 \leq \mathbb{E}(\mathbf{p}) - \mathbb{E}(\mathbf{q}) - \nabla \mathbb{E}(\mathbf{q})^\top (\mathbf{p} - \mathbf{q}) \leq \frac{C_\Omega}{2} \|\mathbf{p} - \mathbf{q}\|_2^2.$$

*Proof.* Let  $\mathbb{E}(\mathbf{p}) := \mathbf{p}^\top \mathbf{S} \mathbf{p}$  denote the energy of the task mixture  $\mathbf{p} \in \Delta$ , where  $\mathbf{S} \in \mathbb{R}^{n \times n}$  is a symmetric positive definite similarity matrix and  $n$  denotes the total number of tasks. We may express the second-order Taylor expansion of  $\mathbb{E}$  as:

$$\mathbb{E}(\mathbf{p}) = \mathbb{E}(\mathbf{q}) + \nabla \mathbb{E}(\mathbf{q})^\top (\mathbf{p} - \mathbf{q}) + \frac{1}{2} (\mathbf{p} - \mathbf{q})^\top \nabla^2 \mathbb{E}(\xi) (\mathbf{p} - \mathbf{q})$$

for some  $\xi$  on the line segment between  $\mathbf{p}$  and  $\mathbf{q}$ .

Since  $\nabla \mathbb{E}(\mathbf{p}) = 2\mathbf{S}\mathbf{p}$  and  $\nabla^2 \mathbb{E}(\mathbf{p}) = 2\mathbf{S}$  is constant over  $\mathbf{p}$ , we simplify the residual energy term:

$$\mathbb{E}(\mathbf{p}) - \mathbb{E}(\mathbf{q}) - \nabla \mathbb{E}(\mathbf{q})^\top (\mathbf{p} - \mathbf{q}) = (\mathbf{p} - \mathbf{q})^\top \mathbf{S} (\mathbf{p} - \mathbf{q})$$

We now invoke spectral bounds on the quadratic form. Let  $\lambda_{\min}(\mathbf{S})$ ,  $\lambda_{\max}(\mathbf{S})$  denote the smallest and largest eigenvalues of  $\mathbf{S}$ . Since  $\mathbf{S} > 0$ , we have:

$$\lambda_{\min}(\mathbf{S}) \|\mathbf{p} - \mathbf{q}\|_2^2 \leq (\mathbf{p} - \mathbf{q})^\top \mathbf{S} (\mathbf{p} - \mathbf{q}) \leq \lambda_{\max}(\mathbf{S}) \|\mathbf{p} - \mathbf{q}\|_2^2$$

Combining with the expression above, we obtain the sandwich bound:

$$\lambda_{\min}(\mathbf{S}) \|\mathbf{p} - \mathbf{q}\|_2^2 \leq \mathbb{E}(\mathbf{p}) - \mathbb{E}(\mathbf{q}) - \nabla \mathbb{E}(\mathbf{q})^\top (\mathbf{p} - \mathbf{q}) \leq \lambda_{\max}(\mathbf{S}) \|\mathbf{p} - \mathbf{q}\|_2^2$$

Defining  $c_\Omega := 2\lambda_{\min}(\mathbf{S})$  and  $L := 2\lambda_{\max}(\mathbf{S})$ , we conclude that  $\mathbb{E}(\mathbf{p})$  is  $(c_\Omega, C_\Omega)$ -restricted strongly convex and smooth over  $\Delta$  in the sense that:

$$\frac{c_\Omega}{2} \|\mathbf{p} - \mathbf{q}\|_2^2 \leq \mathbb{E}(\mathbf{p}) - \mathbb{E}(\mathbf{q}) - \nabla \mathbb{E}(\mathbf{q})^\top (\mathbf{p} - \mathbf{q}) \leq \frac{C_\Omega}{2} \|\mathbf{p} - \mathbf{q}\|_2^2$$

□

**Lemma 3** (Finite RSC and RSM of Eq: 1 Energy Potential). *Let  $\mathbf{S} \in \mathbb{R}^{n \times n}$  be a symmetric positive definite similarity matrix. Then the quadratic function  $\mathbb{E}(\mathbf{p}) = -\Psi_{\text{un}}^\top \mathbf{p} + \frac{1}{2} \mathbf{p}^\top \Psi_{\text{pair}} \mathbf{p}$  satisfies Restricted Strong Convexity (RSC) with parameter  $c_\Omega$  and Restricted Smoothness (RSM) with parameter  $C_\Omega$  over the probability simplex  $\Delta_n = \{\mathbf{p} \in \mathbb{R}^n : \mathbf{p} \geq 0, \|\mathbf{p}\|_1 = 1\}$  with finite constants  $c_\Omega > 0$  and  $C_\Omega > 0$ , respectively. That is, for all  $\mathbf{p}, \mathbf{q} \in \Delta_n$ ,*

$$\frac{c_\Omega}{2} \|\mathbf{p} - \mathbf{q}\|_2^2 \leq \mathbb{E}(\mathbf{p}) - \mathbb{E}(\mathbf{q}) - \nabla \mathbb{E}(\mathbf{q})^\top (\mathbf{p} - \mathbf{q}) \leq \frac{C_\Omega}{2} \|\mathbf{p} - \mathbf{q}\|_2^2.$$

*Proof.* We begin by analyzing the structure of the energy function  $\mathbb{E} : \mathbb{R}^n \rightarrow \mathbb{R}$ , defined as

$$\mathbb{E}(\mathbf{p}) = -\Psi_{\text{un}}^\top \mathbf{p} + \frac{1}{2} \mathbf{p}^\top \Psi_{\text{pair}} \mathbf{p}.$$

This function is a standard quadratic form, with gradient and Hessian given by

$$\nabla \mathbb{E}(\mathbf{p}) = \Psi_{\text{pair}} \mathbf{p} - \Psi_{\text{un}}, \quad \nabla^2 \mathbb{E}(\mathbf{p}) = \Psi_{\text{pair}}.$$

Since  $\Psi_{\text{pair}}$  is symmetric positive definite, it admits an eigenvalue decomposition  $\Psi_{\text{pair}} = \mathbf{U} \Lambda \mathbf{U}^\top$  with eigenvalues  $0 < \lambda_1 \leq \dots \leq \lambda_n$ . Let  $c_\Omega := \lambda_{\min}(\Psi_{\text{pair}})$  and  $C_\Omega := \lambda_{\max}(\Psi_{\text{pair}})$ .

We now apply the standard second-order Taylor expansion of  $\mathbb{E}$  at  $\mathbf{q} \in \Delta$  evaluated at  $\mathbf{p} \in \Delta$ :

$$\mathbb{E}(\mathbf{p}) = \mathbb{E}(\mathbf{q}) + \nabla \mathbb{E}(\mathbf{q})^\top (\mathbf{p} - \mathbf{q}) + \frac{1}{2} (\mathbf{p} - \mathbf{q})^\top \Psi_{\text{pair}} (\mathbf{p} - \mathbf{q}),$$

and hence,

$$\mathbb{E}(\mathbf{p}) - \mathbb{E}(\mathbf{q}) - \nabla \mathbb{E}(\mathbf{q})^\top (\mathbf{p} - \mathbf{q}) = \frac{1}{2} (\mathbf{p} - \mathbf{q})^\top \Psi_{\text{pair}} (\mathbf{p} - \mathbf{q}).$$

Applying the Rayleigh quotient bounds for the positive definite matrix  $\Psi_{\text{pair}}$ , we obtain

$$c_\Omega \|\mathbf{p} - \mathbf{q}\|_2^2 \leq (\mathbf{p} - \mathbf{q})^\top \Psi_{\text{pair}} (\mathbf{p} - \mathbf{q}) \leq C_\Omega \|\mathbf{p} - \mathbf{q}\|_2^2,$$

and thus

$$\frac{c_\Omega}{2} \|\mathbf{p} - \mathbf{q}\|_2^2 \leq \mathbb{E}(\mathbf{p}) - \mathbb{E}(\mathbf{q}) - \nabla \mathbb{E}(\mathbf{q})^\top (\mathbf{p} - \mathbf{q}) \leq \frac{C_\Omega}{2} \|\mathbf{p} - \mathbf{q}\|_2^2.$$

This establishes that  $\mathbb{E}$  is  $c_\Omega$ -strongly convex and  $C_\Omega$ -smooth over the probability simplex  $\Delta$ , with constants determined by the minimal and maximal eigenvalues of  $\Psi_{\text{pair}}$ . □

*Note:* In any case even if  $\Psi_{\text{pair}}$  is non-psd, psd correction via Spectral Shifting can be utilised to make it a psd matrix.

