

TabEmbed: Benchmarking and Learning Generalist Embeddings for Tabular Understanding

Anonymous ACL submission

Abstract

Foundation models have established unified representations for natural language processing, yet this paradigm remains largely unexplored for tabular data. Existing methods face fundamental limitations: LLM-based approaches lack retrieval-compatible vector outputs, whereas text embedding models often fail to capture tabular structure and numerical semantics. To bridge this gap, we first introduce the Tabular Embedding Benchmark (TabBench), a comprehensive suite designed to evaluate the tabular understanding capability of embedding models. We then propose TabEmbed, the first generalist embedding model that unifies tabular classification and retrieval within a shared embedding space. By reformulating diverse tabular tasks as semantic matching problems, TabEmbed leverages large-scale contrastive learning with positive-aware hard negative mining to discern fine-grained structural and numerical nuances. Experimental results on TabBench demonstrate that TabEmbed significantly outperforms state-of-the-art text embedding models, establishing a new baseline for universal tabular representation learning.

1 Introduction

Recently, foundation models have achieved remarkable success in establishing universal representations for Natural Language Processing (Wang et al., 2024; Yang et al., 2025), such as Retrieval-Augmented Generation (RAG) (Qiang et al., 2025), where dense text embeddings enable efficient semantic search through vector similarity computation. However, this unified representation paradigm has not been effectively adapted to tabular data. Existing research (Ye et al., 2025; Qu et al., 2025; Mueller et al., 2025) typically treats tabular classification and retrieval as distinct problems requiring specialized models. Consequently, the tabular domain lacks a shared embedding space capable of

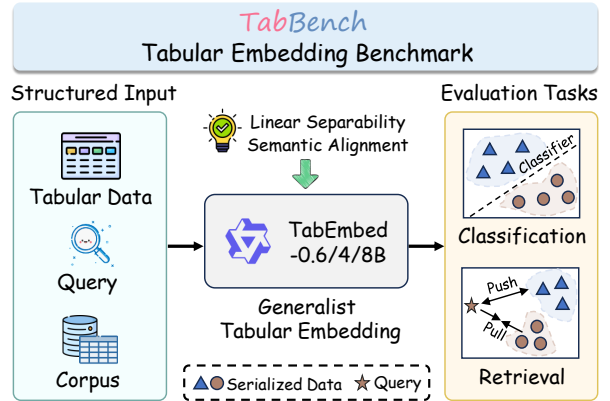


Figure 1: Overview of **TabBench** and **TabEmbed**.

simultaneously addressing all tabular understanding tasks without task-specific architectures.

Traditional tree-based models excel at tabular classification tasks but are constrained by fixed schemas, rendering them incompatible with zero-shot transfer and retrieval scenarios. Recent advances in large language models have shown considerable promise for tabular tasks (Gardner et al., 2024; Ye et al., 2025; Qu et al., 2025). However, these methods do not produce the dense, fixed-dimensional vectors required for vector databases and downstream retrieval applications. While general-purpose text embedding models (Zhang et al., 2025a; Yu et al., 2025; Zhang et al., 2025b) can generate such embeddings with remarkable success in text domains, they treat serialized tables as unstructured text, often failing to capture essential structural logic such as numerical magnitude and column-specific semantics. These constraints motivate the development of a generalist tabular embedding model that inherently understands tabular structure to handle various tabular understanding tasks within a shared embedding space.

However, training such a tabular embedding model presents three significant challenges. First, the absence of benchmarks specifically designed

for tabular embeddings hinders systematic evaluation of embedding models, especially given that tabular data often exhibits fine-grained distinctions where rows differing in only a few critical columns may belong to opposite classes. Second, the tabular domain is constrained by the scarcity of large-scale training data, as tabular datasets rarely contain natural language triplets that explicitly describe numerical or logical constraints for contrastive learning. Finally, unifying classification and retrieval within a shared embedding space is non-trivial. Retrieval relies on semantic ranking to identify relevant data, whereas classification requires precise decision boundaries for label prediction.

