

FLARE: Task-Agnostic Embedding Model Evaluation via Normalizing Flows

Anonymous ACL submission

Abstract

Despite the widespread adoption of text embedding models, selecting the optimal model for a specific target corpus remains challenging due to the lack of task-specific labels. While task-agnostic evaluation offers a promising solution by relying on unlabeled data, existing approaches based on kernel estimators or Gaussian mixtures fail to model high-dimensional distributions effectively, resulting in unstable rankings. To address this limitation, we propose **FLARE** (Flow-based Label-free Assessment of Representation Embeddings), which leverages normalizing flows to estimate information sufficiency in high-dimensional spaces. By learning invertible transformations, flows enable exact density estimation while mitigating the instability inherent in distance-based methods. We provide theoretical guarantees showing that our estimation error depends on the data’s intrinsic structure rather than its raw dimensionality. Experiments across 11 datasets demonstrate that FLARE achieves a strong Spearman’s ρ (up to 0.90) with supervised benchmarks, remaining robust even for high-dimensional embeddings ($d \geq 3,584$).

1 Introduction

Recent advances in text embeddings have produced powerful semantic representation models such as Qwen3 Embedding (Yang et al., 2025) and Gemini Embedding (Lee et al., 2025). However, as the number of available models grows, each with different architectures, training objectives, and pretraining corpora, selecting the most suitable model for a given corpus has become increasingly challenging. The standard approach relies on supervised benchmarks like MTEB (Muennighoff et al., 2023), ranking models by their performance on annotated tasks.

This approach requires labeled data, which is often unavailable in practice. Consider deploying a retrieval system over proprietary documents such

as legal contracts, medical records, or financial reports. These collections have no existing labeled query-document pairs, and creating annotations requires significant time and domain expertise. Specialized corpora may also differ substantially from public benchmarks in vocabulary, style, and topic distribution (Tang and Yang, 2025). An embedding model that ranks highly on MTEB may perform poorly on domain-specific text, but without labels we cannot measure this gap. Benchmark contamination further undermines public leaderboards: as test sets appear in pretraining data, scores become inflated. This raises our central question: *how can we evaluate embedding models without labels?*

Recent work has explored task-agnostic evaluation using only unlabeled corpora. One approach analyzes geometric properties of embedding models like uniformity and alignment (Wang and Isola, 2020; Rudman et al., 2022). However, these metrics measure the embedding hypersphere structure rather than semantic content. A random projection can achieve perfect uniformity while preserving no information. A more principled alternative estimates mutual information between embeddings, quantifying how much information the embedding retains. Existing implementations use non-parametric estimators like Kernel Density Estimation (KDE) or Gaussian Mixtures (Darrin et al., 2024). These methods suffer from the curse of dimensionality (Beirlant et al., 1997): as embedding dimension grows (modern models often exceed $d = 3,000$), reliable density estimation requires exponentially more data. In high-dimensional space, these estimators become unstable and fail to predict downstream performance.

In this work, we propose Flow-based Label-free Assessment of Representation Embeddings (**FLARE**), a framework grounded in **information-theoretic sufficiency** (Darrin et al., 2024). Specifically, we quantify embedding quality by measuring the reduction in uncertainty about input data

given the embedding. The key insight is to use Normalizing Flows (Durkan et al., 2019), deep generative models that learn invertible transformations from complex distributions to simple base densities. Flows enable exact log-likelihood estimation via the change-of-variables formula, effectively mitigating the curse of dimensionality inherent to distance-based estimators. Our finite-sample bounds provide theoretical justification: the estimation error depends on the intrinsic effective dimension of the data manifold, not the embedding dimension. This ensures that FLARE remains reliable as embedding dimensions scale toward the high-dimensional space of modern LLM-based embedding models.

We validate FLARE on 11 datasets across diverse tasks. Our method achieves strong rank correlation with supervised performance (Spearman’s ρ up to 0.90), significantly outperforming geometry-based and information theoretical-based baselines. Importantly, FLARE remains stable in high-dimensional settings ($d \geq 3,584$) where information theoretical-based method fails.

Our contributions are summarized as follows:

- We introduce FLARE, a task-agnostic text embedding evaluation framework using normalizing flows to estimate information sufficiency without labeled data.
- We derive finite-sample generalization bounds showing robustness in high dimensions through dependence on effective rather than high-dimensional space.
- We demonstrate that FLARE reliably predicts downstream performance and outperforms existing baselines, especially for high-dimensional LLM-based embeddings.

2 Related Work

Embedding Evaluation. Modern text embedding models, often leveraging Large Language Model (LLM) architectures, provide high-dimensional representations that generalize across diverse semantic tasks without requiring task-specific fine-tuning (Neelakantan et al., 2022; Wang et al., 2024). Currently, the Massive Text Embedding Benchmark (MTEB) (Muennighoff et al., 2023) and its specialized variants like FinMTEB (Tang and Yang, 2025) serve as the primary evaluation standards. These benchmarks aggregate performance across diverse labeled tasks including

clustering, retrieval, and semantic textual similarity (STS). However, these methods depend heavily on ground-truth annotations. This dependency makes them unsuitable for assessing model utility on proprietary, dynamic, or out-of-distribution corpora where labels are absent, highlighting the need for unsupervised metrics that can quantify representation quality directly. To bridge this gap, we propose FLARE, a task-agnostic framework that quantifies embedding quality directly from the data distribution, offering a reliable performance proxy without the need for supervision.

Task-Agnostic Approaches. Task-agnostic metrics offer a long-standing alternative to supervised benchmarks by eliminating labeling costs. Classic indices such as the Silhouette Score (Rousseeuw, 1987) evaluate cluster separation, while more recent studies emphasize spectral and geometric properties. Examples include Effective Rank (Roy and Vetterli, 2007) for dimensionality estimation along with Uniformity (Wang and Isola, 2020) and IsoScore (Rudman et al., 2022) for spatial distribution analysis. A primary limitation of these metrics is their focus on geometric structure rather than semantic content. Such methods often rely on global priors like isotropy, which may not represent the intrinsic low-dimensional manifold structure characteristic of text embeddings. To address this, EMIR (Darrin et al., 2024) introduce the concept of Information Sufficiency to quantify how well one embedding model can reconstruct another. However, standard implementations of EMIR utilize Gaussian Mixture Models (GMM), which struggle with high-dimensional data and lack the generalization guarantees required for modern LLM-based embeddings.

Density Estimation. Calculating information-theoretic measures like differential entropy requires an accurate model of the underlying probability density. Traditional non-parametric approaches, notably Kernel Density Estimation (KDE), are fundamentally limited by the curse of dimensionality (Silverman, 2018; Beirlant et al., 1997). In high-dimensional spaces, these methods become statistically inefficient because the sample size needed to control estimation error grows exponentially with dimensionality. While neural variational estimators such as MINE (Ishmael Belghazi et al., 2018) and CLUB (Cheng et al., 2020) improve scalability, they optimize variational bounds instead of exact likelihoods. These bounds often

suffer from a severe bias-variance trade-off, leading to loose estimates and optimization instability (Poole et al., 2019). Normalizing Flows provide a robust alternative by learning a sequence of invertible transformations that map complex data to a simple base distribution (Rezende and Mohamed, 2015; Dinh et al., 2016). Crucially, they allow for exact log-likelihood computation through the change-of-variables formula (Papamakarios et al., 2021). Our work leverages this property to formulate embedding evaluation as a precise density estimation problem. By integrating flows into the information-sufficiency framework (Darrin et al., 2024), we ensure that our metrics remain theoretically grounded and empirically stable even for high-dimensional embeddings.

3 Method

To evaluate embedding quality without task-specific labels, we propose FLARE, which quantifies representation quality by measuring information sufficiency using normalizing flows.

3.1 Problem Formulation

Given an unlabeled corpus \mathcal{X} , we consider a set of candidate embedding models $\mathcal{E} = \{E_1, \dots, E_K\}$. Each model $E \in \mathcal{E}$ maps an input text $x \in \mathcal{X}$ to a high-dimensional representation:

$$z = E(x) \in \mathbb{R}^d. \quad (1)$$

To evaluate a specific model E_a , we pair it with a reference model E_b ($b \neq a$). We denote the embedding being evaluated as the source U , and the reference embedding as the target V :

$$U = E_a(x), \quad V = E_b(x). \quad (2)$$

Our objective is to derive a task-agnostic score for E_a based solely on these representations such that the resulting model ranking aligns with downstream supervised performance.

3.2 Information-Sufficiency Score

We build upon the information-sufficiency framework of (Darrin et al., 2024). The core intuition is that a high-quality source embedding U should act as a sufficient representation of the semantic space, enabling the reconstruction of the target representations V . We formalize this via Information-Sufficiency (I_s), which measures the reduction in uncertainty of V once U is observed.

Let \mathcal{F} be a family of marginal densities and \mathcal{K} be a family of conditional densities. The directional I_s score from U to V is defined as the difference between marginal and conditional entropy:

$$I_s(U \rightarrow V) = \underbrace{\inf_{f \in \mathcal{F}} \mathbb{E}_v[-\log f(v)]}_{H(V)} - \underbrace{\mathbb{E}_u \left[\inf_{M \in \mathcal{K}} \mathbb{E}_{v|u}[-\log M(v|u)] \right]}_{H(V|U)}. \quad (3)$$

To obtain a single quality score for model E_a , we compute the normalized median of its pairwise scores against all other models in the pool:

$$I_{s\text{norm}}(E_a) = \text{median}_{b \neq a} \frac{I_s(U_a \rightarrow U_b)}{\dim(U_b)}. \quad (4)$$

Normalization by the target dimension $\dim(U_b)$ is essential for comparability across reference models with varying output sizes, as raw entropy naturally scales with dimensionality.

3.3 Normalizing-Flow Implementation

We instantiate the density families in Eq. 3 using normalizing flows (Durkan et al., 2019). Flows enable exact log-likelihood computation via invertible transformations, providing stable estimation.

Two-Stage Training. As shown in Figure 1, we employ a progressive training strategy. We first train a marginal flow $p_\phi(v)$ on target embeddings to model distribution. Next, we initialize a conditional flow $p_\theta(v|u)$ by copying the marginal backbone weights. This warm-start strategy ensures the conditional model begins from a well-defined density baseline.