## C.6 WEAK SUBMODULARITY OF SET FUNCTION $f$

**Theorem 1.** (Weak Submodularity) *The set function  $f(\tilde{A}) := \max_{\mathbf{p} \in \Delta_{n_1}^{\tilde{A}}; \text{supp}(\mathbf{p}) \subseteq \tilde{A}} \mathbb{E}(\mathbf{p})$  in Eq (4) is weakly submodular where  $\mathbb{E}(\mathbf{p}) = -\Psi_{\text{un}}^\top \mathbf{p} + \frac{1}{2} \mathbf{p}^\top \Psi_{\text{pair}} \mathbf{p}$  with the submodularity ratio  $\gamma > 0$ .*

*Proof.* Let  $L, S \subseteq [n_1]$  be disjoint sets and define  $m = |L| + |S|$ . Let  $\zeta(L) = \arg \max_{\mathbf{p} \in \Delta^{\mathbb{R}}, \text{supp}(\mathbf{p}) \subseteq L} \mathbb{E}(\mathbf{p})$  and similarly define  $\zeta(L \cup S)$  for the superset.

By the Restricted Strong Convexity (RSC) and Restricted Smoothness (RSM) of  $\mathbb{E}$  over the probability simplex (proved previously), we have for constants  $c_\Omega > 0$ ,  $C_\Omega > 0$ , and for any  $\mathbf{p}, \mathbf{q}$  supported in a set of size  $m$ ,

$$\frac{c_\Omega}{2} \|\mathbf{p} - \mathbf{q}\|_2^2 \leq \mathbb{E}(\mathbf{p}) - \mathbb{E}(\mathbf{q}) - \nabla \mathbb{E}(\mathbf{q})^\top (\mathbf{p} - \mathbf{q}) \leq \frac{C_\Omega}{2} \|\mathbf{p} - \mathbf{q}\|_2^2.$$

756 Let us upper bound the total gain from adding  $S$  to  $L$ :

$$757 \quad f(L \cup S) - f(L) = \mathbb{E}(\zeta(L \cup S)) - \mathbb{E}(\zeta(L)).$$

759 By the descent lemma and RSM,

$$760 \quad \mathbb{E}(\zeta(L \cup S)) - \mathbb{E}(\zeta(L)) \leq \langle \nabla \mathbb{E}(\zeta(L)), \zeta(L \cup S) - \zeta(L) \rangle - \frac{c_\Omega}{2} \|\zeta(L \cup S) - \zeta(L)\|^2.$$

763 We upper bound the inner product using the point  $\mathbf{v}$  defined as the projected optimal update within the support  $L \cup S$ . That is,

$$764 \quad v_{L \cup S} = \max \left\{ \frac{1}{c_\Omega} \nabla \mathbb{E}_{L \cup S}(\zeta(L)) + \zeta(L)_{L \cup S}, 0 \right\}.$$

767 Since  $\zeta(L \cup S)$  maximizes  $\mathbb{E}$  over support  $L \cup S$ , and  $\mathbf{v}$  is a feasible direction, we can use:

$$768 \quad \mathbb{E}(\zeta(L \cup S)) - \mathbb{E}(\zeta(L)) \leq \langle \nabla \mathbb{E}(\zeta(L)), \mathbf{v} - \zeta(L) \rangle - \frac{c_\Omega}{2} \|\mathbf{v} - \zeta(L)\|^2.$$

772 Now consider the coordinate-wise marginal gains. For each  $j \in S$ , we define the directional gain from adding  $j$  to  $L$  as:

$$773 \quad f(L \cup \{j\}) - f(L) \geq \max_{\alpha \geq 0} \left[ \langle \nabla_j \mathbb{E}(\zeta(L)), \alpha \rangle - \frac{L}{2} \alpha^2 \right] = \frac{1}{2C_\Omega} [\nabla_j \mathbb{E}(\zeta(L))]_+^2.$$

777 Summing over  $j \in S$  where  $\nabla_j \mathbb{E}(\zeta(L)) > 0$ , we get

$$778 \quad \sum_{j \in S} f(L \cup \{j\}) - f(L) \geq \frac{1}{2C_\Omega} \|\nabla_S^+ \mathbb{E}(\zeta(L))\|^2.$$

781 From the earlier upper bound, we had

$$782 \quad f(L \cup S) - f(L) \leq \langle \nabla \mathbb{E}(\zeta(L)), \mathbf{v} - \zeta(L) \rangle - \frac{c_\Omega}{2} \|\mathbf{v} - \zeta(L)\|^2.$$

785 The maximizer of this expression occurs at:

$$786 \quad v_j = \max \left\{ \frac{1}{c_\Omega} \nabla_j \mathbb{E}(\zeta(L)), 0 \right\}.$$

789 This gives:

$$790 \quad f(L \cup S) - f(L) \leq \frac{1}{2c_\Omega} \|\nabla_S^+ \mathbb{E}(\zeta(L))\|^2.$$

793 Combining the lower and upper bounds:

$$794 \quad \sum_{j \in S} f(L \cup \{j\}) - f(L) \geq \frac{1}{2C_\Omega} \|\nabla_S^+ \mathbb{E}(\zeta(L))\|^2, \quad f(L \cup S) - f(L) \leq \frac{1}{2c_\Omega} \|\nabla_S^+ \mathbb{E}(\zeta(L))\|^2.$$

797 Hence,

$$798 \quad \sum_{j \in S} f(L \cup \{j\}) - f(L) \geq \frac{c_\Omega}{C_\Omega} (f(L \cup S) - f(L)),$$

800 which proves weak submodularity with submodularity ratio  $\gamma = c_\Omega/C_\Omega > 0$ . □

801  
802  
803  
804  
805  
806  
807  
808  
809

## D COMPARATIVE ANALYSIS ACROSS VARIOUS NOTIONS OF TASK SIMILARITIES

### D.1 SIMILARITY ACROSS TASK VECTORS VIA LINEARIZED FINETUNING

Large-scale pretrained language models (PLMs) such as GPT-2 are widely adapted to downstream tasks via full-model fine-tuning. However, multi-task or per-task retraining remains computationally burdensome. *Task arithmetic* ? introduces a simple yet effective approach: given a pretrained checkpoint initialization  $\theta_0$  and task-specific fine-tuned weights  $\theta_t^*$ , the *task vector* is defined as:

$$\tau_t := \theta_t^* - \theta_0$$

These vectors enable model editing via linear composition:

- **Addition:**  $\theta_0 + \sum_{t \in T} \tau_t$  synthesizes multi-task behaviors.
- **Negation:**  $\theta_0 - \tau_s$  induces task-specific forgetting.

While effective, the underlying mechanisms behind this arithmetic remain poorly understood.

**Linearized Fine-Tuning:** (?) posit that *tangent-space fine-tuning* disentangles task behaviors more effectively by constraining updates to the local linear approximation of the model. Let  $f(x; \theta)$  denote a PLM with parameters  $\theta \in \mathbb{R}^m$ , the corresponding **nonlinear task vector** is given by  $\tau_t^{\text{nl}} := \theta_t^* - \theta_0$ .

In contrast, *linearized fine-tuning* restricts optimization to the first-order Taylor expansion:

$$f_{\text{lin}}(x; \theta) := f(x; \theta_0) + \nabla_{\theta} f(x; \theta_0)^{\top} (\theta - \theta_0)$$

This surrogate is optimized using Jacobian-vector products (JVP), yielding a linearized task vector:

$$\tau_t^{\text{lin}} := \theta_t^{\text{lin}*} - \theta_0$$

Task vectors are generally useful as they can enable model editing as well provide a well defined representation of the finetuning task at hand, dependent on the model parameters. Ideally, the goal would be to select multiple linearly independent task vectors such that they represent generalizably well across a range of IFT datasets and does generalizably well across different benchmark datasets. The algorithm is presented as Algorithm 1 in Section D.2.

### Similarity Structure of Task Embeddings

Directly computing any similarity metric over  $m \sim 10^6$  to  $10^9$  parameters, is computationally expensive. Thus, we first isolate the most informative layer (chosen via task-vector analysis using **layer-wise subsetting** and then project its high-dimensional slice task vector  $\tau \in \mathbb{R}^m$  to a much lower-dimensional vector  $\tilde{\tau} = R\tau \in \mathbb{R}^k$  using a **Gaussian random matrix**  $R \in \mathbb{R}^{k \times m}$  with  $k \ll m$ . This projection technique is known to preserve similarity distances in expectation, providing a reliable and efficient approximation for comparing vector directions in the reduced space.