To address these challenges, as shown in Figure 1, we first introduce Tabular Embedding Benchmark (TabBench), a comprehensive evaluation suite designed to assess numerical reasoning and retrieval capabilities across diverse structural complexities. Then we propose TabEmbed, an embedding model that unifies classification and retrieval within a shared embedding space. By reformulating diverse tabular tasks as semantic matching problems, TabEmbed enables large-scale contrastive learning on the T4 corpus. To ensure fine-grained discrimination, we further employ positive-aware hard negative mining that compels the model to distinguish subtle schema differences. Extensive experiments in our proposed TabBench demonstrate that TabEmbed significantly outperforms state-of-the-art text embedding models. Our work establishes a new baseline for tabular understanding, demonstrating the effectiveness of generalist embedding models in the tabular domain.

2 Related Work

2.1 Embedding Models

Text embedding research has evolved from encoder-based architectures to Large Language Models (LLMs). Early works adapted BERT (Devlin et al., 2019) and T5 (Raffel et al., 2020) via contrastive learning (Reimers and Gurevych, 2019; Gao et al., 2021), with recent models like E5 (Wang et al., 2022), BGE (Xiao et al., 2024), and GTE (Li et al., 2023) achieving scalability through multi-stage training. However, limited capacity for complex schemas prompted a shift toward decoder-only LLMs. Approaches like SGPT (Muennighoff, 2022) and E5-Mistral (Wang et al., 2024) demonstrated the efficacy of generative backbones, while subsequent innovations (Lee et al., 2024; Zhang

et al., 2025b) further optimized bidirectional information flow. Despite these advancements, existing generalist models typically process serialized tables as unstructured text, lacking the specific optimization for structural reasoning and numerical understanding required for precise tabular retrieval.

2.2 Benchmarks for Text and Tabular Tasks

Current evaluation protocols remain bifurcated between unstructured text and supervised tabular classification. In the text domain, standards like BEIR (Thakur et al., 2021) and MTEB (Muennighoff et al., 2023) lack dedicated structured data scenarios. Conversely, tabular benchmarks such as OpenML-CC18 (Bischl et al., 2017) and Grinsztajn (Grinsztajn et al., 2022) focus primarily on comparing decision trees against neural networks on fixed classification splits. While recent massive corpora (Eggert et al., 2023; Gardner et al., 2024) have enabled transfer learning suites like TabZilla (McElfresh et al., 2023) and GTL (Wen et al., 2024), these remain rooted in generative or predictive paradigms. To bridge this critical gap, we introduce TabBench for comprehensive evaluation and propose TabEmbed, a generalist model that unifies these diverse tabular tasks.

3 The Tabular Embedding Benchmark

To rigorously evaluate the capabilities of embedding models in tabular understanding, we introduce the **Tabular Embedding Benchmark (TabBench)**. Building upon the high-quality data curation of the tabula-8b-eval-suite (Gardner et al., 2024), TabBench provides a comprehensive framework to assess two critical dimensions of tabular representation: linear separability (via classification) and semantic alignment (via retrieval). The benchmark aggregates diverse datasets from four authoritative repositories: **Grinsztajn** (Grinsztajn et al., 2022), **OpenML-CC18** (Bischl et al., 2017), **OpenML-CTR23** (Fischer et al., 2023), and **UniPredict** (Wang et al., 2023). The detailed composition of TabBench is illustrated in Figure 2. We implement a standardized pipeline for data serialization, task construction, and quality filtering.

3.1 Data Serialization

Bridging the modality gap between structured tabular data and large language models requires an effective serialization strategy (Hegselmann et al., 2023; Gardner et al., 2024). Formally, let a tabular row be represented as a set of feature-value pairs

$\mathbf{x} = \{(k_j, v_j)\}_{j=1}^M$, where k_j denotes the column header and v_j is the cell value. We define a serialization function $\mathcal{S} : \mathcal{X} \rightarrow \mathcal{T}$ that maps \mathbf{x} to a natural language sequence via concatenation:

$$\mathcal{S}(\mathbf{x}) = \bigoplus_{j=1}^M \text{“The } k_j \text{ is } \tilde{v}_j\text{.”}, \quad (1)$$

where \tilde{v}_j represents the pre-processed value. To balance numerical precision with token efficiency, continuous values are rounded to two decimal places, while categorical values retain their original string representation. We filter out rows that surpass the predefined token limit, aligning with the context constraints of mainstream embedding models.

3.2 Evaluation Tasks

We formulate two distinct tasks to comprehensively evaluate the versatility of the learned embeddings within a shared vector space.