Low-Rank Conditioning. Standard conditional flows often use hypernetworks with $O(d^2)$ complexity, which is prohibitive for high-dimensional embeddings. We instead inject the source information u through a parameter-efficient low-rank residual branch. Let \mathbf{h}_{base} be the intermediate features of the target flow. The conditional feature is computed as:

$$\mathbf{h}_{\text{cond}} = \mathbf{h}_{\text{base}} + B(A(u)), \quad (5)$$

where $A \in \mathbb{R}^{r \times d}$ projects the d -dimensional source embedding to a low-rank bottleneck of dimension $r = 64$, and B maps from the bottleneck to the output space. This mechanism allows the source u

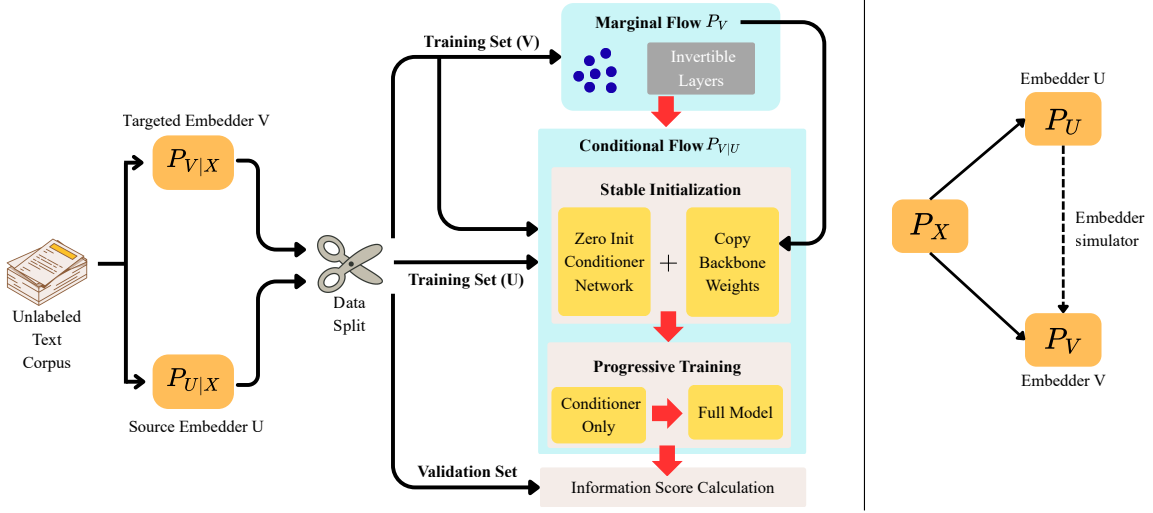


Figure 1: Overview of our two stage flow based estimation pipeline. **Stage 1:** Train a marginal flow $p_\phi(v)$ on target embeddings V to model their intrinsic distribution. **Stage 2:** Initialize a conditional flow $p_\theta(v|u)$ by copying the marginal backbone weights and adding a zero-initialized low-rank conditioning branch. This branch is then trained to capture the dependency between source embeddings U and target embeddings V , enabling computation of the information-sufficiency score via Eq. 3.

to adjust the transformation trajectory of the target v with minimal parameter overhead.

Zero Initialization. We initialize B to zeros and A with random initialization. At the start of training, the conditioning term $B(A(u))$ is zero, making the conditional flow $p_\theta(v|u)$ initially equivalent to the pre-trained marginal flow $p_\phi(v)$. The dependence on the source U is learned gradually, which prevents gradient instability and improves convergence speed.

4 Theoretical Justification

Motivation. In real-world deployment, embedding models must generalize to vast amounts of unseen data that extend far beyond the validation set. Purely empirical evaluation on a fixed dataset cannot theoretically guarantee reliability on the underlying population distribution. To address this, we establish a theoretical framework relying on two pillars: the spectral stability of our flow architecture and the low-dimensional manifold hypothesis.

Core Assumptions. Our theory relies on two core assumptions (Appendix A.2): a low intrinsic dimension (Assumption 1) to enable scaling to high-dimensional embeddings, and layer-wise ap-

proximate independence (Assumption 2) to ensure stable gradient propagation. The latter, theoretically motivated by (Cohen et al., 2021), is practically realized by our Zero-Initialization strategy, which ensures the network exhibits stable, linear growth rather than exponential instability.

Finite-Sample Generalization Bound. Building on this stability, we verify the reliability of FLARE by bounding the gap between the training and validation losses.

Theorem 1 (Finite-sample generalization bound). *Under Assumptions 1 and 2, for any fixed flow model p_θ and any $\delta \in (0, 1)$, with probability at least $1 - \delta$, the train-validation gap satisfies:*

$$\begin{aligned}
 |\hat{L}_{\text{val}}(\theta) - \hat{L}_{\text{train}}(\theta)| &\leq \frac{2\tilde{C}_{\text{Rad}} L \bar{\sigma} \sqrt{d_{\text{eff}}}}{\sqrt{m}} \\
 &+ M_{\text{val}} \sqrt{\frac{\log(2/\delta)}{2m_{\text{val}}}} \\
 &+ 3M_{\text{train}} \sqrt{\frac{\log(2/\delta)}{2m}}.
 \end{aligned} \tag{6}$$

Interpretation. The bound scales with the intrinsic dimension d_{eff} rather than the raw embedding dimension d . Since $d_{\text{eff}} \ll d$ in semantic representations, this property allows for reliable evaluation

on high-dimensional embeddings using only moderate sample sizes, effectively mitigating the curse of dimensionality. Furthermore, the error depends linearly on the depth L and the spectral stability $\bar{\sigma}$. By minimizing this term via our Zero-Initialization strategy, we explicitly control the model’s complexity to prevent overfitting. These factors jointly justify why FLARE remains robust, whereas traditional kernel-based estimators often degrade in high-dimensional spaces.

5 Experiments

We design our experiments to evaluate the empirical utility of the Flow-based estimator across three primary dimensions: **(Q1)** its reliability in predicting model rankings across diverse downstream tasks; **(Q2)** its comparative performance against kernel-based baselines in high-dimensional embedding spaces; and **(Q3)** the alignment between empirical observations and our theoretical generalization guarantees.

5.1 Embedders and Datasets

Embedding Models. To ensure a comprehensive evaluation, we select eight representative embedding models spanning various architectures and dimensionalities. We focus specifically on high-dimensional space, including BGE-Multilingual-Gemma2 ($d = 3,584$) (Chen et al., 2024a), gte-Qwen2-7B-instruct ($d = 3,584$) (Li et al., 2023), and four Mistral-7B-based models ($d = 4,096$): Zeta-Alpha-E5-Mistral¹, Grit7B (Muennighoff et al., 2024), SFR-Embedding-Mistral (Rui Meng, 2024), and Linq-Embed-Mistral (Junseong Kim, 2024). For lower-dimensional benchmarks, we include stella-base-en-v2² ($d = 768$) and all-MiniLM-L6-v2 (Reimers and Gurevych, 2019) ($d = 384$).

Dataset Selection and Task Categorization. To simulate realistic enterprise scenarios where models are evaluated on novel, internal corpora, we consciously deviate from standard benchmarks like MTEB. Since public benchmarks frequently permeate pre-training corpora, leading to data leakage and inflated scores, we instead employ a strict temporal selection strategy. We choose 11 datasets from Hugging Face released *after* the training cut-

¹<https://huggingface.co/zeta-alpha-ai/Zeta-Alpha-E5-Mistral>

²<https://huggingface.co/infgrad/stella-base-en-v2>

off dates of all candidate models, thereby ensuring zero data contamination and mimicking the challenge of genuinely unseen distributions. Crucially, while FLARE operates in a fully unsupervised manner, we specifically select datasets containing ground-truth labels solely to serve as an oracle for validating the accuracy of our predicted rankings. Detailed dataset statistics are provided in Appendix B. These datasets cover four core task categories:

- **Classification:** Apt-eval (Saha and Feizi, 2025) (safety/robustness), GT-FintechLab (Shah et al., 2025) (finance), and BhashaBench-Finance (Devane et al., 2025) (multilingual finance).
- **Retrieval:** AIR-Bench-Finance (Chen et al., 2024b), LIMIT (Weller et al., 2025) (instruction-following), and ArXiv-Abstracts 2025³ (scientific literature).
- **Semantic Textual Similarity (STS):** Augmented STS-B⁴, LivNLP-STS (Zhang et al., 2025), and Philosophical-STS⁵.
- **Clustering:** Clustered-FunPang Medical⁶, and Reasoning-Clustering⁷.

Evaluation Protocol. We assess the reliability of our unsupervised Information Sufficiency (IS) metric by measuring its alignment with ground-truth supervised rankings. Ground-truth performance is established using standard MTEB metrics (Muennighoff et al., 2023): accuracy for classification, nDCG@10 for retrieval, Spearman correlation for STS, and V-measure for clustering. We quantify ranking alignment using Spearman’s rank correlation (ρ) and Pearson correlation (r) between the predicted IS scores and the supervised metrics.

Baselines. We benchmark FLARE against existing unsupervised metrics, explicitly categorized into two types of unsupervised evaluation methods: geometric-based and information-theoretic

³https://huggingface.co/datasets/almanach/arxiv_abstracts_2025

⁴https://huggingface.co/datasets/maiammar/augmented_stsb_multi_mt

⁵<https://huggingface.co/datasets/johnnyboycurtis/Philosophical-STS-Text-Pairs>

⁶https://huggingface.co/datasets/mukulb/clustered_FUNPANG_dataset_with_groups

⁷https://huggingface.co/datasets/Ibisbill/Clustering_deduplicated_reasoning

based. For the geometric-based methods, we considered (1) **Uniformity** (Wang and Isola, 2020) and (2) **IsoScore** (Rudman et al., 2022); and (3) **Silhouette Score** (Rousseeuw, 1987). For the information-theoretic based methods, we considered (4) **EMIR** (Darrin et al., 2024). Our work aligns with the information-theoretic based evaluation (specifically the framework established by EMIR), as we share the fundamental objective of quantifying representation quality via information retention. However, we diverge by employing normalizing flows to robustly estimate densities in high-dimensional spaces where the GMMs used in EMIR may fail.

5.2 Main Results

Table 1 summarizes the ranking correlations across eleven representative datasets, revealing three key findings. Details experiment results are shown in Appendix B

Flow-Based Estimation Succeeds Where Kernel Methods Fail. FLARE achieves an average Spearman correlation of $\rho = 0.69$, substantially outperforming all unsupervised baselines. Most strikingly, the kernel-based EMIR which shares our information-sufficiency framework but relies on Gaussian Mixture Model (GMM) density estimation—yields a negative average correlation ($\rho = -0.12$), indicating systematic ranking inversions. This failure is not incidental: in high-dimensional spaces ($d \geq 3,584$), the curse of dimensionality causes kernel densities to become vanishingly sparse, making distance-based estimates unreliable. Our flow-based approach circumvents this issue by learning an explicit parametric density model that adapts to the intrinsic manifold structure, maintaining positive correlations across *all* task categories.