**Cosine Similarity across Task Vectors:** To analyze inter-task relationships, we examine the cosine similarity between task vectors:

$$\text{sim}(\tau_A, \tau_B) := \frac{\tau_A^{\top} \tau_B}{\|\tau_A\|_2 \cdot \|\tau_B\|_2} \in [-1, 1]$$

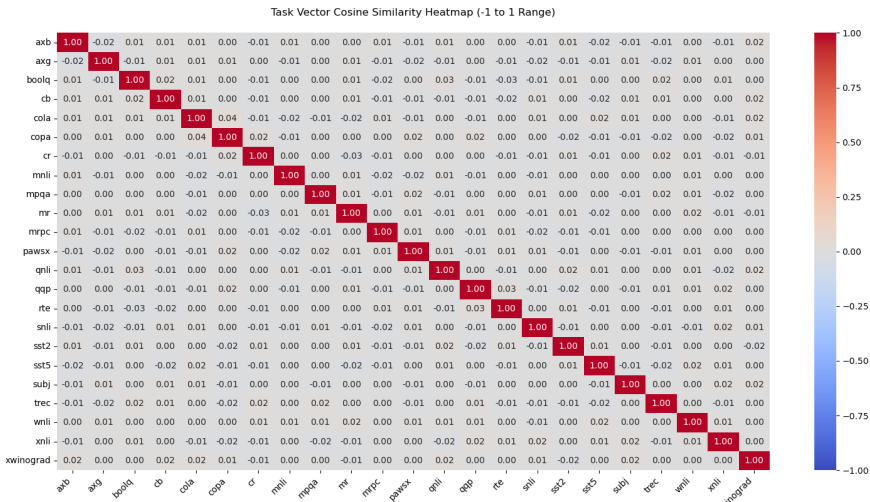
This metric probes the angular alignment between task-specific directions in parameter space. High similarity indicates shared representational updates; near-orthogonality suggests disentangled task pathways.

**Analyzing Task Vector Relationships via Cosine Similarity, PMI and JSD:** To analyze inter-task relationships, we work with Cosine Similarity, PMI, and JSD. While **Cosine Similarity** is a commonly used metric for comparing vector representations, it falls short in capturing nuanced differences in model behavior when applied to classification probability distributions. Cosine only measures the angular similarity between two vectors and is therefore invariant to vector magnitude. Hence, two models assigning vastly different probabilities but in the same proportional direction can still yield a high cosine score, misleadingly implying strong similarity. This limitation becomes evident in our experimental heatmap (Figure 2a), where task relationships are not clearly differentiated as many unrelated tasks appear spuriously similar due to their shared vector directionality. Moreover, cosine

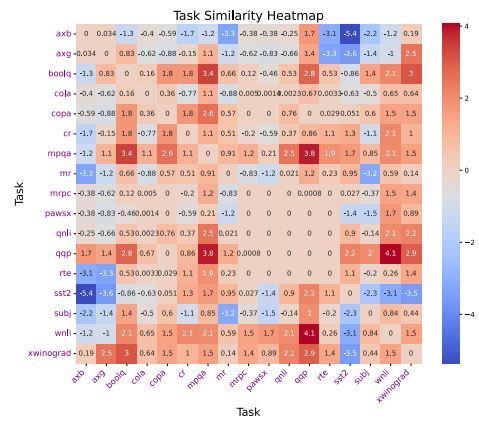
similarity does not adequately account for uncertainty or confidence in model outputs. To address these issues, we used Pointwise Mutual Information (PMI) and Jensen-Shannon Divergence (JSD), which offer better theoretical grounding and practical discriminability. As shown in Figures 2b and 2c, PMI captures directional alignment of model predictions with respect to task-specific specialization, while JSD provides a symmetric and robust comparison of output distributions. These metrics yield much more interpretable heatmaps where related tasks cluster more meaningfully and task-specific behaviors are more distinctly captured.

Concretely, the cosine heatmap appears overly uniform—masking important task groupings—whereas the PMI and JSD maps each expose clear blocks of high intra-group similarity and low inter-group coupling. These results confirm that, for fine-grained task-similarity assessment in large models, information-theoretic measures substantially outperform simple angular alignment.

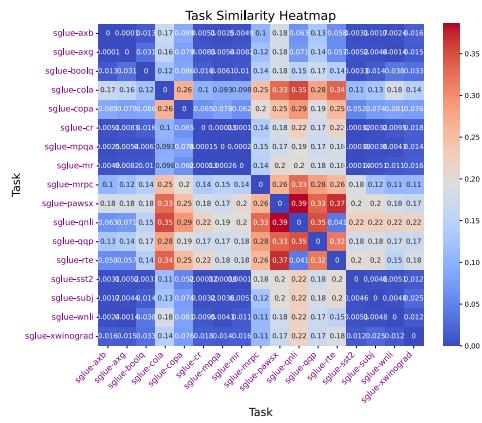
Below figure 2 visualizes the effects (on SGLUE tasks), comparing cosine, PMI, and JSD heatmaps to illustrate their differing sensitivity to inter-task relationships.



(a) Cosine Similarity



(b) PMI



## D.2 ALGORITHMS FOR COMPUTING PMI AND JSD

**Algorithm 1** performs fine-tuning by linearizing the model around its pretrained parameters. Instead of recomputing the full forward pass, it uses a Jacobian-vector product (JVP) to approximate the effect of parameter updates, allowing faster gradient-based updates in the “tangent space” of the original model.

**Algorithm 1: Linearized (Tangent-Space) Fine-Tuning**

**Require:** Pretrained weights  $\theta_0$ , dataset  $D_t$

```

1: Initialize  $\theta \leftarrow \theta_0$ 
2: while not converged do
3:   Sample mini-batch  $(x, y) \sim D_t$ 
4:   Compute base output  $o_0 = f(x; \theta_0)$ 
5:   Compute JVP:  $g = \text{JVP}(f(\cdot; \theta_0), \theta - \theta_0; x)$ 
6:    $\hat{o} = o_0 + g$ 
7:    $\theta \leftarrow \theta - \eta \nabla_{\theta} \ell(\hat{o}, y)$ 
8: end while
9: return  $\theta_t^{\text{lin*}}$ 

```

**Algorithm 2** quantifies how similarly two models  $M_A$  and  $M_B$  score the same labeled examples, using a pointwise mutual information (PMI)–inspired score. By averaging the log-ratio of predicted probabilities on each other’s held-out data, it produces a symmetric similarity score  $S_{AB}$ .

**Algorithm 2: PMI-Based Inter-Model Similarity  $S_{AB}$** 

**Require:** Models  $M_A, M_B$ ; datasets  $D^A, D^B$

**Ensure:** Similarity score  $S_{AB}$

```

1: Initialize accumulator  $\text{sum}_B \leftarrow 0$ 
2: for all  $(x, y) \in D^B$  do
3:   Compute  $p_A \leftarrow M_A(x)$  and extract  $p_A(y)$ 
4:   Compute  $p_B \leftarrow M_B(x)$  and extract  $p_B(y)$ 
5:   Update  $\text{sum}_B += \log\left(\frac{p_A(y)}{p_B(y)}\right)$ 
6: end for
7: Set  $\Delta_B \leftarrow \frac{1}{|D^B|} \cdot \text{sum}_B$ 
8: Initialize accumulator  $\text{sum}_A \leftarrow 0$ 
9: for all  $(x, y) \in D^A$  do
10:  Compute  $p_A \leftarrow M_A(x)$  and extract  $p_A(y)$ 
11:  Compute  $p_B \leftarrow M_B(x)$  and extract  $p_B(y)$ 
12:  Update  $\text{sum}_A += \log\left(\frac{p_B(y)}{p_A(y)}\right)$ 
13: end for
14: Set  $\Delta_A \leftarrow \frac{1}{|D^A|} \cdot \text{sum}_A$ 
15: return  $S_{AB} \leftarrow \frac{1}{2}(\Delta_A + \Delta_B)$ 

```

**Algorithm 3** computes the average Jensen–Shannon divergence between the predictive distributions of two models  $M_A$  and  $M_B$  across a shared dataset. Uses softmax outputs to measure how differently the models assign probabilities.