Tabular Classification This task evaluates the linear separability of the embeddings. We construct the evaluation suite by treating each source dataset as an independent classification task, where the original label column serves as the prediction target and each row is serialized as a text instance. To ensure evaluation quality, we apply a strict filtering protocol based on the statistical properties of labels. Specifically, datasets are excluded if the number of unique classes exceeds 50 or the label-to-sample ratio exceeds 0.1. These criteria help identify valid classification targets while filtering out regression-like targets. For qualified datasets, we employ stratified sampling to partition data into training and testing splits, guaranteeing a minimum of two samples per class to mitigate cold-start issues for rare classes. We evaluate performance using a linear probing setup, training a Logistic Regression classifier on top of frozen embeddings.

Tabular Retrieval Unlike classification, which assesses intra-dataset separability, the retrieval task evaluates the model’s ability to align natural language queries with serialized rows across a heterogeneous global corpus. We construct a unified retrieval corpus by aggregating rows from all datasets, capping each dataset’s contribution at 10,000 samples to prevent distribution dominance. To simulate realistic user intent, we propose a *seed-based query generation* pipeline. For a given “seed row” in the corpus, we generate a natural language query q that

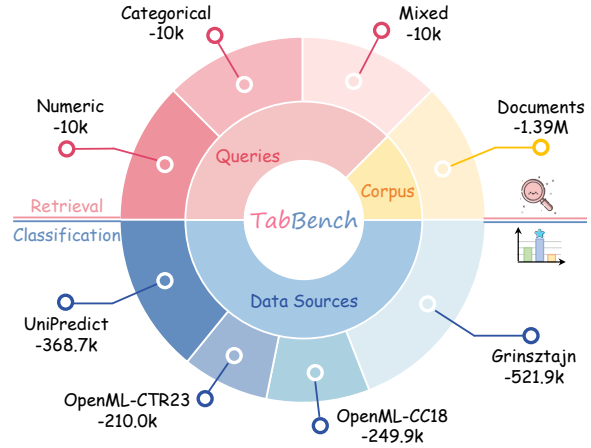


Figure 2: Data composition and statistics of TabBench.

describes a subset of the row’s attributes. The tabular retrieval task is to retrieve the original row (and other matching rows) given query q . The queries follow the template: “*Find records where C_1 and \dots and C_k* ”, where each C_i represents a constraint condition. Based on the type of constraints, we define three query types of increasing complexity:

- **Categorical Queries:** Assess exact-match semantics. Each constraint C_i enforces strict equality on discrete features (e.g., “*Status is Active*”).
- **Numeric Queries:** Test the understanding of magnitude and ranges. Each constraint C_i is generated by sampling an operator $op \in \{>, <, =\}$ and perturbing the original feature value (e.g., “*Price less than 50.25*”).
- **Mixed Queries:** Evaluate complex reasoning by combining numeric and categorical constraints derived from the same row (e.g., “*Status is Active and Price less than 50.25*”).

To ensure benchmark validity, we perform symbolic verification for every generated query, retaining only those that match at least 5 rows in the global corpus. This process yields a balanced evaluation set, where each query contains 1 to 3 conditions, covering diverse query complexities.

4 TabEmbed: Unified Tabular Embedding Learning

To bridge the gap between structured data and semantic representation, we propose **TabEmbed**, a generalist embedding model trained within a unified framework that learns tabular representations by casting disparate downstream tasks into a shared