Geometric Baselines Exhibit Inconsistent Performance. Table 1 reveals that existing unsupervised metrics struggle to maintain consistent correlations across task types. Uniformity achieves near-zero average correlation ($\rho = 0.05$), suggesting that embedding space uniformity alone is not predictive of downstream quality. IsoScore exhibits negative correlations across nearly all categories (Avg $\rho = -0.27$), as it penalizes anisotropy under the assumption that uniformity maximizes information capacity. However, high-quality semantic spaces are inherently anisotropic—meaningful concepts naturally form dense, non-uniform clusters

on the manifold rather than populating the hypersphere uniformly. Silhouette shows task-specific success only on STS ($\rho = 0.56$) but fails elsewhere, indicating that geometric cohesion does not generalize as a universal quality indicator. As visualized in Figure 2, these rigid geometric assumptions lead to extreme variance and frequent ranking inversions. Consequently, we propose that measuring Information Sufficiency via mutual information offers a more reliable basis for evaluation than geometric properties alone. By aligning with this information-theoretic objective, FLARE adapts to the intrinsic data density without enforcing conflicting geometric shapes, maintaining robust alignment with ground truth.

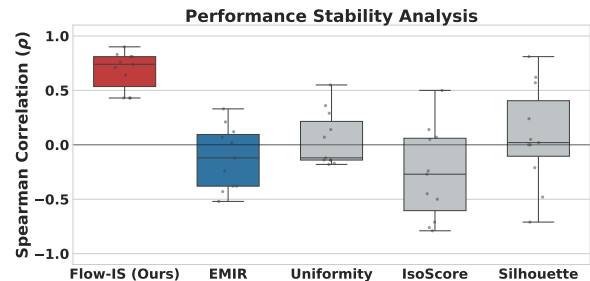


Figure 2: **Stability Analysis.** Distribution of Spearman correlations across all datasets. While geometric baselines (grey) and kernel-based EMIR (blue) exhibit high variance and negative correlations, our FLARE (red) maintains consistent positive alignment with ground truth, demonstrating superior robustness.

Consistent Performance Across Task Diversity. Unlike baselines that excel only on specific task types but fail to generalize, FLARE demonstrates robust generalization: $\rho = 0.83$ on Clustering, $\rho = 0.72$ on Retrieval, $\rho = 0.70$ on STS, and $\rho = 0.56$ on Classification. This cross-task stability indicates that our information-sufficiency metric captures a fundamental quality of embeddings: the ability to preserve semantic information that transfers across diverse downstream objectives. The per-task breakdown in Table 6 (Appendix B) reveals no negative correlations for our method, a property not shared by any baseline. Notably, we observed stronger results on STS, Retrieval, and Clustering tasks compared to classification. We attribute this to the nature of these tasks: they directly leverage embedding geometry through similarity measures without training additional models. In contrast, classification involves training a separate predictor that may exploit task-specific features beyond the overall distributional structure captured

Task	Spearman’s ρ					Pearson’s r				
	Uni.	Iso.	Sil.	EMIR	Ours	Uni.	Iso.	Sil.	EMIR	Ours
Class.	0.18	-0.40	-0.06	-0.06	0.56	0.35	-0.58	-0.08	-0.20	0.31
STS	0.01	-0.33	0.56	-0.06	0.70	0.70	-0.14	0.53	-0.08	0.70
Retr.	0.08	-0.45	-0.03	-0.22	0.72	0.43	-0.37	0.39	-0.18	0.64
Clust.	-0.14	0.29	-0.24	-0.16	0.83	-0.43	0.05	-0.25	-0.10	0.69
Avg	0.05	-0.27	0.08	-0.12	0.69	0.33	-0.29	0.18	-0.14	0.58

Table 1: Task-aggregated Comparison with Unsupervised Baselines.

Task Type	Bound Ratio	Rademacher %
Classification	11.0×	92.4%
STS	21.1×	94.5%
Retrieval	21.2×	95.5%
Clustering	18.4×	98.5%
Average	17.9×	95.2%

Table 2: **Theoretical validation.** Bound Ratio ($\Delta_{\text{theo}}/\Delta_{\text{emp}}$) and Rademacher complexity contribution, grouped by task type.

by our unsupervised metric.

Ranking Stability under Subsampling. To assess the sensitivity of FLARE to sample size, we evaluate the stability of I_s scores by subsampling the evaluation sets at ratios $\alpha \in \{0.05, 0.1, \dots, 1.0\}$. FLARE maintains a remarkably strong correlation with ground-truth rankings even when the sample size is reduced to 20% of the original corpus. This stability indicates that the normalizing flow successfully captures the global manifold structure from a limited set of observations, effectively bypassing the high sample complexity required by traditional estimators. From a practical perspective, this result suggests that FLARE enables reliable model selection on small-scale specialized datasets without the need for extensive data collection. Comprehensive results and per-dataset stability curves are provided in Appendix C.

5.3 Generalization Bound Analysis

To empirically validate the guarantee provided by Theorem 1, we compare the derived theoretical bound against the observed empirical generalization gap. This validation is crucial for addressing the practical challenges of enterprise deployment, ensuring that the evaluation method generalizes

reliably as new data continuously arrives.

Setup. Since the Information Sufficiency estimator decomposes into marginal and conditional components, its total estimation error is bounded by the sum of their respective generalization gaps. We therefore compute the empirical gap as the sum of the absolute differences between training and validation Negative Log-Likelihoods (NLL) for both flows, and compare this empirical quantity against the theoretical bound derived from our model architecture and sample complexity.

Analysis. As reported in Table 2, the theoretical bound Δ_{theo} upper-bounds the empirical gap Δ_{emp} by a margin ranging from 11.0× to 21.2× across task types (average ratio 17.9×). This confirms that our conservative linear bound remains informative, providing meaningful generalization guarantees without being vacuously loose. The Rademacher complexity term accounts for 92.4% to 98.5% of the total bound (average 95.2%), which aligns with its dependence on the effective intrinsic dimension d_{eff} . As expected, higher-dimensional embedding spaces incur larger complexity penalties. Notably, classification tasks exhibit tighter bounds (11.0×) compared to retrieval and STS tasks ($\sim 21\times$), a discrepancy that reflects the richer representational capacity required for fine-grained semantic matching in the latter. Practically, these findings certify that FLARE provides a trustworthy and theoretically grounded signal on unseen data.

5.4 Ablation Study

Shuffle ablation. We partially shuffle the correspondence between source embeddings U and target embeddings V at a ratio $p \in [0, 1]$ while keeping the remaining pairs unchanged. As visu-

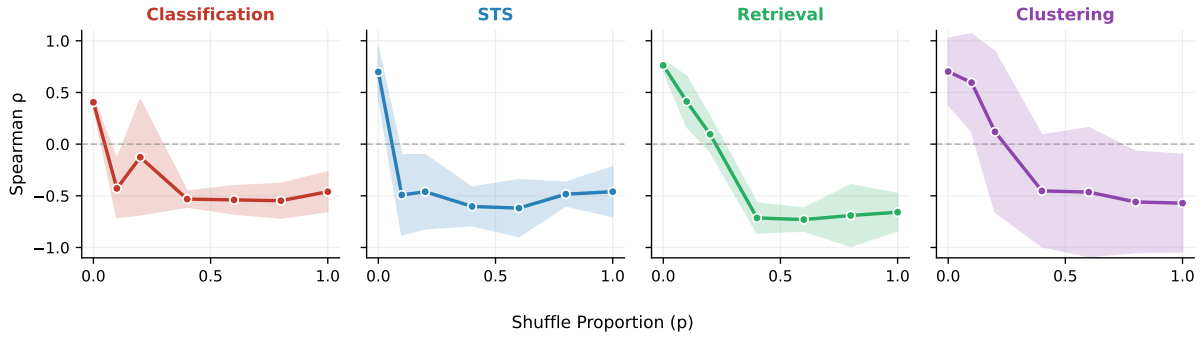


Figure 3: Partial shuffle ablation by task type. Increasing the shuffle proportion p causes correlation to degrade from positive ($p = 0$) to negative ($p = 1$), confirming that the metric relies on correct alignment.

alized in Figure 3 and detailed in Appendix D.1, the Spearman correlation with downstream performance degrades significantly as p increases, transitioning from positive to negative around $p = 0.2$ to 0.4 . The degradation pattern varies across task types: Retrieval and Clustering exhibit a gradual decline, whereas Classification and STS show sharper transitions at lower shuffle ratios. At full shuffle ($p = 1.0$), all categories converge to negative correlations (averaging $\rho \approx -0.4$). This confirms that our metric relies heavily on the correct **alignment between source and target embeddings** rather than marginal statistics, and the varying sensitivity across tasks may reflect differences in the underlying embedding structure.

Conditional-only ablation. We investigate the contribution of the marginal term by comparing the full metric against a conditional-only variant, defined as $I_{\text{cond}}(u, v) = \log p_{\phi}(v | u)$. While this variant may be competitive when the target marginal distribution varies little, it lacks consistent reliability across diverse tasks. As quantified in Table 8 and visualized in Figure 5 (see Appendix D.2), removing the marginal term $\mathbb{E}[\log p(Y)]$ reduces the average Spearman correlation substantially from 0.69 to 0.21. This empirical evidence confirms that both components are essential: the conditional term captures the cross-model mapping, while the marginal term accounts for the intrinsic structure of the target space. Consequently, we retain the full score as the default to ensure robust alignment with downstream performance.

Aggregation strategy. We adopt the median rather than the mean to compute the final score, as the median is more robust to outliers and less sensitive to skewed distributions of model performances. When many models exhibit similar behav-

ior, the mean can be biased by minor perturbations, whereas the median remains stable by focusing on the central tendency of the ranked distribution. We further investigate whether trimmed aggregation mitigates this bias. The results are summarized in Table 9. Across all task types, the median consistently achieves the highest or near-highest correlation with ground-truth rankings, while the mean shows greater variance and occasional degradation. Based on these findings, we adopt the median as the default aggregation method for its robustness and simplicity.

6 Conclusion

In this work, we addressed the critical challenge of evaluating representation quality in scenarios where downstream labels are inaccessible. To this end, we proposed FLARE, a principled framework bridging information theory with deep generative modeling. By leveraging normalizing flows to model complex embedding densities, FLARE precisely quantifies information sufficiency to rank models effectively without supervision. Beyond empirical results, we established the theoretical foundations of our approach by deriving finite-sample generalization guarantees based on Rademacher complexity. These bounds validate the estimator’s robustness in high-dimensional spaces, explicitly linking performance to the low intrinsic dimension of semantic manifolds. Extensive experiments across 11 datasets confirm that FLARE consistently correlates with ground-truth performance and maintains stability across diverse tasks. Ultimately, FLARE offers a scalable and mathematically grounded solution for automated model selection, significantly reducing the community’s reliance on expensive annotated benchmarks.