972  
973  
974  
975  
976  
977  
978  
979  
980  
981  
982  
983  
984  
985  
986  
987  
988  
989  
990  
991  
992  
993  
994  
995  
996  
997  
998  
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

### Algorithm 3: Jensen–Shannon Divergence (JSD) for Model Comparison

**Require:** Two models  $M_A, M_B$ ; dataset  $D$

**Ensure:** Average JSD value  $J\bar{S}D$

```

1: Initialize total_jsd  $\leftarrow 0$ 
2: for each input  $(x, y) \in D$  do
3:    $P \leftarrow \text{softmax}(M_A(x))$            {Predictive distribution from  $M_A$ }
4:    $Q \leftarrow \text{softmax}(M_B(x))$        {Predictive distribution from  $M_B$ }
5:    $M \leftarrow \frac{1}{2}(P + Q)$          {Mixture distribution}
6:    $KL_P \leftarrow \sum_i P_i \log\left(\frac{P_i}{M_i}\right)$ 
7:    $KL_Q \leftarrow \sum_i Q_i \log\left(\frac{Q_i}{M_i}\right)$ 
8:    $JSD(x) \leftarrow \frac{1}{2}(KL_P + KL_Q)$ 
9:   total_jsd  $\leftarrow$  total_jsd +  $JSD(x)$ 
10: end for
11: return  $J\bar{S}D \leftarrow \frac{\text{total\_jsd}}{|D|}$ 

```

## E EXPERIMENTAL DETAILS

All the experiments are conducted in a standardized and uniform environment to ensure reproducibility and cost-effectiveness. We finetune the models for one epoch on each dataset split, leveraging 8 NVIDIA H100 GPUs in bf16 precision. We use a per-device train batch size of 1, and using AdamW optimizer with a learning rate of  $2 \times 10^{-5}$ , weight decay 0.01, and gradient accumulation of 1. A linear learning-rate decay schedule is applied with a linear warmup over the first 3 % of total steps. To maximize memory efficiency, we enable gradient checkpointing and used DDP. The workloads are largely of 3 types, specifications and details of each are listed below.

**Fine-Tuning on Task Pool Datasets:** The objective of the approach is to find a final mixture from a large set of datasets which target different tasks. The pre-trained causal language model was used as the base model that was fine-tuned on each individual task. This stage follows the same configuration, with the following modification: models are finetuned for 3 epochs using an effective batch size of 64 and a cosine learning rate decay. A higher weight decay of 0.1 was applied, and all 8 GPUs were utilized in a Data Parallel setting. The goal is to train individual models on 316 distinct task drawn from diverse target sub-mixtures (T0, Flan2021, CoT, TULU, SGLue). All fine-tunings are full-parameter with no freezing or adapters.

**Similarity Matrix Computation :** We propose the use of two primary metrics, namely, 1) PMI and 2) JSD, although we arrive at the same by exhaustive experiments and analysis of other similarity measures and conclude with the efficacy of the two metrics. The **PMI** matrix computation, as illustrated in Algorithm 2 in Section D.2, is implemented similarly with optimizations at the PyTorch GPU and CPU multiprocessing level to speed up the computation of pairwise similarity scores due to the higher number of inferences required. We acquired the **JSD** matrix following the procedure outlined in Algorithm 3 in Section D.2. To optimize computation, we first precompute and store each model’s self-distribution ( $P_{X \rightarrow X}$ ) and cross-distribution ( $P_{X \rightarrow Y}$ ) across all tasks to prevent redundant forward passes. Distribution computation is vectorized by batching samples per task into single forward passes and all pairwise JSD values were calculated in parallel. A total of  $\frac{n(n-1)}{2}$  pairs were computed in both the cases, due to the inherent symmetric nature of the metric matrices, where  $n$  is the number of tasks.

**Fine-Tuning on Final Mixture :** This phase follows the same environment and base hyperparameters configuration described earlier, with modifications tailored to the final mixture evaluation. The mixture dataset acquired from the set of tasks using our proposed solution has to be evaluated against recognized benchmarks, for which the mixture dataset is used to fine-tune a Llama-2-7B model for a single epoch with an effective batch size of 8, a learning rate of  $2 \times 10^{-5}$  and gradient accumulation at every 8th step. A weight decay of 0.01 was used along with cosine learning rate decay and all 8 GPUs were utilized in a Data Parallel setting. Same hyperparameters and environment configuration was used when fine-tuning on Mistral-7B to showcase the relevance of the base model in the experimental results from our proposed mixture. We further explore mixture scale by evaluating training on subsets of varying sizes (25K, 50K and 100K) and examine performance sensitivity to batch size by comparing runs with effective batch size of 8. Additionally, for the 25K and 50K subsets, we conducted experiments with different values of  $\beta$  and  $\lambda$  to analyze their influence on mixture composition in both PMI-based and JSD-based submix selection strategies.

F ADDITIONAL RESULTS

Table 3: Llama-2-7b: Instruction-tuning performance on MMLU and Leaderboard subsets with  $\beta = 20, \lambda = 10$  using batch size 8.

Dataset		MMLU			Leaderboard			
Size	Method	BBH	GPQA	IFEval	Math	MMLU-Pro	MUSR	
<b>25K</b>								
25K	Random	0.3913 $\pm$ 0.0040	0.3482 $\pm$ 0.0059	0.2626 $\pm$ 0.0128	0.3729 $\pm$ N/A	0.0098 $\pm$ 0.0027	0.1877 $\pm$ 0.0036	0.3677 $\pm$ 0.0172
25K	Uniform	0.3479 $\pm$ 0.0039	0.3501 $\pm$ 0.0059	0.2701 $\pm$ 0.0129	0.3501 $\pm$ N/A	0.0151 $\pm$ 0.0034	0.1768 $\pm$ 0.0035	0.4127 $\pm$ 0.0175
25K	EPM	0.3802 $\pm$ 0.0040	0.3593 $\pm$ 0.0059	0.2601 $\pm$ 0.0127	0.3405 $\pm$ N/A	0.0151 $\pm$ 0.0033	0.1836 $\pm$ 0.0035	0.4286 $\pm$ 0.0177
25K	Ours (PMI)	0.4242 $\pm$ 0.0040	0.3598 $\pm$ 0.0059	0.2718 $\pm$ 0.0129	0.3561 $\pm$ N/A	0.0136 $\pm$ 0.0032	0.1877 $\pm$ 0.0036	0.4008 $\pm$ 0.0174
25K	Ours (JSD)	0.3926 $\pm$ 0.0040	0.3454 $\pm$ 0.0059	0.2785 $\pm$ 0.0130	0.3465 $\pm$ N/A	0.0151 $\pm$ 0.0034	0.1790 $\pm$ 0.0035	0.4021 $\pm$ 0.0175
<b>50K</b>								
50K	Random	0.4108 $\pm$ 0.0040	0.3565 $\pm$ 0.0060	0.2668 $\pm$ 0.0128	0.3681 $\pm$ N/A	0.0144 $\pm$ 0.0033	0.1881 $\pm$ 0.0036	0.3770 $\pm$ 0.0172
50K	Uniform	0.3725 $\pm$ 0.0040	0.3480 $\pm$ 0.0059	0.2785 $\pm$ 0.0130	0.4041 $\pm$ N/A	0.0181 $\pm$ 0.0037	0.1896 $\pm$ 0.0036	0.4206 $\pm$ 0.0176
50K	EPM	0.3801 $\pm$ 0.0040	0.3532 $\pm$ 0.0059	0.2634 $\pm$ 0.0128	0.3507 $\pm$ N/A	0.0128 $\pm$ 0.0031	0.1799 $\pm$ 0.0035	0.4206 $\pm$ 0.0176
50K	Ours (PMI)	0.4156 $\pm$ 0.0040	0.3619 $\pm$ 0.0060	0.2794 $\pm$ 0.0130	0.3417 $\pm$ N/A	0.0189 $\pm$ 0.0037	0.1856 $\pm$ 0.0035	0.3876 $\pm$ 0.0174
50K	Ours (JSD)	0.4074 $\pm$ 0.0040	0.3624 $\pm$ 0.0060	0.2802 $\pm$ 0.0130	0.3525 $\pm$ N/A	0.0098 $\pm$ 0.0027	0.1927 $\pm$ 0.0036	0.4206 $\pm$ 0.0176
<b>100K</b>								
100K	Random	0.3816 $\pm$ 0.0040	0.3458 $\pm$ 0.0059	0.2621 $\pm$ 0.0128	0.3705 $\pm$ N/A	0.0113 $\pm$ 0.0029	0.1893 $\pm$ 0.0036	0.4101 $\pm$ 0.0176
100K	Uniform	0.3953 $\pm$ 0.0040	0.3569 $\pm$ 0.0060	0.2710 $\pm$ 0.0129	0.3801 $\pm$ N/A	0.0189 $\pm$ 0.0037	0.1890 $\pm$ 0.0036	0.3730 $\pm$ 0.0172
100K	EPM	0.3915 $\pm$ 0.0040	0.3439 $\pm$ 0.0059	0.2844 $\pm$ 0.0131	0.3717 $\pm$ N/A	0.0098 $\pm$ 0.0027	0.1873 $\pm$ 0.0036	0.4259 $\pm$ 0.0176
100K	Ours (PMI)	0.4021 $\pm$ 0.0040	0.3633 $\pm$ 0.0059	0.2626 $\pm$ 0.0127	0.3525 $\pm$ N/A	0.0166 $\pm$ 0.0035	0.1894 $\pm$ 0.0035	0.3902 $\pm$ 0.0174
100K	Ours (JSD)	0.4256 $\pm$ 0.0040	0.3598 $\pm$ 0.0060	0.2894 $\pm$ 0.0131	0.3769 $\pm$ N/A	0.0273 $\pm$ 0.0035	0.1923 $\pm$ 0.0036	0.4101 $\pm$ 0.0176

Table 4: Mistral-7B: Instruction-tuning performance on MMLU and Leaderboard subsets with  $\beta = 20, \lambda = 10$  using batch size 8.