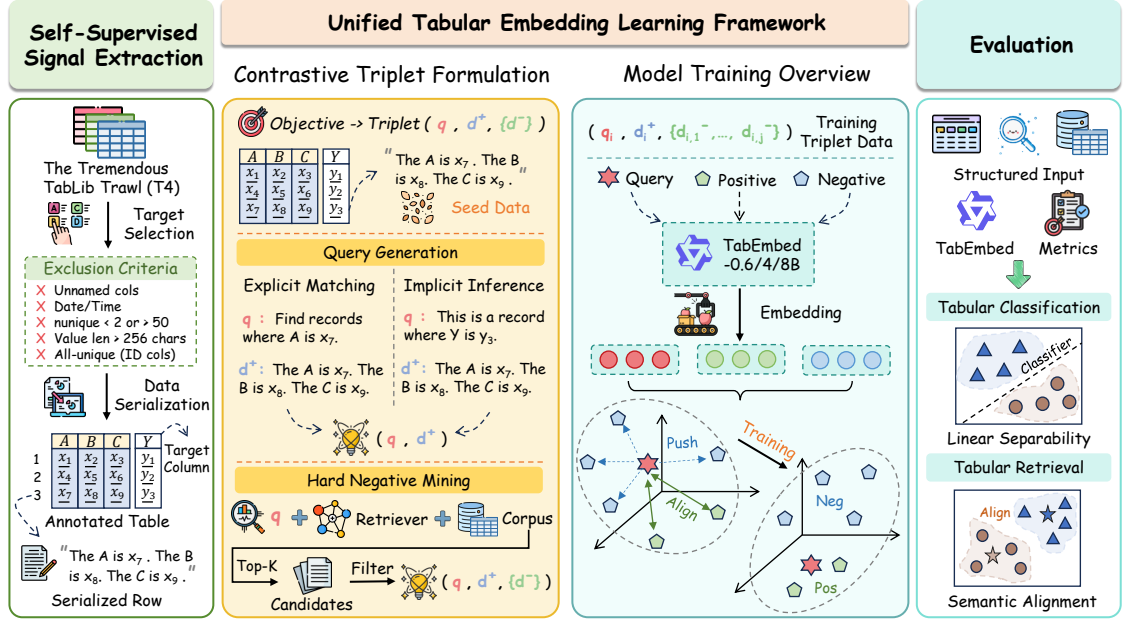


Figure 3: The overall framework of TabEmbed.

contrastive paradigm. The overall framework is illustrated in Figure 3. Leveraging the massive scale of the **T4** dataset (Gardner et al., 2024), we employ a self-supervised contrastive learning approach. By training on diverse contrastive triplets constructed from unannotated tables, we unify disparate downstream capabilities into a shared semantic space.

4.1 Self-Supervised Signal Extraction

Since the T4 corpus lacks explicit task annotations, we employ an automated pipeline to transform raw tables into self-supervised training instances. We first dynamically identify a target column y within each table to serve as the prediction signal. To ensure the quality of these self-supervised signals, we apply a rigorous filtering protocol to exclude non-informative attributes (e.g., identifiers, timestamps) and prioritize targets with clear semantic boundaries. The Detailed pipeline is provided in Appendix F. To prevent information leakage and compel the model to learn latent dependencies, we apply a *target-masked serialization* strategy. Specifically, we strictly exclude the selected target y from the feature set and apply the serialization function S (defined in Section 3.1) to the remaining columns. This yields the serialized row $d = S(\mathbf{x}_{-y})$, ensuring that the embedding captures the row’s semantic content without revealing the ground truth label.

4.2 Contrastive Triplet Formulation

We formulate tabular representation learning as optimizing similarity within triplets $(q, d^+, \{d^-\})$,

where q is a task-specific query, d^+ is the positive serialized row, and $\{d^-\}$ are hard negatives. We construct these components to cover both explicit signal matching and implicit semantic inference.

4.2.1 Task-Adaptive Query Generation

We generate synthetic queries q to model two complementary tasks using a shared data format:

Tabular Retrieval (Explicit Matching) The retrieval task aligns natural language constraints with rows that satisfy them. We leverage the query generation pipeline detailed in Section 3.2, which samples subsets of attributes from \mathbf{x}_{-y} to form logical conditions spanning both numerical and categorical fields. Formally, for a serialized row d^+ , we generate a query q_{ret} describing specific attribute constraints (e.g., “Find records where Status is Active and Price less than 50.25”). This forces the model to align natural language constraints with specific attribute values present in the input.

Tabular Classification (Implicit Inference) The classification task aligns abstract label descriptions with rows that imply those labels. Unlike retrieval, the query content (the value of target y) is absent from the input d^+ and must be inferred solely from the correlations among the remaining features. For a hidden target column y with value v , we construct a descriptive label query q_{cls} (e.g., “This is a record where y is v ”). This formulation encourages the model to cluster rows based on latent predictive features rather than surface-level token overlap.

Table 1: Main results on TabBench. We evaluate TabEmbed against state-of-the-art generalist text embedding models across three parameter scales. The best results within each parameter group are highlighted in **bold**.