614 Limitations

615 Despite its promising results, FLARE has limita-
616 tions. First, as a learning-based evaluation method,
617 FLARE incurs higher computational costs than
618 training-free statistical measures. While we miti-
619 gate inference latency by adopting discrete normal-
620 izing flows instead of continuous-time approaches
621 like flow matching, the requirement to train a gen-
622 erative model remains a bottleneck for resource-
623 constrained scenarios.

624 Second, we observe a performance discrep-
625 ancy across task types. While FLARE excels on
626 geometry-centric tasks like STS and Retrieval, its
627 correlation is lower for Classification ($\rho = 0.56$).
628 This suggests that global information sufficiency
629 captures the overall semantic manifold but may
630 miss fine-grained, class-specific features required
631 for linear separability.

632 Third, our theoretical analysis relies on simplify-
633 ing assumptions. Specifically, regarding Assump-
634 tion 2, parameter coupling during end-to-end back-
635 propagation inevitably introduces inter-layer cor-
636 relations, making the approximate independence
637 assumption an idealization that is computationally
638 intractable to verify in high dimensions.

639 References

640 Alessio Ansuini, Alessandro Laio, Jakob H Macke, and
641 Davide Zoccolan. 2019. Intrinsic dimension of data
642 representations in deep neural networks. *Advances*
643 *in Neural Information Processing Systems*, 32.

644 Peter L Bartlett and Shahar Mendelson. 2002.
645 Rademacher and gaussian complexities: Risk bounds
646 and structural results. *Journal of machine learning*
647 *research*, 3(Nov):463–482.

648 Jan Beirlant, Edward J Dudewicz, László Györfi, Ed-
649 ward C Van der Meulen, and 1 others. 1997. Non-
650 parametric entropy estimation: An overview. *Inter-
651 national Journal of Mathematical and Statistical*
652 *Sciences*, 6(1):17–39.

653 Jianlv Chen, Shitao Xiao, Peitian Zhang, Kun Luo, Defu
654 Lian, and Zheng Liu. 2024a. [Bge m3-embedding:
655 Multi-lingual, multi-functionality, multi-granularity
656 text embeddings through self-knowledge distillation.](#)
657 *Preprint*, arXiv:2402.03216.

658 Jianlyu Chen, Nan Wang, Chaofan Li, Bo Wang, Shi-
659 tao Xiao, Han Xiao, Hao Liao, Defu Lian, and
660 Zheng Liu. 2024b. [Air-bench: Automated hetero-
661 geneous information retrieval benchmark.](#) *Preprint*,
662 arXiv:2412.13102.

663 Pengyu Cheng, Weituo Hao, Shuyang Dai, Jiachang
664 Liu, Zhe Gan, and Lawrence Carin. 2020. Club:

A contrastive log-ratio upper bound of mutual in-
665 formation. In *International conference on machine*
666 *learning*, pages 1779–1788. PMLR. 667

668 Alain-Sam Cohen, Rama Cont, Alain Rossier, and
669 Renyuan Xu. 2021. [Scaling properties of deep resid-
670 ual networks.](#) In *Proceedings of the 38th Interna-
671 tional Conference on Machine Learning*, volume 139
672 of *Proceedings of Machine Learning Research*, pages
673 2039–2048. PMLR.

674 Maxime Darrin, Philippe Formont, Ismail Ayed,
675 Jackie CK Cheung, and Pablo Piantanida. 2024.
676 When is an embedding model more promising than
677 another? *Advances in Neural Information Process-
678 ing Systems*, 37:68330–68379.

679 Vijay Devane, Mohd Nauman, Bhargav Patel,
680 Aniket Mahendra Wakchoure, Yogeshkumar Sant,
681 Shyam Pawar, Viraj Thakur, Ananya Godse, Sunil Pa-
682 tra, Neha Maurya, Suraj Racha, Nitish Kamal Singh,
683 Ajay Nagpal, Piyush Sawarkar, Kundeshwar Vijayrao
684 Pundalik, Rohit Saluja, and Ganesh Ramakrishnan.
685 2025. [Bhashabench v1: A comprehensive bench-
686 mark for the quadrant of indic domains.](#) *Preprint*,
687 arXiv:2510.25409.

688 Laurent Dinh, Jascha Sohl-Dickstein, and Samy Ben-
689 gio. 2016. Density estimation using real nvp. *arXiv*
690 *preprint arXiv:1605.08803*.

691 Conor Durkan, Artur Bekasov, Iain Murray, and George
692 Papamakarios. 2019. Neural spline flows. *Advances*
693 *in neural information processing systems*, 32.

694 Wassily Hoeffding. 1994. *Probability Inequalities for*
695 *Sums of Bounded Random Variables*, pages 409–426.
696 Springer New York, New York, NY.

697 Mohamed Ishmael Belghazi, Aristide Baratin, Sai
698 Rajeswar, Sherjil Ozair, Yoshua Bengio, Aaron
699 Courville, and R Devon Hjelm. 2018. Mine: mu-
700 tual information neural estimation. *arXiv e-prints*,
701 pages arXiv–1801.

702 Jihoon Kwon Sangmo Gu Yejin Kim Minkyung Cho
703 Jy-yong Sohn Chanyeol Choi Junseong Kim, Seol-
704 hwa Lee. 2024. [Linq-embed-mistral: elevating text
705 retrieval with improved gpt data through task-specific
706 control and quality refinement.](#) Linq AI Research
707 Blog.

708 Jinhyuk Lee, Feiyang Chen, Sahil Dua, Daniel
709 Cer, Madhuri Shanbhogue, Iftekhar Naim, Gus-
710 tavo Hernández Ábrego, Zhe Li, Kaifeng Chen, Hen-
711 rique Schechter Vera, and 1 others. 2025. Gemini
712 embedding: Generalizable embeddings from gemini.
713 *arXiv preprint arXiv:2503.07891*.

714 Zehan Li, Xin Zhang, Yanzhao Zhang, Dingkun Long,
715 Pengjun Xie, and Meishan Zhang. 2023. Towards
716 general text embeddings with multi-stage contrastive
717 learning. *arXiv preprint arXiv:2308.03281*.

718	Niklas Muennighoff, Hongjin Su, Liang Wang, Nan Yang, Furu Wei, Tao Yu, Amanpreet Singh, and Douwe Kiela. 2024. Generative representational instruction tuning . <i>Preprint</i> , arXiv:2402.09906.	
719		
720		
721		
722	Niklas Muennighoff, Nouamane Tazi, Loic Magne, and Nils Reimers. 2023. MTEB: Massive text embedding benchmark . In <i>Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics</i> , pages 2014–2037, Dubrovnik, Croatia. Association for Computational Linguistics.	
723		
724		
725		
726		
727		
728	Arvind Neelakantan, Tao Xu, Raul Puri, Alec Radford, Jesse Michael Han, Jerry Tworek, Qiming Yuan, Nikolas Tezak, Jong Wook Kim, Chris Hallacy, and 1 others. 2022. Text and code embeddings by contrastive pre-training. <i>arXiv preprint arXiv:2201.10005</i> .	
729		
730		
731		
732		
733		
734	George Papamakarios, Eric Nalisnick, Danilo Jimenez Rezende, Shakir Mohamed, and Balaji Lakshminarayanan. 2021. Normalizing flows for probabilistic modeling and inference. <i>Journal of Machine Learning Research</i> , 22(57):1–64.	
735		
736		
737		
738		
739	Ben Poole, Sherjil Ozair, Aaron Van Den Oord, Alex Alemi, and George Tucker. 2019. On variational bounds of mutual information. In <i>International conference on machine learning</i> , pages 5171–5180. PMLR.	
740		
741		
742		
743		
744	Nils Reimers and Iryna Gurevych. 2019. Sentence-BERT: Sentence embeddings using Siamese BERT-networks . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 3982–3992, Hong Kong, China. Association for Computational Linguistics.	
745		
746		
747		
748		
749		
750		
751		
752	Danilo Rezende and Shakir Mohamed. 2015. Variational inference with normalizing flows. In <i>International conference on machine learning</i> , pages 1530–1538. PMLR.	
753		
754		
755		
756	Peter J Rousseeuw. 1987. Silhouettes: a graphical aid to the interpretation and validation of cluster analysis. <i>Journal of computational and applied mathematics</i> , 20:53–65.	
757		
758		
759		
760	Olivier Roy and Martin Vetterli. 2007. The effective rank: A measure of effective dimensionality. In <i>2007 15th European signal processing conference</i> , pages 606–610. IEEE.	
761		
762		
763		
764	William Rudman, Nate Gillman, Taylor Rayne, and Carsten Eickhoff. 2022. Isoscore: Measuring the uniformity of embedding space utilization. In <i>Findings of the Association for Computational Linguistics: ACL 2022</i> , pages 3325–3339.	
765		
766		
767		
768		
769	Shafiq Rayhan Joty Caiming Xiong Yingbo Zhou Semih Yavuz Rui Meng, Ye Liu. 2024. Sfr-embedding-mistral:enhance text retrieval with transfer learning . Salesforce AI Research Blog.	
770		
771		
772		
	Shoumik Saha and Soheil Feizi. 2025. Almost AI, almost human: The challenge of detecting AI-polished writing . In <i>Findings of the Association for Computational Linguistics: ACL 2025</i> , pages 25414–25431, Vienna, Austria. Association for Computational Linguistics.	773 774 775 776 777 778
	Agam Shah, Siddhant Sukhani, Huzaifa Pardawala, Saketh Budideti, Riya Bhadani, Rudra Gopal, Sid-dhartha Somani, Rutwik Routu, Michael Galarnyk, Soungmin Lee, Arnav Hiray, Akshar Ravichandran, Eric Kim, Pranav Aluru, Joshua Zhang, Sebastian Jaskowski, Veer Guda, Meghaj Tarte, Liqin Ye, and 8 others. 2025. Words that unite the world: A unified framework for deciphering central bank communications globally . <i>Preprint</i> , arXiv:2505.17048.	779 780 781 782 783 784 785 786 787
	Bernard W Silverman. 2018. <i>Density estimation for statistics and data analysis</i> . Routledge.	788 789
	Yixuan Tang and Yi Yang. 2025. FinMTEB: Finance massive text embedding benchmark . In <i>Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing</i> , pages 3620–3638, Suzhou, China. Association for Computational Linguistics.	790 791 792 793 794 795
	Liang Wang, Nan Yang, Xiaolong Huang, Binxing Jiao, Linjun Yang, Daxin Jiang, Rangan Majumder, and Furu Wei. 2024. Text embeddings by weakly-supervised contrastive pre-training . <i>Preprint</i> , arXiv:2212.03533.	796 797 798 799 800
	Tongzhou Wang and Phillip Isola. 2020. Understanding contrastive representation learning through alignment and uniformity on the hypersphere. In <i>International conference on machine learning</i> , pages 9929–9939. PMLR.	801 802 803 804 805
	Orion Weller, Michael Boratko, Iftexhar Naim, and Jinhyuk Lee. 2025. On the theoretical limitations of embedding-based retrieval . <i>Preprint</i> , arXiv:2508.21038.	806 807 808 809
	An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, and 1 others. 2025. Qwen3 technical report. <i>arXiv preprint arXiv:2505.09388</i> .	810 811 812 813 814
	Gaifan Zhang, Yi Zhou, and Danushka Bollegala. 2025. Annotating training data for conditional semantic textual similarity measurement using large language models . <i>Preprint</i> , arXiv:2509.14399.	815 816 817 818
	A Theoretical Assumptions and Proofs	819
	A.1 Problem setup	820
	Let \mathcal{D} be an unknown distribution on \mathbb{R}^d . We observe an i.i.d. training sample	821 822
	$S_{\text{train}} = \{v_1^{\text{train}}, \dots, v_m^{\text{train}}\} \sim \mathcal{D}^m.$ (7)	823
	A normalizing flow T_θ induces a density	824
	$p_\theta(v) = p_Z(T_\theta(v)) \det J_{T_\theta}(v) ,$ (8)	825