Dataset		MMLU			Leaderboard			
Size	Method	BBH	GPQA	IFEval	Math	MMLU-Pro	MUSR	
<b>25K</b>								
25K	Random	0.4539 $\pm$ 0.0041	0.3701 $\pm$ 0.0060	0.2760 $\pm$ 0.0130	0.4197 $\pm$ N/A	0.0120 $\pm$ 0.0031	0.1762 $\pm$ 0.0035	0.4101 $\pm$ 0.0175
25K	Uniform	0.4376 $\pm$ 0.0041	0.3628 $\pm$ 0.0060	0.2601 $\pm$ 0.0127	0.4029 $\pm$ N/A	0.0159 $\pm$ 0.0034	0.1735 $\pm$ 0.0035	0.4458 $\pm$ 0.0178
25K	EPM	0.4364 $\pm$ 0.0041	0.3355 $\pm$ 0.0060	0.2869 $\pm$ 0.0131	0.4281 $\pm$ N/A	0.0121 $\pm$ 0.0030	0.1498 $\pm$ 0.0036	0.3492 $\pm$ 0.0168
25K	Ours (PMI)	0.3903 $\pm$ 0.0040	0.3244 $\pm$ 0.0058	0.2626 $\pm$ 0.0128	0.3297 $\pm$ N/A	0.0128 $\pm$ 0.0031	0.1503 $\pm$ 0.0033	0.3836 $\pm$ 0.0174
25K	Ours (JSD)	0.3783 $\pm$ 0.0040	0.3420 $\pm$ 0.0060	0.2878 $\pm$ 0.0131	0.3525 $\pm$ N/A	0.0129 $\pm$ 0.0031	0.1742 $\pm$ 0.0034	0.3929 $\pm$ 0.0174
<b>50K</b>								
50K	Random	0.4177 $\pm$ 0.0040	0.3446 $\pm$ 0.0059	0.2659 $\pm$ 0.0128	0.4113 $\pm$ N/A	0.0106 $\pm$ 0.0028	0.1733 $\pm$ 0.0035	0.3836 $\pm$ 0.0175
50K	Uniform	0.4452 $\pm$ 0.0041	0.3479 $\pm$ 0.0059	0.2651 $\pm$ 0.0128	0.4161 $\pm$ N/A	0.0151 $\pm$ 0.0033	0.1799 $\pm$ 0.0035	0.3823 $\pm$ 0.0172
50K	EPM	0.4405 $\pm$ 0.0041	0.3413 $\pm$ 0.0059	0.2701 $\pm$ 0.0129	0.4293 $\pm$ N/A	0.0174 $\pm$ 0.0036	0.1871 $\pm$ 0.0036	0.4034 $\pm$ 0.0174
50K	Ours (PMI)	0.4228 $\pm$ 0.0040	0.3492 $\pm$ 0.0058	0.2735 $\pm$ 0.0129	0.3094 $\pm$ N/A	0.0174 $\pm$ 0.0036	0.1758 $\pm$ 0.0035	0.4259 $\pm$ 0.0176
50K	Ours (JSD)	0.4138 $\pm$ 0.0040	0.3498 $\pm$ 0.0059	0.2567 $\pm$ 0.0127	0.4065 $\pm$ N/A	0.0159 $\pm$ 0.0034	0.1898 $\pm$ 0.0035	0.3890 $\pm$ 0.0173
<b>100K</b>								
100K	Random	0.4476 $\pm$ 0.0041	0.3416 $\pm$ 0.0060	0.2542 $\pm$ 0.0126	0.4388 $\pm$ N/A	0.0186 $\pm$ 0.0038	0.1730 $\pm$ 0.0034	0.4048 $\pm$ 0.0175
100K	Uniform	0.4486 $\pm$ 0.0041	0.3532 $\pm$ 0.0059	0.2661 $\pm$ 0.0128	0.3741 $\pm$ N/A	0.0174 $\pm$ 0.0036	0.1724 $\pm$ 0.0034	0.3810 $\pm$ 0.0173
100K	EPM	0.4505 $\pm$ 0.0041	0.3578 $\pm$ 0.0060	0.2466 $\pm$ 0.0125	0.4388 $\pm$ N/A	0.0174 $\pm$ 0.0036	0.1859 $\pm$ 0.0035	0.4074 $\pm$ 0.0175
100K	Ours (PMI)	0.5476 $\pm$ 0.0040	0.3388 $\pm$ 0.0058	0.2508 $\pm$ 0.0126	0.3369 $\pm$ N/A	0.0136 $\pm$ 0.0032	0.1810 $\pm$ 0.0035	0.4081 $\pm$ 0.0176
100K	Ours (JSD)	0.5301 $\pm$ 0.0040	0.3591 $\pm$ 0.0060	0.2667 $\pm$ 0.0127	0.4137 $\pm$ N/A	0.0189 $\pm$ 0.0037	0.1953 $\pm$ 0.0035	0.4140 $\pm$ 0.0175

  highest accuracy    
   2nd highest accuracy    
   3rd highest accuracy.

We observe that for **LLaMA** as the base model, increasing the number of instances in the mixture has negligible impact on performance when using  $\beta = 20$  and  $\lambda = 10$ . However, in the case of **Mistral**, the same configuration leads to a substantial improvement, where our **PMI-based method** yields **at least 10% higher accuracy on MMLU** compared to heuristic-driven methods. This strongly indicates that **PMI scales more effectively** with larger mixtures, leveraging the increased data volume to improve instruction tuning performance.

On the other hand, methods that rely on **heuristics** tend to perform better with **smaller instance sizes**. The reduced size helps **control the randomness** in mixture construction, suggesting that such heuristic approaches **do not scale well** as the number of instances increases. This confirms that their design may lack robustness in high-complexity or large-scale scenarios, where principled methods like PMI show a clear advantage.

Table 5: Llama-2-7b: Instruction-tuning performance on MMLU and Leaderboard subsets with varying  $\beta=20$ ,  $\lambda=10$  using batch size 64.

Dataset		MMLU		Leaderboard				
Size	Method	BBH	GPQA	IFEval	Math	MMLU-Pro	MUSR	
<b>25K</b>								
25K	Random	<b>0.4004</b> $\pm 0.0040$	<b>0.3602</b> $\pm 0.0059$	<b>0.2928</b> $\pm 0.0132$	0.3357 $\pm N/A$	<b>0.0166</b> $\pm 0.0035$	<b>0.1924</b> $\pm 0.0036$	0.3889 $\pm 0.0173$
25K	Uniform	0.3987 $\pm 0.0040$	0.3525 $\pm 0.0059$	0.2710 $\pm 0.0129$	0.3441 $\pm N/A$	0.0159 $\pm 0.0034$	0.1832 $\pm 0.0035$	<b>0.4220</b> $\pm 0.0176$
25K	EPM	0.3970 $\pm 0.0040$	0.3468 $\pm 0.0059$	0.2685 $\pm 0.0128$	<b>0.3681</b> $\pm N/A$	0.0128 $\pm 0.0031$	0.1853 $\pm 0.0035$	0.4140 $\pm 0.0176$
25K	Ours (PMI)	0.3917 $\pm 0.0040$	0.3479 $\pm 0.0059$	0.2676 $\pm 0.0128$	0.3477 $\pm N/A$	0.0106 $\pm 0.0028$	0.1918 $\pm 0.0036$	0.3889 $\pm 0.0174$
25K	Ours (JSD)	0.4057 $\pm 0.0040$	0.3517 $\pm 0.0059$	0.2886 $\pm 0.0131$	0.3273 $\pm N/A$	0.0121 $\pm 0.0030$	0.1849 $\pm 0.0035$	0.3995 $\pm 0.0175$
<b>50K</b>								
50K	Random	0.3761 $\pm 0.0040$	0.3515 $\pm 0.0059$	0.2810 $\pm 0.0130$	0.3549 $\pm N/A$	0.0113 $\pm 0.0029$	0.1845 $\pm 0.0035$	0.3796 $\pm 0.0172$
50K	Uniform	0.3923 $\pm 0.0040$	<b>0.3612</b> $\pm 0.0059$	0.2710 $\pm 0.0129$	0.3693 $\pm N/A$	0.0121 $\pm 0.0030$	0.1875 $\pm 0.0036$	0.4206 $\pm 0.0177$
50K	EPM	<b>0.4029</b> $\pm 0.0040$	0.3461 $\pm 0.0059$	0.2710 $\pm 0.0129$	<b>0.3885</b> $\pm N/A$	0.0151 $\pm 0.0034$	0.1869 $\pm 0.0036$	0.4325 $\pm 0.0177$
50K	Ours (PMI)	0.3748 $\pm 0.0040$	0.3562 $\pm 0.0059$	<b>0.2878</b> $\pm 0.0131$	0.3441 $\pm N/A$	<b>0.0159</b> $\pm 0.0034$	0.1896 $\pm 0.0036$	0.3929 $\pm 0.0174$
50K	Ours (JSD)	0.3758 $\pm 0.0040$	0.3543 $\pm 0.0060$	0.2676 $\pm 0.0128$	0.3741 $\pm N/A$	0.0136 $\pm 0.0032$	<b>0.1902</b> $\pm 0.0036$	<b>0.4220</b> $\pm 0.0176$
<b>100K</b>								
100K	Random	0.3816 $\pm 0.0040$	0.3458 $\pm 0.0059$	0.2651 $\pm 0.0128$	0.3705 $\pm N/A$	0.0113 $\pm 0.0029$	0.1893 $\pm 0.0036$	0.4101 $\pm 0.0176$
100K	Uniform	0.3953 $\pm 0.0040$	0.3569 $\pm 0.0060$	0.2710 $\pm 0.0129$	<b>0.3801</b> $\pm N/A$	<b>0.0189</b> $\pm 0.0037$	0.1890 $\pm 0.0036$	0.3730 $\pm 0.0172$
100K	EPM	0.3915 $\pm 0.0040$	0.3439 $\pm 0.0059$	<b>0.2844</b> $\pm 0.0131$	0.3717 $\pm N/A$	0.0098 $\pm 0.0027$	0.1873 $\pm 0.0036$	<b>0.4259</b> $\pm 0.0176$
100K	Ours (PMI)	0.4017 $\pm 0.0040$	0.3591 $\pm 0.0060$	0.2827 $\pm 0.0131$	0.3213 $\pm N/A$	0.0166 $\pm 0.0035$	0.1854 $\pm 0.0035$	0.4021 $\pm 0.0175$
100K	Ours (JSD)	<b>0.4165</b> $\pm 0.0040$	<b>0.3609</b> $\pm 0.0060$	0.2827 $\pm 0.0131$	0.3585 $\pm N/A$	0.0151 $\pm 0.0034$	<b>0.1893</b> $\pm 0.0036$	0.3981 $\pm 0.0176$

Table 6: Mistral-7B: Instruction-tuning performance on MMLU and Leaderboard subsets with varying  $\beta = 20$ ,  $\lambda = 10$  using batch size 64.

Dataset		MMLU		Leaderboard				
Size	Method	BBH	GPQA	IFEval	Math	MMLU-Pro	MUSR	
<b>25K</b>								
25K	Random	0.5541 $\pm 0.0040$	<b>0.4227</b> $\pm 0.0061$	0.2903 $\pm 0.0132$	0.4544 $\pm N/A$	0.0227 $\pm 0.0041$	<b>0.2578</b> $\pm 0.0040$	<b>0.4484</b> $\pm 0.0179$
25K	Uniform	<b>0.5600</b> $\pm 0.0040$	0.4055 $\pm 0.0061$	0.2685 $\pm 0.0128$	<b>0.4592</b> $\pm N/A$	0.0242 $\pm 0.0042$	0.2557 $\pm 0.0040$	0.4259 $\pm 0.0176$
25K	EPM	0.5449 $\pm 0.0040$	0.4152 $\pm 0.0062$	0.2735 $\pm 0.0129$	0.4376 $\pm N/A$	0.0219 $\pm 0.0040$	0.2345 $\pm 0.0039$	0.4418 $\pm 0.0178$
25K	Ours (PMI)	0.5383 $\pm 0.0040$	0.4171 $\pm 0.0062$	0.2626 $\pm 0.0128$	0.3921 $\pm N/A$	0.0211 $\pm 0.0039$	0.2485 $\pm 0.0039$	0.3876 $\pm 0.0175$
25K	Ours (JSD)	0.5400 $\pm 0.0040$	0.4180 $\pm 0.0062$	<b>0.2928</b> $\pm 0.0132$	0.4544 $\pm N/A$	<b>0.0264</b> $\pm 0.0044$	0.2462 $\pm 0.0039$	0.4180 $\pm 0.0178$
<b>50K</b>								
50K	Random	0.5524 $\pm 0.0040$	0.4044 $\pm 0.0061$	<b>0.2878</b> $\pm 0.0131$	0.4844 $\pm N/A$	<b>0.0272</b> $\pm 0.0045$	0.2620 $\pm 0.0040$	0.3995 $\pm 0.0174$
50K	Uniform	<b>0.5585</b> $\pm 0.0040$	0.4062 $\pm 0.0061$	0.2727 $\pm 0.0129$	<b>0.4940</b> $\pm N/A$	0.0257 $\pm 0.0043$	<b>0.2702</b> $\pm 0.0040$	0.4272 $\pm 0.0178$
50K	EPM	0.5541 $\pm 0.0040$	0.4294 $\pm 0.0062$	0.2861 $\pm 0.0131$	0.4676 $\pm N/A$	0.0189 $\pm 0.0037$	0.2612 $\pm 0.0040$	<b>0.4458</b> $\pm 0.0179$
50K	Ours (PMI)	0.5499 $\pm 0.0040$	<b>0.4135</b> $\pm 0.0061$	0.2861 $\pm 0.0131$	0.3825 $\pm N/A$	0.0181 $\pm 0.0037$	0.2479 $\pm 0.0039$	0.4378 $\pm 0.0178$
50K	Ours (JSD)	0.5389 $\pm 0.0040$	0.3543 $\pm 0.0060$	0.2676 $\pm 0.0128$	0.3741 $\pm N/A$	0.0136 $\pm 0.0032$	0.1902 $\pm 0.0036$	0.4220 $\pm 0.0176$
<b>100K</b>								
100K	Random	0.4476 $\pm 0.0041$	0.3416 $\pm 0.0060$	0.2542 $\pm 0.0126$	<b>0.4388</b> $\pm N/A$	0.0196 $\pm 0.0038$	0.1730 $\pm 0.0034$	0.4048 $\pm 0.0175$
100K	Uniform	0.4486 $\pm 0.0041$	0.3532 $\pm 0.0059$	0.2668 $\pm 0.0128$	0.3741 $\pm N/A$	0.0174 $\pm 0.0036$	0.1724 $\pm 0.0034$	0.3810 $\pm 0.0173$
100K	EPM	0.4505 $\pm 0.0041$	0.3578 $\pm 0.0060$	0.2466 $\pm 0.0125$	<b>0.4388</b> $\pm N/A$	0.0174 $\pm 0.0036$	0.1859 $\pm 0.0035$	0.4074 $\pm 0.0175$
100K	Ours (PMI)	<b>0.5476</b> $\pm 0.0040$	<b>0.4161</b> $\pm 0.0062$	0.2701 $\pm 0.0129$	0.3501 $\pm N/A$	<b>0.0234</b> $\pm 0.0041$	<b>0.2558</b> $\pm 0.0040$	0.4101 $\pm 0.0176$
100K	Ours (JSD)	0.5301 $\pm 0.0040$	0.3784 $\pm 0.0059$	<b>0.2768</b> $\pm 0.0130$	0.4257 $\pm N/A$	0.0204 $\pm 0.0039$	0.2342 $\pm 0.0039$	<b>0.4206</b> $\pm 0.0176$

We demonstrate that a **small adjustment in batch size**-specifically increasing it to **64**-in conjunction with the use of **Mistral**, allows us to achieve performance that is **comparable to a 100K instance mixture** trained with **BS=8**, while using only a **25K instance mixture**. This setup delivers a **7% boost in performance on BBH and MMLU-Pro**, thereby validating the **efficacy of our mixture strategy**. These results suggest that, when provided with the right computational environment, our mixture formulation has the **potential to match or surpass** much larger-scale setups on major benchmarks. Furthermore, our **JSD-based mixture** shows a remarkable **13% improvement over its LLaMA variant** when deployed with Mistral and BS=64. This emphasizes the importance of **careful hyperparameter tuning** in fully realizing the benefits of the proposed mixtures.