where p_Z is a fixed base density and $J_{T_\theta}(v)$ is the Jacobian of the flow transformation. We use the negative log-likelihood loss $\ell_\theta(v) = -\log p_\theta(v)$. The empirical training risk is

$$\hat{L}_{\text{train}}(\theta) = \frac{1}{m} \sum_{i=1}^m \ell_\theta(v_i^{\text{train}}). \quad (9)$$

The population risk is

$$L(\theta) = \mathbb{E}_{v \sim \mathcal{D}} [\ell_\theta(v)]. \quad (10)$$

Given a validation set

$$S_{\text{val}} = \{v_1^{\text{val}}, \dots, v_{m_{\text{val}}}^{\text{val}}\} \sim \mathcal{D}^{m_{\text{val}}}, \quad (11)$$

the empirical validation risk is

$$\hat{L}_{\text{val}}(\theta) = \frac{1}{m_{\text{val}}} \sum_{j=1}^{m_{\text{val}}} \ell_\theta(v_j^{\text{val}}). \quad (12)$$

We aim to bound the generalization gap between the validation risk (observable proxy for performance) and the training risk (optimization objective):

$$\Delta(\theta) = |\hat{L}_{\text{val}}(\theta) - \hat{L}_{\text{train}}(\theta)|. \quad (13)$$

Note that bounding this gap involves controlling the deviation of both empirical risks from the population risk $L(\theta)$.

We first consider a marginal flow for a fixed embedder, then apply the same argument to a conditional flow. The Information Sufficiency (I_s) score is defined as

$$IS = L_{\text{marg}}(V) - L_{\text{cond}}(V | U), \quad (14)$$

and the empirical \widehat{IS} score replaces the population risks with their validation counterparts. Establishing a finite-sample bound on the generalization gaps of the marginal and conditional flows directly yields a bound for the estimation error of the I_s score.

A.2 Core assumptions

Assumption 1 (Low intrinsic dimension). *The embedding distribution is supported on a compact subset $\mathcal{M} \subset \mathbb{R}^d$ with intrinsic dimension $d_{\text{eff}} \ll d$ and bounded diameter. This type of low-dimensional structure for learned representations is consistent with empirical measurements of intrinsic dimensionality in deep features (Ansuini et al., 2019).*

Assumption 2 (Approximate layer independence). *The flow T_θ is a composition of L invertible blocks with **perturbation rank** at most r . Assuming a regime of near-identity mappings (e.g., via residual connections or zero-initialization), the Jacobians of different blocks are approximately independent in their dominant singular directions. Consequently, the overall Lipschitz behaviour of the composition **grows linearly** in depth rather than exhibiting exponential growth. This assumption is consistent with analyses of signal propagation and dynamical isometry in deep architectures (Cohen et al., 2021).*

Remark: Validity via Zero-Initialization.

While assuming approximate independence between layer Jacobians is non-trivial for generic deep networks, our Zero-Initialization strategy (Section 3.3) provides a rigorous structural justification. Specifically: (1) **Identity Start**: By initializing the conditioner’s projection matrix B to zero, the conditional flow starts as an identity transformation relative to the backbone ($T_\theta = T_\phi$), ensuring that initial layer-wise Jacobians satisfy $J_l = I$. (2) **Linear Accumulation**: During progressive training, the deviation Δ_l from identity (where $J_l = I + \Delta_l$) is explicitly constrained via L_2 regularization. In this small-perturbation regime, the norm of the composed Jacobian follows a first-order approximation $\|\prod_{l=1}^L (I + \Delta_l)\| \approx \|I + \sum_{l=1}^L \Delta_l\|$, which implies that the Lipschitz constant grows **linearly** with depth L rather than exponentially. This design effectively aligns the flow’s architectural behavior with the stability postulated in Assumption 2.

A.3 Rademacher complexity on a manifold

Let \mathcal{F} be a class of real-valued functions on \mathcal{M} . Given a sample $S = \{v_1, \dots, v_m\} \subset \mathcal{M}$, the empirical Rademacher complexity of \mathcal{F} is

$$\mathcal{R}_m(\mathcal{F}) = \mathbb{E}_\sigma \left[\sup_{f \in \mathcal{F}} \frac{1}{m} \sum_{i=1}^m \sigma_i f(v_i) \right], \quad (15)$$

where σ_i are i.i.d. Rademacher variables.

Lemma 1 (Rademacher complexity on a manifold). *Let \mathcal{F} be the class of **linear functions** $f(v) = \langle w, v \rangle$ restricted to \mathcal{M} , where $\|w\| \leq L_f$. Under Assumption 1, there exists a constant $C_{\text{Rad}} > 0$ such that*

$$\mathcal{R}_m(\mathcal{F}) \leq \frac{C_{\text{Rad}} L_f D \sqrt{d_{\text{eff}}}}{\sqrt{m}}. \quad (16)$$

909 *Proof.* The manifold \mathcal{M} admits a covering number
910 bound of the form

$$911 \quad \mathcal{N}(\mathcal{M}, \varepsilon) \leq C_0 \left(\frac{D}{\varepsilon} \right)^{d_{\text{eff}}} \quad \text{for all } \varepsilon > 0. \quad (17)$$

912 Here, $\mathcal{N}(\mathcal{M}, \varepsilon)$ represents the covering number
913 of the manifold, and C_0 is a geometry-dependent
914 constant.

915 For the class of linear functions \mathcal{F} with bounded
916 norm $\|w\| \leq L_f$, the covering number of the func-
917 tion space is controlled by the covering number of
918 the domain \mathcal{M} via a duality argument. Specifically,
919 distinguishing two linear functions on \mathcal{M} is equiv-
920 alent to covering the domain at a finer scale ε/L_f .
921 Consequently, we have:

$$922 \quad \mathcal{N}(\mathcal{F}, \varepsilon) \leq \mathcal{N}\left(\mathcal{M}, \frac{\varepsilon}{L_f}\right) \quad (18)$$

$$923 \quad \leq C_0 \left(\frac{DL_f}{\varepsilon} \right)^{d_{\text{eff}}}. \quad (19)$$

924 Taking logarithms yields

$$925 \quad \log \mathcal{N}(\mathcal{F}, \varepsilon) \leq \log C_0 + d_{\text{eff}} \log\left(\frac{DL_f}{\varepsilon}\right). \quad (20)$$

926 Let $F = L_f D$ be the uniform bound on $|f(v)|$.
927 Dudley's entropy integral gives

$$928 \quad \mathcal{R}_m(\mathcal{F}) \leq \frac{12}{\sqrt{m}} \int_0^F \sqrt{\log \mathcal{N}(\mathcal{F}, \varepsilon)} d\varepsilon. \quad (21)$$

929 Substituting the covering number bound and sim-
930 plifying (absorbing $\log C_0$ into the constant factor
931 for asymptotic behavior), we obtain:

$$932 \quad \mathcal{R}_m(\mathcal{F}) \leq \frac{12L_f}{\sqrt{m}} \int_0^D \sqrt{d_{\text{eff}} \log\left(\frac{D}{t}\right)} dt, \quad (22)$$

933 where we utilized the substitution $t = \varepsilon/L_f$. Let
934 $I = \int_0^D \sqrt{\log(D/t)} dt$. Using the change of vari-
935 ables $u = \log(D/t)$, we have $t = De^{-u}$, which
936 yields $I = D \int_0^\infty u^{1/2} e^{-u} du = D \cdot \Gamma(3/2) = D \frac{\sqrt{\pi}}{2}$.
937 Substituting this back yields:

$$938 \quad \mathcal{R}_m(\mathcal{F}) \leq \frac{12L_f \sqrt{d_{\text{eff}}}}{\sqrt{m}} \left(\frac{D\sqrt{\pi}}{2} \right) = \frac{6\sqrt{\pi} L_f D \sqrt{d_{\text{eff}}}}{\sqrt{m}}. \quad (23)$$

939 By setting $C_{\text{Rad}} = 6\sqrt{\pi}$, we recover the bound in
940 (16). \square

941 **A.4 Architectural stability of the flow**

942 We now turn to the flow T_θ . Let $J_\ell(v)$ be the Ja-
943 cobian of the ℓ -th invertible block at input v , and
944 let $J_{\text{tot}}(v)$ denote the input–output Jacobian of T_θ .
945 The Lipschitz constant is defined as

$$946 \quad \text{Lip}(T_\theta) = \sup_v \|J_{\text{tot}}(v)\|_2. \quad (24)$$

947 **Lemma 2** (Architectural stability via Zero-Ini-
948 tialization). *Consider the flow T_θ composed of L*
949 *blocks. Under the Zero-Initialization strategy, the*
950 *Jacobian of the ℓ -th block takes the form $J_\ell = I + \Delta_\ell$.*
951 *Let $\bar{\sigma} = \frac{1}{L} \sum_{\ell=1}^L \mathbb{E}[\|\Delta_\ell\|_2]$ be the mean spectral*
952 *norm of the perturbations. Consistent with the*
953 *conservative estimation used in our implementa-*
954 *tion, the expected Lipschitz constant of the flow is*
955 *bounded by:*

$$956 \quad \mathbb{E}[\text{Lip}(T_\theta)] \leq 1 + L \cdot \bar{\sigma}. \quad (25)$$