We also observe a **consistent gain of 5–13%** across several **leaderboard benchmarks**, including **BBH, IFEval, and Math**, when the instance size is scaled from **25K to 50K** using Mistral. However, the same scaling yields only a modest **1–2% improvement with LLaMA**. Notably, increasing the instance size to **100K results in negligible performance gains** across most benchmarks for both Mistral and LLaMA, suggesting a possible **diminishing return** beyond a certain mixture size threshold.

Table 7: Llama-2-7b: Instruction-tuning performance on MMLU and Leaderboard subsets with varying  $\beta$  and  $\lambda$  using batch size 8.

Dataset Size(Method)	MMLU		Leaderboard				
	BBH	GPQA	IFEval	Math	MMLU-Pro	MUSR	
<b>25K (PMI)</b>							
$\beta=14954; \lambda=263$	0.4098 $\pm$ 0.0040	0.3637 $\pm$ 0.0059	0.2685 $\pm$ 0.0128	0.3405 $\pm$ N/A	0.0159 $\pm$ 0.0034	0.1869 $\pm$ 0.0036	0.3849 $\pm$ 0.0173
$\beta=5273; \lambda=195$	0.4045 $\pm$ 0.0040	0.3536 $\pm$ 0.0059	0.2718 $\pm$ 0.0129	0.3609 $\pm$ N/A	<b>0.0166</b> $\pm$ 0.0035	0.1823 $\pm$ 0.0035	0.4021 $\pm$ 0.0174
$\beta=2535; \lambda=196$	<b>0.4258</b> $\pm$ 0.0040	<b>0.3659</b> $\pm$ 0.0059	0.2701 $\pm$ 0.0129	0.3357 $\pm$ N/A	0.0128 $\pm$ 0.0031	<b>0.1890</b> $\pm$ 0.0036	0.3929 $\pm$ 0.0173
$\beta=307; \lambda=60$	0.3977 $\pm$ 0.0040	0.3605 $\pm$ 0.0059	<b>0.2735</b> $\pm$ 0.0129	<b>0.3681</b> $\pm$ N/A	0.0159 $\pm$ 0.0034	0.1881 $\pm$ 0.0036	<b>0.4074</b> $\pm$ 0.0174
$\beta=19; \lambda=5$	0.3827 $\pm$ 0.0040	0.3576 $\pm$ 0.0059	0.2693 $\pm$ 0.0129	0.3381 $\pm$ N/A	0.0121 $\pm$ 0.0030	0.1872 $\pm$ 0.0036	<b>0.4246</b> $\pm$ 0.0177
<b>25K (JSD)</b>							
$\beta=14954; \lambda=263$	0.3929 $\pm$ 0.0040	0.3486 $\pm$ 0.0059	0.2626 $\pm$ 0.0128	0.3357 $\pm$ N/A	<b>0.0189</b> $\pm$ 0.0037	0.1828 $\pm$ 0.0035	0.3995 $\pm$ 0.0176
$\beta=5273; \lambda=195$	0.3793 $\pm$ 0.0040	<b>0.3614</b> $\pm$ 0.0059	0.2693 $\pm$ 0.0129	<b>0.3657</b> $\pm$ N/A	0.0166 $\pm$ 0.0035	0.1769 $\pm$ 0.0035	<b>0.3981</b> $\pm$ 0.0174
$\beta=2535; \lambda=196$	0.3978 $\pm$ 0.0040	0.3574 $\pm$ 0.0059	<b>0.2794</b> $\pm$ 0.0130	0.3573 $\pm$ N/A	0.0113 $\pm$ 0.0029	0.1844 $\pm$ 0.0035	<b>0.4048</b> $\pm$ 0.0175
$\beta=307; \lambda=60$	<b>0.4188</b> $\pm$ 0.0040	0.3522 $\pm$ 0.0060	<b>0.2794</b> $\pm$ 0.0130	0.3453 $\pm$ N/A	0.0144 $\pm$ 0.0033	<b>0.1913</b> $\pm$ 0.0036	0.4021 $\pm$ 0.0176
$\beta=19; \lambda=5$	0.3890 $\pm$ 0.0040	0.3545 $\pm$ 0.0060	0.2743 $\pm$ 0.0129	0.3525 $\pm$ N/A	0.0144 $\pm$ 0.0033	0.1823 $\pm$ 0.0035	0.3942 $\pm$ 0.0174
<b>50K (PMI)</b>							
$\beta=14954; \lambda=263$	0.3731 $\pm$ 0.0040	0.3527 $\pm$ 0.0059	0.2903 $\pm$ 0.0132	<b>0.3489</b> $\pm$ N/A	<b>0.0166</b> $\pm$ 0.0035	0.1864 $\pm$ 0.0036	0.3968 $\pm$ 0.0174
$\beta=5273; \lambda=195$	<b>0.4031</b> $\pm$ 0.0040	0.3609 $\pm$ 0.0059	0.2836 $\pm$ 0.0131	0.3429 $\pm$ N/A	0.0159 $\pm$ 0.0034	0.1869 $\pm$ 0.0036	0.4048 $\pm$ 0.0174
$\beta=2535; \lambda=196$	0.4164 $\pm$ 0.0040	0.3567 $\pm$ 0.0060	0.2735 $\pm$ 0.0129	0.3537 $\pm$ N/A	0.0151 $\pm$ 0.0034	0.1906 $\pm$ 0.0036	0.4127 $\pm$ 0.0174
$\beta=307; \lambda=60$	0.3919 $\pm$ 0.0040	0.3637 $\pm$ 0.0059	0.2743 $\pm$ 0.0129	<b>0.3489</b> $\pm$ N/A	0.0091 $\pm$ 0.0026	0.1797 $\pm$ 0.0035	<b>0.4140</b> $\pm$ 0.0176
$\beta=19; \lambda=5$	0.4004 $\pm$ 0.0040	<b>0.3690</b> $\pm$ 0.0060	<b>0.2936</b> $\pm$ 0.0132	0.3393 $\pm$ N/A	0.0159 $\pm$ 0.0034	<b>0.1921</b> $\pm$ 0.0036	0.3862 $\pm$ 0.0174
<b>50K (JSD)</b>							
$\beta=14954; \lambda=263$	0.4212 $\pm$ 0.0040	0.3536 $\pm$ 0.0059	0.2659 $\pm$ 0.0128	0.3513 $\pm$ N/A	0.0151 $\pm$ 0.0034	<b>0.1902</b> $\pm$ 0.0036	<b>0.4233</b> $\pm$ 0.0176
$\beta=5273; \lambda=195$	<b>0.4219</b> $\pm$ 0.0040	0.3584 $\pm$ 0.0060	0.2659 $\pm$ 0.0128	0.3561 $\pm$ N/A	<b>0.0204</b> $\pm$ 0.0039	0.1877 $\pm$ 0.0036	0.3929 $\pm$ 0.0175
$\beta=2535; \lambda=196$	0.4205 $\pm$ 0.0040	0.3545 $\pm$ 0.0060	<b>0.2810</b> $\pm$ 0.0130	0.3513 $\pm$ N/A	0.0121 $\pm$ 0.0030	0.1895 $\pm$ 0.0036	0.4074 $\pm$ 0.0176
$\beta=307; \lambda=60$	0.4025 $\pm$ 0.0040	0.3600 $\pm$ 0.0059	0.2643 $\pm$ 0.0128	0.3585 $\pm$ N/A	0.0144 $\pm$ 0.0033	0.1813 $\pm$ 0.0035	0.3823 $\pm$ 0.0174
$\beta=19; \lambda=5$	0.4039 $\pm$ 0.0040	<b>0.3650</b> $\pm$ 0.0060	0.2668 $\pm$ 0.0128	<b>0.3561</b> $\pm$ N/A	0.0166 $\pm$ 0.0035	0.1846 $\pm$ 0.0035	0.4074 $\pm$ 0.0175

Table 8: Mistral-7b: Instruction-tuning performance on MMLU and Leaderboard subsets with varying  $\beta$  and  $\lambda$  using batch size 8.

Dataset Size(Method)	MMLU		Leaderboard				
	BBH	GPQA	IFEval	Math	MMLU-Pro	MUSR	
<b>25K (PMI)</b>							
$\beta=14954; \lambda=263$	<b>0.4582</b> $\pm$ 0.0041	0.3579 $\pm$ 0.0059	0.2685 $\pm$ 0.0128	0.3237 $\pm$ N/A	0.0144 $\pm$ 0.0033	0.1893 $\pm$ 0.0036	0.3981 $\pm$ 0.0175
$\beta=5273; \lambda=195$	0.4368 $\pm$ 0.0040	0.3498 $\pm$ 0.0059	0.2592 $\pm$ 0.0127	0.3621 $\pm$ N/A	0.0136 $\pm$ 0.0032	0.1863 $\pm$ 0.0035	<b>0.4696</b> $\pm$ 0.0179
$\beta=2535; \lambda=196$	0.4221 $\pm$ 0.0040	0.3579 $\pm$ 0.0060	0.2525 $\pm$ 0.0126	0.2842 $\pm$ N/A	0.0136 $\pm$ 0.0032	<b>0.1912</b> $\pm$ 0.0036	0.3677 $\pm$ 0.0172
$\beta=307; \lambda=60$	0.4568 $\pm$ 0.0041	<b>0.3612</b> $\pm$ 0.0059	0.2584 $\pm$ 0.0127	0.3177 $\pm$ N/A	0.0128 $\pm$ 0.0031	0.1877 $\pm$ 0.0036	0.4127 $\pm$ 0.0176
$\beta=19; \lambda=5$	0.4262 $\pm$ 0.0041	0.3532 $\pm$ 0.0060	<b>0.2693</b> $\pm$ 0.0129	<b>0.3669</b> $\pm$ N/A	<b>0.0151</b> $\pm$ 0.0033	0.1841 $\pm$ 0.0035	0.4114 $\pm$ 0.0175
<b>25K (JSD)</b>							
$\beta=14954; \lambda=263$	0.4327 $\pm$ 0.0040	0.3909 $\pm$ 0.0061	0.2601 $\pm$ 0.0127	0.4077 $\pm$ N/A	0.0128 $\pm$ 0.0031	0.1796 $\pm$ 0.0035	0.4378 $\pm$ 0.0177
$\beta=5273; \lambda=195$	0.4313 $\pm$ 0.0041	0.3637 $\pm$ 0.0060	0.2617 $\pm$ 0.0127	0.3453 $\pm$ N/A	0.0166 $\pm$ 0.0035	0.1784 $\pm$ 0.0035	0.4220 $\pm$ 0.0176
$\beta=2535; \lambda=196$	0.4453 $\pm$ 0.0041	<b>0.3706</b> $\pm$ 0.0060	0.2601 $\pm$ 0.0127	<b>0.3969</b> $\pm$ N/A	<b>0.0128</b> $\pm$ 0.0031	0.1728 $\pm$ 0.0034	0.4220 $\pm$ 0.0176
$\beta=307; \lambda=60$	<b>0.4568</b> $\pm$ 0.0041	0.3604 $\pm$ 0.0060	<b>0.2810</b> $\pm$ 0.0130	0.3933 $\pm$ N/A	0.0181 $\pm$ 0.0037	<b>0.1762</b> $\pm$ 0.0035	0.4259 $\pm$ 0.0173
$\beta=19; \lambda=5$	0.3957 $\pm$ 0.0040	0.3581 $\pm$ 0.0059	0.2450 $\pm$ 0.0125	0.4137 $\pm$ N/A	0.0136 $\pm$ 0.0032	0.1615 $\pm$ 0.0034	<b>0.4418</b> $\pm$ 0.0176
<b>50K (PMI)</b>							
$\beta=14954; \lambda=263$	0.4325 $\pm$ 0.0041	0.3340 $\pm$ 0.0058	0.2668 $\pm$ 0.0128	0.3070 $\pm$ N/A	0.0166 $\pm$ 0.0035	0.1902 $\pm$ 0.0036	0.3968 $\pm$ 0.0176
$\beta=5273; \lambda=195$	<b>0.4379</b> $\pm$ 0.0041	0.3539 $\pm$ 0.0059	<b>0.2886</b> $\pm$ 0.0131	<b>0.3705</b> $\pm$ N/A	<b>0.0181</b> $\pm$ 0.0037	<b>0.1911</b> $\pm$ 0.0036	0.3717 $\pm$ 0.0170
$\beta=2535; \lambda=196$	0.4009 $\pm$ 0.0040	<b>0.3581</b> $\pm$ 0.0060	0.2819 $\pm$ 0.0130	0.3501 $\pm$ N/A	0.0174 $\pm$ 0.0036	0.1661 $\pm$ 0.0034	0.4127 $\pm$ 0.0175
$\beta=307; \lambda=60$	0.4336 $\pm$ 0.0041	0.3508 $\pm$ 0.0059	0.2785 $\pm$ 0.0130	0.3177 $\pm$ N/A	0.0136 $\pm$ 0.0032	0.1818 $\pm$ 0.0035	0.3810 $\pm$ 0.0171
$\beta=19; \lambda=5$	0.4134 $\pm$ 0.0040	0.3520 $\pm$ 0.0059	0.2601 $\pm$ 0.0127	0.3777 $\pm$ N/A	0.0128 $\pm$ 0.0031	0.1661 $\pm$ 0.0034	<b>0.4193</b> $\pm$ 0.0175
<b>50K (JSD)</b>							
$\beta=14954; \lambda=263$	0.4295 $\pm$ 0.0041	<b>0.3631</b> $\pm$ 0.0060	0.2592 $\pm$ 0.0127	0.4233 $\pm$ N/A	0.0196 $\pm$ 0.0038	0.1687 $\pm$ 0.0034	0.3929 $\pm$ 0.0173
$\beta=5273; \lambda=195$	0.4372 $\pm$ 0.0041	0.3623 $\pm$ 0.0060	0.2727 $\pm$ 0.0129	0.4233 $\pm$ N/A	0.0159 $\pm$ 0.0034	<b>0.1768</b> $\pm$ 0.0035	0.4418 $\pm$ 0.0177
$\beta=2535; \lambda=196$	0.4350 $\pm$ 0.0041	0.3505 $\pm$ 0.0059	0.2617 $\pm$ 0.0127	<b>0.4424</b> $\pm$ N/A	0.0219 $\pm$ 0.0040	0.1669 $\pm$ 0.0034	0.3836 $\pm$ 0.0172
$\beta=307; \lambda=60$	<b>0.4400</b> $\pm$ 0.0041	0.3373 $\pm$ 0.0059	<b>0.2735</b> $\pm$ 0.0129	0.4137 $\pm$ N/A	<b>0.0227</b> $\pm$ 0.0041	0.1750 $\pm$ 0.0035	0.3717 $\pm$ 0.0173
$\beta=19; \lambda=5$	0.4285 $\pm$ 0.0041	0.3444 $\pm$ 0.0058	0.2424 $\pm$ 0.0124	0.4257 $\pm$ N/A	0.0106 $\pm$ 0.0028	0.1743 $\pm$ 0.0035	<b>0.4484</b> $\pm$ 0.0179

We observe a notable improvement in convergence for Mistral over LLaMA, reflected in a consistent **2–5% boost in benchmark performance**. This underscores Mistral’s enhanced compatibility with our mixture strategies.

1242 Among the evaluated configurations, the JSD-based mixture with  $\beta = 307$  and  $\lambda = 60$  emerges as  
1243 **the most reliable**, frequently achieving either the best or near-best results across a diverse range of  
1244 datasets and evaluation metrics.

1245 Our analysis also reveals that **PMI and JSD excel in distinct areas**. While **JSD outperforms in**  
1246 **leaderboard subsets**—notably on **IFEval and Math**—the **PMI method leads on MMLU tasks**,  
1247 demonstrating that each method has specialized strengths.

1248 Interestingly, we find that **leaderboard metrics benefit from larger instance mixtures**, whereas  
1249 **MMLU-related tasks such as BBH and GPQA plateau or even degrade** in performance when too  
1250 many instances are included. This may be due to overfitting to harder instances or increased noise  
1251 from larger mixtures.

1252 We also identify that a **balanced ratio of  $\frac{\beta}{\lambda}$** , such as  $\beta = 307$ ,  $\lambda = 60$ , tends to **consistently**  
1253 **outperform** other configurations. In contrast, **higher ratios** offer **strong MMLU performance but**  
1254 **underperform on leaderboard metrics**, while **lower ratios** result in weaker performance across  
1255 BBH, GPQA, and most benchmarks, likely due to their similarity to a near-uniform distribution.

1256  
1257  
1258  
1259  
1260  
1261  
1262  
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 G CODE

1297

1298 We provide access to anonymous version of our code: <sup>1</sup>Anonymous Code

1299

1300

1301

1302

1303

1304

1305

1306

1307

1308

1309

1310

1311

1312

1313

1314

1315

1316

1317

1318

1319

1320

1321

1322

1323

1324

1325

1326

1327

1328

1329

1330

1331

1332

1333

1334

1335

1336

1337

1338

1339

1340

1341

1342

1343

1344

1345

1346

1347

1348

1349

---

<sup>1</sup><https://anonymous.4open.science/r/task-mixtures-62D3>