957 *Proof.* The total Jacobian J_{tot} is the product of
958 layer-wise Jacobians:

$$959 \quad J_{\text{tot}} = \prod_{\ell=1}^L J_\ell = \prod_{\ell=1}^L (I + \Delta_\ell). \quad (26)$$

960 We aim to bound the Lipschitz constant $\text{Lip}(T_\theta) =$
961 $\|J_{\text{tot}}\|_2$. Recall that $J_{\text{tot}} = \prod_{\ell=1}^L (I + \Delta_\ell)$. Perform-
962 ing a first-order expansion of this product yields:

$$963 \quad J_{\text{tot}} \approx I + \sum_{\ell=1}^L \Delta_\ell. \quad (27)$$

964 Applying the triangle inequality separates the iden-
965 tity transformation from the perturbations:

$$966 \quad \|J_{\text{tot}}\|_2 \approx \left\| I + \sum_{\ell=1}^L \Delta_\ell \right\|_2 \leq \|I\|_2 + \sum_{\ell=1}^L \|\Delta_\ell\|_2. \quad (28)$$

967 Noting that $\|I\|_2 = 1$, we take the expectation:

$$968 \quad \mathbb{E}[\|J_{\text{tot}}\|_2] \leq 1 + \sum_{\ell=1}^L \mathbb{E}[\|\Delta_\ell\|_2] \quad (29)$$

$$969 \quad = 1 + L \cdot \left(\frac{1}{L} \sum_{\ell=1}^L \mathbb{E}[\|\Delta_\ell\|_2] \right). \quad (30)$$

970 Substituting the definition of the mean spectral
971 norm $\bar{\sigma}$, we directly obtain:

$$972 \quad \mathbb{E}[\text{Lip}(T_\theta)] \leq 1 + L \cdot \bar{\sigma}. \quad (31)$$

973 This confirms that under the small-perturbation
974 regime enforced by zero-initialization, the Lips-
975 chitz constant grows linearly with depth L , aligning
976 with the conservative bound used in our implemen-
977 tation. \square

978 A.5 Finite-sample generalization bound

979 We now derive a finite-sample bound that matches
980 the decomposition used in the main report. For
981 clarity, we distinguish the training and validation
982 samples

$$983 S_{\text{train}} = \{v_1^{\text{train}}, \dots, v_m^{\text{train}}\}, \quad S_{\text{val}} = \{v_1^{\text{val}}, \dots, v_{m_{\text{val}}}^{\text{val}}\}, \quad (32)$$

984 and define the corresponding empirical risks

$$985 \hat{L}_{\text{train}}(\theta) = \frac{1}{m} \sum_{i=1}^m \ell_{\theta}(v_i^{\text{train}}). \quad (33)$$

$$986 \hat{L}_{\text{val}}(\theta) = \frac{1}{m_{\text{val}}} \sum_{j=1}^{m_{\text{val}}} \ell_{\theta}(v_j^{\text{val}}). \quad (34)$$

987 The loss class is

$$988 \mathcal{L} = \{\ell_{\theta}(\cdot) = -\log p_{\theta}(\cdot) : \theta \in \Theta\}. \quad (35)$$

989 By Assumption 2 and Lemma 2, the loss functions
990 are Lipschitz with respect to v . Using the architec-
991 tural bound derived in Lemma 2 (which establishes
992 linear growth with depth L), the Rademacher com-
993 plexity of \mathcal{L} admits the bound:

$$994 \mathcal{R}_m(\mathcal{L}) \leq \frac{\tilde{C}_{\text{Rad}} L \bar{\sigma} \sqrt{d_{\text{eff}}}}{\sqrt{m}}, \quad (36)$$

995 where \tilde{C}_{Rad} collects the constants from Lemma 1
996 and the affine identity term. Note that this scales
997 linearly with L , consistent with our implementa-
998 tion.

999 We now control the difference between training
1000 and validation risks. Introduce the population risk

$$1001 L(\theta) = \mathbb{E}_{v \sim \mathcal{D}}[\ell_{\theta}(v)]. \quad (37)$$

1002 By the triangle inequality,

$$1003 \left| \hat{L}_{\text{val}}(\theta) - \hat{L}_{\text{train}}(\theta) \right| \leq \left| L(\theta) - \hat{L}_{\text{train}}(\theta) \right| + \left| \hat{L}_{\text{val}}(\theta) - L(\theta) \right|. \quad (38)$$

1004 The two terms on the right-hand side are treated
1005 separately.

1006 **Training to population.** A standard Rademacher
1007 complexity bound (see for example (Bartlett and
1008 Mendelson, 2002)) states that if $|\ell_{\theta}(v)| \leq M_{\text{train}}$
1009 for all $\theta \in \Theta$ and all v in the support of the training
1010 distribution, then for any $\delta \in (0, 1)$, with probabil-
1011 ity at least $1 - \delta$ over the draw of S_{train} one has

$$1012 \left| L(\theta) - \hat{L}_{\text{train}}(\theta) \right| \leq 2\mathcal{R}_m(\mathcal{L}) + 3M_{\text{train}} \sqrt{\frac{\log(2/\delta)}{2m}}. \quad (39)$$

Combining (39) with (36) gives

$$\left| L(\theta) - \hat{L}_{\text{train}}(\theta) \right| \leq \frac{2\tilde{C}_{\text{Rad}} L \bar{\sigma} \sqrt{d_{\text{eff}}}}{\sqrt{m}} + 3M_{\text{train}} \sqrt{\frac{\log(2/\delta)}{2m}}. \quad (40)$$

1015 **Validation to population.** For the validation set,
1016 the model parameter θ is fixed, so we only need a
1017 concentration bound for bounded random variables.
1018 Let $X_j = \ell_{\theta}(v_j^{\text{val}})$ with $X_j \in [a, b]$ and define
1019 $M_{\text{val}} = b - a$. Hoeffding's inequality (Hoeffding,
1020 1994) yields

$$1021 \mathbb{P}\left(\left|\hat{L}_{\text{val}}(\theta) - L(\theta)\right| \geq t\right) \leq 2 \exp\left(-\frac{2m_{\text{val}}t^2}{M_{\text{val}}^2}\right). \quad (41)$$

1022 Setting the right-hand side to δ and solving for t
1023 gives that, with probability at least $1 - \delta$,

$$\left| \hat{L}_{\text{val}}(\theta) - L(\theta) \right| \leq M_{\text{val}} \sqrt{\frac{\log(2/\delta)}{2m_{\text{val}}}}. \quad (42)$$

1025 **Final bound.** Combining (38), (40) and (42), and
1026 applying a union bound, we obtain the following
1027 result.

1028 **Theorem 1** (Finite-sample generalization bound).
1029 *Under Assumptions 1 and 2, for any fixed flow*
1030 *model p_{θ} and any $\delta \in (0, 1)$, with probability*
1031 *at least $1 - \delta$ over the draws of the training set*
1032 *of size m and the validation set of size m_{val} , the*
1033 *train-validation gap satisfies*

$$\left| \hat{L}_{\text{val}}(\theta) - \hat{L}_{\text{train}}(\theta) \right| \leq \frac{2\tilde{C}_{\text{Rad}} L \bar{\sigma} \sqrt{d_{\text{eff}}}}{\sqrt{m}} + M_{\text{val}} \sqrt{\frac{\log(2/\delta)}{2m_{\text{val}}}} + 3M_{\text{train}} \sqrt{\frac{\log(2/\delta)}{2m}}. \quad (43)$$

1035 A.6 Extension to the IS score

1036 Let L_{marg} and L_{cond} denote the population loss of
1037 the marginal and conditional flows, and let \hat{L}_{marg}
1038 and \hat{L}_{cond} be the corresponding validation losses.
1039 The population I_S score is

$$1040 IS = L_{\text{marg}}(V) - L_{\text{cond}}(V|U), \quad (44)$$

1041 and the empirical I_S score is

$$1042 \widehat{IS} = \hat{L}_{\text{marg}}(V) - \hat{L}_{\text{cond}}(V|U). \quad (45)$$

The difference is

$$IS - \widehat{IS} = (L_{\text{marg}} - \hat{L}_{\text{marg}}) - (L_{\text{cond}} - \hat{L}_{\text{cond}}). \quad (46)$$

By the triangle inequality,

$$|IS - \widehat{IS}| \leq |L_{\text{marg}} - \hat{L}_{\text{marg}}| + |L_{\text{cond}} - \hat{L}_{\text{cond}}|. \quad (47)$$

Applying Theorem 1 separately to the marginal and conditional flows yields a finite-sample upper bound for the generalization error of the I_s estimator as the sum of the two individual generalization gaps.

B Experiment Detail

In this section, we will provide the necessary experimental details to reproduce these experiments.

B.1 Evaluated Model and Datasets Details

In Table 3, we provide the metadata of the evaluated models and their score on the 11 datasets. We provide the statistics of the datasets used to evaluate \bar{I}_s in table 4.

Model	Dim.	\bar{I}_s
Zeta_Alpha_E5_Mistral	4096	0.20
Linq_Embed_Mistral	4096	0.20
SFR_Embedding_Mistral	4096	0.19
bge_multilingual_gemma2	3584	0.19
GritLM_7B	4096	0.18
gte_Qwen2_7B_instruct	3584	0.17
all_MiniLM_L6_v2	768	0.13
stella_base_en_v2	384	0.13

Table 3: Information sufficiency of the evaluated models by FLARE .

B.2 Downstream Task Evaluation

We deliberately select datasets that are either newly released or underexplored to minimize the risk of data leakage during embedding model pretraining. As no established performance benchmarks for embedding models exist on these datasets, we evaluate downstream task performance ourselves following the MTEB evaluation (Muennighoff et al., 2023) protocol: F1 Macro for classification, Spearman for STS, nDCG@10 for retrieval, and V-measure for clustering.

Table 5 reports the average downstream task performance of the eight evaluated embedding models, aggregated by task type. Overall, 7B-scale

instruction-tuned models (GritLM, SFR, Linq, Zeta-Alpha) consistently outperform smaller models across all task categories. GritLM-7B achieves the best classification performance (0.61), while SFR-Embedding-Mistral leads on STS (0.67). The two smaller models, stella-base-en-v2 and all-MiniLM-L6-v2, show competitive performance on classification but lag significantly on retrieval and STS tasks.

Notably, retrieval scores exhibit the largest variance, partly due to the inclusion of challenging benchmarks such as LIMIT, which is specifically designed to probe the upper limits of embedding model capabilities and thus yields lower absolute scores across all models.

B.3 Comprehensive Results

The primary goal of our method is to *rank* candidate embedding models so that practitioners can select the best one for a given downstream task without access to labeled data. The correlation coefficients reported in Table 6 validate that the rankings produced by FLARE align with ground-truth task performance.

We compare FLARE against four unsupervised baselines: Uniformity, IsoScore, Silhouette Score, and EMIR. Across 11 datasets spanning classification, STS, retrieval, and clustering tasks, FLARE achieves the highest average Spearman correlation ($\rho = 0.69$), substantially outperforming all baselines. FLARE is the only method that maintains consistently positive correlations across all datasets, demonstrating robust ranking capability regardless of task type.

EMIR achieves an average Spearman correlation of only -0.12 , indicating that its rankings frequently contradict ground-truth performance. We believe this might be due to the fact that the nuclear method performs poorly in high-dimensional situations. On individual datasets, EMIR shows high variance: while achieving moderate positive correlations on some tasks (e.g., $\rho = 0.33$ on LivNLP-STS), it produces strongly negative correlations on others (e.g., $\rho = -0.52$ on FunPang, $\rho = -0.43$ on arXiv’25). This instability limits its practical utility for model selection.

IsoScore exhibits consistently poor performance with an average of $\rho = -0.27$. These results confirm that FLARE provides the most reliable unsupervised signal for embedding model selection.

Dataset	Task	Train	Val	Total
apt-eval	Classification	13,185	1,465	14,650
gtfintechlab	Classification	12,250	2,625	14,875
BhashaBench-Finance	Classification	12,105	1,346	13,451
Augmented-stsb	STS	33,635	3,738	37,373
C-STIS-Reannotated	STS	12,758	1,418	14,176
Philosophical-STIS	STS	58,560	6,507	65,067
AIR-Bench	Retrieval	23,639	2,627	26,266
LIMIT	Retrieval	45,000	5,000	50,000
arXiv-abstracts	Retrieval	2,610	290	2,900
clustered-FUNPANG	Clustering	28,716	3,191	31,907
deduplicated-reasoning	Clustering	50,772	5,642	56,414

Table 4: Statistics of the datasets used in our experiments.

Model	Classification	STS	Retrieval	Clustering
bge_multilingual_gemma2	0.58	0.63	0.46	0.30
Zeta_Alpha_E5_Mistral	0.57	0.65	0.53	0.31
GritLM_7B	0.61	0.64	0.50	0.32
SFR_Embedding_Mistral	0.57	0.67	0.53	0.26
Linq_Embed_Mistral	0.60	0.66	0.52	0.32
gte_Qwen2_7B_instruct	0.57	0.64	0.50	0.28
stella_base_en_v2	0.49	0.45	0.05	0.28
all_MiniLM_L6_v2	0.50	0.51	0.47	0.21

Table 5: Summary of the evaluated embedders and their performance on downstream datasets.

C Ranking Stability under Subsampling

We evaluate whether the proposed FLARE yields stable model rankings when the evaluation set is randomly subsampled. For each dataset, we subsample the evaluation set at ratios $\alpha \in \{5\%, 10\%, 20\%, 40\%, 60\%, 80\%, 100\%\}$ without replacement, and repeat this process 20 times. For each α , we recompute scores using the same pre-trained marginal and conditional flows and obtain a ranking of embedding models. We then compute the Spearman rank correlation $\rho(\alpha)$ between the ranking induced by the subsampled evaluation set and the reference ranking computed on the full evaluation set ($\alpha = 1.0$), and report the deviation

$$\Delta_\rho(\alpha) = |\rho(\alpha) - \rho(1.0)|. \quad (48)$$

We focus on Spearman correlation because our goal is to assess the stability of *model ranking* for model selection, rather than the linear agreement of raw scores. Rank-based measures directly quantify whether the relative ordering of models is preserved

under subsampling, which is the main quantity of interest in this analysis.

Figure 4 shows the ranking deviation $\Delta_\rho(\alpha) = |\rho_\alpha - \rho_{\text{full}}|$ as a function of the subsampling ratio α across all 11 datasets, grouped by task type. Overall, I_s rankings remain highly stable under subsampling: for 8 out of 11 datasets, $\Delta_\rho < 0.05$ even when using only 20% of the validation data.

Among task types, **clustering** exhibits the highest stability, with all three datasets maintaining $\Delta_\rho < 0.025$ across all subsampling ratios. **Classification** and **retrieval** tasks also demonstrate strong robustness, with most datasets showing only minor deviations ($\Delta_\rho < 0.045$) even at very low α .

In contrast, **STS** datasets display larger variance at small sample sizes—*Aug-stsb* shows the highest deviation ($\Delta_\rho \approx 0.12$) at $\alpha = 0.05$, consistent with correlation-based evaluation being more sample-sensitive. However, these deviations diminish rapidly: once $\alpha \geq 0.2$, all STS datasets achieve $\Delta_\rho < 0.07$.

These results demonstrate that I_s selection is

Dataset	Task	Spearman’s ρ					Pearson’s r				
		Uni.	Iso.	Sil.	EMIR	Ours	Uni.	Iso.	Sil.	EMIR	Ours
apt-eval	Class.	0.36	-0.24	-0.21	0.07	0.43	0.50	-0.70	-0.42	-0.01	0.20
gtfintechlab	Class.	-0.12	-0.50	0.00	-0.38	0.43	0.11	-0.42	-0.18	-0.46	0.14
BhashaBench	Class.	0.29	-0.45	0.02	0.12	0.81	0.44	-0.62	0.35	-0.12	0.59
Aug.-stsb	STS	0.14	0.05	0.24	-0.12	0.83	0.51	0.25	0.06	-0.11	0.68
LivNLP-STS	STS	-0.18	-0.27	0.81	0.33	0.83	0.77	-0.41	0.83	0.39	0.94
Philo-STS	STS	0.07	-0.76	0.62	-0.38	0.43	0.82	-0.26	0.70	-0.51	0.49
AIR-Bench	Retr.	-0.14	-0.79	0.05	-0.24	0.71	0.70	-0.27	0.45	0.07	0.69
arXiv ’25	Retr.	0.55	0.14	-0.71	-0.43	0.64	0.72	-0.27	0.10	-0.33	0.62
LIMIT	Retr.	-0.17	-0.71	0.57	0.02	0.81	-0.14	-0.57	0.62	-0.28	0.62
FunPang	Clust.	-0.14	0.50	-0.48	-0.52	0.90	-0.52	0.82	-0.75	-0.57	0.83
Reasoning	Clust.	-0.14	0.07	0.00	0.21	0.76	-0.33	-0.73	0.26	0.37	0.55
Average		0.05	-0.27	0.08	-0.12	0.69	0.33	-0.29	0.18	-0.14	0.58

Table 6: Comparisons with unsupervised baselines. FLARE achieves the highest consistency and average correlation.

robust to evaluation subsampling, with 20–40% of validation data typically sufficient for reliable model ranking across diverse task types.

D Ablation Study

In this section, we conduct a set of ablation studies to better understand the factors that contribute to the effectiveness of our method. Unless otherwise specified, all experiments are conducted using the same evaluation protocol and datasets as in the main experiments.

Appendix D.1. Experimental setup.

For all ablation experiments, we use the same pre-trained models and evaluation datasets as in the main experiments. Unless otherwise stated, we recompute the evaluation scores under modified settings while keeping all other components unchanged. Performance is measured using Spearman correlation between the predicted ranking and the ground-truth ranking.

D.1 Shuffle ablation.

To examine whether the proposed metric truly relies on the correspondence between paired representations, we conduct a shuffle-based ablation with varying shuffle ratios. Specifically, for a given ratio $p \in \{0, 0.1, 0.2, 0.4, 0.6, 0.8, 1.0\}$, we randomly select a p fraction of the evaluation samples

and permute the correspondence between U and V within this subset, while keeping the remaining $(1 - p)$ portion unchanged. This procedure preserves the marginal distributions of both U and V but progressively destroys their pairwise alignment as p increases.

For each shuffle ratio, we recompute the proposed score and evaluate the resulting ranking against the ground-truth downstream performance. As shown in Table 7, the Spearman correlation generally degrades as p increases across datasets. At $p = 0$, all datasets exhibit positive correlations (avg. $\rho = 0.69$); at $p = 1.0$, correlations become predominantly negative (avg. $\rho \approx -0.40$). The transition point varies by task type: Classification and STS datasets tend to flip sign at lower shuffle ratios ($p \approx 0.1$), while Retrieval and Clustering datasets maintain positive correlations longer (up to $p \approx 0.2$ – 0.4). This behavior confirms that the proposed metric critically depends on correct U - V alignment rather than marginal statistics alone.

D.2 Full I_s vs Conditional-Only.

To investigate the contribution of the marginal term in the I_s , we conduct an ablation study comparing the full I_s formulation against a conditional-only variant that omits the marginal likelihood component.

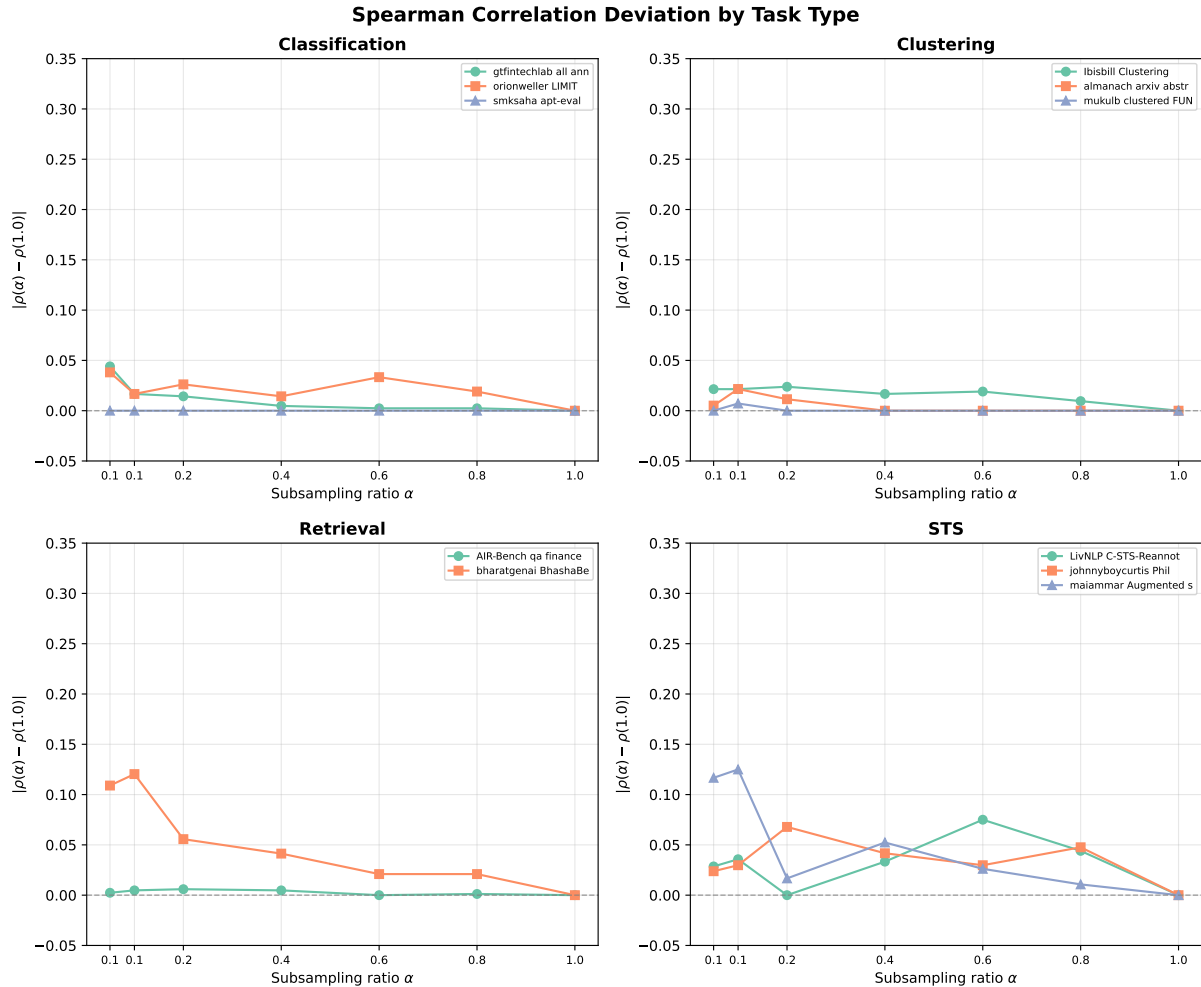


Figure 4: **Ranking stability under evaluation subsampling.** We report the deviation $\Delta_\rho(\alpha) = |\rho(\alpha) - \rho(1.0)|$, where $\rho(\alpha)$ is the Spearman rank correlation between the model ranking induced by IS scores computed on a subsampled evaluation set (ratio α) and the ranking computed on the full evaluation set ($\alpha = 1.0$). Smaller values indicate more stable rankings.

Dataset	Task	$p = 0$	$p = 0.1$	$p = 0.2$	$p = 0.4$	$p = 0.6$	$p = 0.8$	$p = 1.0$
apt-eval	Cls	0.43	-0.36	0.48	-0.48	-0.43	-0.48	-0.48
gtfintech	Cls	0.43	-0.19	-0.24	-0.62	-0.69	-0.74	-0.64
BhashaBench	Cls	0.81	-0.74	-0.62	-0.50	-0.50	-0.43	-0.26
Aug-STs	STS	0.83	-0.83	-0.81	-0.81	-0.93	-0.57	-0.71
LivNLP-STs	STS	0.83	-0.07	-0.10	-0.55	-0.52	-0.52	-0.43
Philo-STs	STS	0.43	-0.57	-0.48	-0.45	-0.40	-0.36	-0.24
AIR-Bench	Ret	0.71	0.76	0.24	0.69	-0.02	0.02	-0.55
arXiv '25	Ret	0.64	0.71	0.36	0.64	0.43	0.50	0.48
LIMIT	Ret	0.81	0.31	0.29	-0.81	-0.67	-0.36	-0.48
FunPang	Clust	0.90	0.93	0.67	-0.83	-0.90	-0.90	-0.90
Reasoning	Clust	0.76	0.26	-0.43	-0.07	-0.02	-0.21	-0.24

Table 7: Per-dataset Spearman correlation (ρ) under partial shuffle ablation. The column $p = 0$ corresponds to the full method without shuffling. As shuffle proportion p increases, correlation with downstream performance generally degrades.

Results. Table 8 and Figure 5 present the comparison between Full I_s and Cond Only across all 11 datasets. The full I_s formulation achieves an average Spearman correlation of $\rho = 0.69$, substantially outperforming the Cond Only variant ($\rho = 0.21$). Notably, Full I_s outperforms Cond Only on 8 out of 11 datasets. The advantage is particularly pronounced on retrieval and clustering tasks: on LIMIT, Full I_s achieves $\rho = 0.81$ compared to -0.21 for Cond Only; on FunPang, the gap is 0.90 vs. -0.36 . These results indicate that the condition-only alone captures partial information about embedding quality. It measures how well the source embedding predicts the target, but fails to account for the intrinsic structure of the target embedding space.

Interestingly, on a few datasets (apt-eval, arXiv '25, gtfintechlab), Cond Only slightly outperforms Full I_s . However, its performance is highly inconsistent, with five datasets showing negative correlations. In contrast, Full I_s maintains positive correlations across all datasets, demonstrating greater robustness.

These findings confirm that both terms in the I_s formulation are necessary: the marginal term captures target-side embedding quality, while the conditional term measures cross-model information transfer. Their combination yields a more reliable and consistent signal for unsupervised model selection.

Dataset	Full IS	Cond Only
apt-eval	0.43	0.64
LIMIT	0.81	-0.21
Aug.-stsb	0.83	0.67
LivNLP-STS	0.83	0.12
FunPang	0.90	-0.36
Philo-STS	0.43	-0.12
AIR-Bench	0.71	-0.07
arXiv '25	0.64	0.71
BhashaBench	0.81	0.64
gtfintechlab	0.43	0.50
Reasoning	0.76	-0.21
Average	0.69	0.21

Table 8: Ablation comparing Full I_s with Cond Only. Full I_s achieves substantially higher correlation with ground truth rankings (avg. $\rho = 0.69$ vs. 0.21), confirming the importance of the marginal term.

D.3 Aggregation strategy.

We investigate the impact of different aggregation strategies for combining IS scores into a single model score. We compare three methods: (1) arithmetic mean, (2) median, and (3) 10% trimmed mean (Trim10), which discards the top and bottom 10% of values before averaging.

As shown in Table 9, median aggregation consistently outperforms alternatives, achieving the highest Spearman correlation on 9 of 11 datasets (average $\rho = 0.69$ vs. 0.52 for mean, a relative improvement of 32.7%). For Pearson correlation, median also leads with an average of $r = 0.56$ compared to 0.54 for mean and 0.53 for trimmed mean.

The advantage of median aggregation is particularly pronounced on STS tasks, where it achieves $\rho = 0.83$ on both Aug.-stsb and LivNLP-STS, compared to $\rho \leq 0.45$ for mean. This robustness stems from the median’s insensitivity to outlier model pairs that may exhibit anomalously high or low IS scores due to training instabilities or distribution mismatches between certain embedding spaces.

Exceptions include Philo-STS and arXiv '25, where mean or trimmed mean aggregation outperforms the median. For Philo-STS, mean aggregation ($\rho = 0.62$) surpasses median ($\rho = 0.43$), likely due to the smaller dataset size reducing the prevalence of outlier scores. Similarly, on arXiv '25, trimmed mean achieves the highest correlation ($\rho = 0.84$), suggesting that while some outliers exist, the distribution tails might contain valuable signal for this specific retrieval task. Despite these cases, median aggregation remains the most robust strategy overall.

Based on these findings, we adopt median aggregation as the default strategy throughout our experiments.

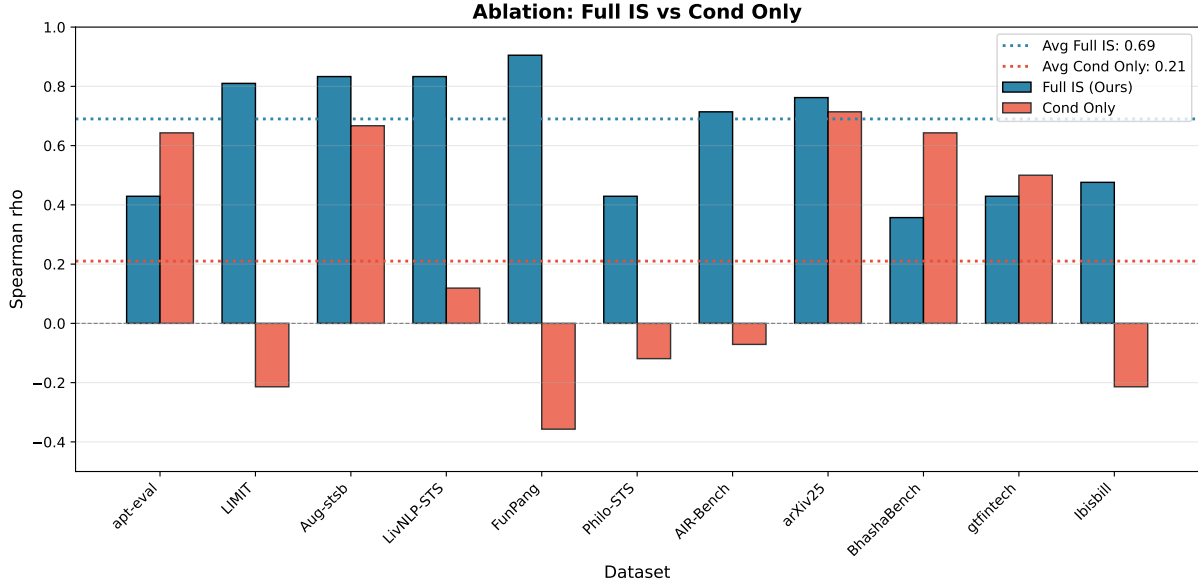


Figure 5: Ablation: Comparison of correlation of Full I_s vs. conditional-only component.

Dataset	Task	Spearman ρ			Pearson r		
		Mean	Median	Trim10	Mean	Median	Trim10
apt-eval	Class.	0.29	0.43	0.29	0.26	0.20	0.16
gtfintechlab	Class.	0.38	0.43	0.17	0.36	0.14	-0.05
BhashaBench	Class.	0.53	0.81	0.36	0.69	0.59	0.45
Aug.-stsb	STS	0.40	0.83	0.69	0.54	0.68	0.63
LivNLP-ST5	STS	0.45	0.83	0.52	0.55	0.94	0.53
Philo-ST5	STS	0.62	0.43	0.52	0.66	0.49	0.63
AIR-Bench	Retr.	0.43	0.71	0.48	0.52	0.70	0.80
arXiv '25	Retr.	0.70	0.64	0.84	0.25	0.46	0.56
LIMIT	Retr.	0.62	0.81	0.81	0.61	0.62	0.65
FunPang	Clust.	0.81	0.90	0.88	0.95	0.83	0.93
Reasoning	Clust.	0.52	0.76	0.60	0.56	0.55	0.57
Average		0.52	0.69	0.56	0.54	0.56	0.53

Table 9: Aggregation ablation comparing mean, median, and trimmed mean (10%). The median aggregation (Ours) achieves the highest consistency ($\rho = 0.69$).