Model	#Params	Overall	Classification		Retrieval	
			Accuracy	F1	MRR@10	nDCG@10
Jina-Embeddings-v3	0.6B	41.48	60.33	46.11	32.49	26.98
Jasper-Token-Compression	0.6B	42.75	61.25	47.69	33.56	28.50
Qwen3-Embedding-0.6B	0.6B	44.92	62.81	50.32	36.00	30.56
TabEmbed-0.6B (Ours)	0.6B	65.27	67.16	56.56	71.72	65.64
F2LLM-4B	4B	48.02	64.92	52.48	40.60	34.08
Qwen3-Embedding-4B	4B	48.91	65.09	52.72	42.04	35.76
TabEmbed-4B (Ours)	4B	70.71	69.51	59.75	79.33	74.25
SFR-Embedding-Mistral	7B	49.42	64.28	50.75	44.23	38.41
Linq-Embed-Mistral	7B	50.74	66.06	53.33	44.65	38.92
GTE-Qwen2-7B-Instruct	7B	51.27	64.67	51.76	47.44	41.19
Qwen3-Embedding-8B	8B	48.03	65.08	52.81	40.06	34.16
TabEmbed-8B (Ours)	8B	71.62	69.88	60.19	80.58	75.83

4.2.2 Positive-Aware Hard Negative Mining

Simple in-batch negatives are insufficient for distinguishing numerically similar values or closely related classes. We implement an offline Hard Negative Mining strategy using a lightweight dense retriever (Qwen3-Embedding-0.6B). For every query q , we retrieve the Top- K candidates from the global corpus. Crucially, we employ a Positive-Aware Filtering mechanism: we strictly retain only those candidates that possess high semantic similarity to the query but explicitly violate the retrieval condition or belong to a different class label. These mined hard negatives d^- constitute the set of samples that are most easily confused with the positive d^+ , ensuring the model learns sharp decision boundaries.

4.3 Training Objective

We optimize our model using the contrastive learning loss. Given a batch \mathcal{B} containing N triplets $(q_i, d_i^+, \{d_{i,j}^-\}_{j=1}^M)$, where M is the number of mined hard negatives, the objective for query q_i is defined as:

$$\mathcal{L}_i = -\log \frac{e^{s_i^+/\tau}}{e^{s_i^+/\tau} + \sum_{j \in \mathcal{N}_i} e^{s_i^j/\tau}}, \quad (2)$$

where $s_i^+ = \text{sim}(q_i, d_i^+)$ is the positive similarity, \mathcal{N}_i includes both M hard negatives and in-batch negatives from other queries in \mathcal{B} , $s_i^j = \text{sim}(q_i, d_j)$ denotes the similarity score, $\text{sim}(\cdot, \cdot)$ denotes cosine similarity, and τ is a temperature hyperparameter. This unified objective fosters a shared

embedding space capable of generalizing across heterogeneous tabular understanding tasks.

5 Experiments

5.1 Implementation Details

We initialize TabEmbed using the Qwen3-Embedding family (Zhang et al., 2025b) across three scales: 0.6B, 4B, and 8B parameters. This selection allows us to evaluate the scalability of our unified training paradigm across varying computational regimes. The models are optimized using a contrastive learning objective within the Sentence-Transformers framework. We conduct evaluations on our proposed TabBench, with dataset statistics detailed in Figure 2. To construct the training data, we curate a balanced mixture of 500,000 retrieval and 100,000 classification contrastive triplets from the T4 dataset. For evaluation metrics, we report Accuracy and F1-Score for the tabular prediction task, and MRR@10 and nDCG@10 for the tabular retrieval task. Further implementation details and evaluation protocols are provided in Appendix A.

5.2 Main Results

We evaluate TabEmbed on TabBench against a comprehensive suite of state-of-the-art generalist text embedding models across three parameter scales: the **0.6B scale**, including Jina-Embeddings-v3 (Sturua et al., 2024), Jasper-Token-Compression (Zhang et al., 2025a), and Qwen3-Embedding-0.6B (Zhang et al., 2025b); the **4B**

scale, including F2LLM-4B (Zhang et al., 2025c) and Qwen3-Embedding-4B; and the **7B-8B scale**, featuring SFR-Embedding-Mistral (Rui Meng, 2024), Linq-Embed-Mistral (Junseong Kim, 2024), GTE-Qwen2-7B-Instruct (Li et al., 2023), and Qwen3-Embedding-8B. Detailed specifications for all baseline models are provided in Appendix H.

Table 1 presents the performance evaluation. The results demonstrate that TabEmbed achieves state-of-the-art performance across all parameter scales, significantly surpassing existing text embedding models. In **Tabular Retrieval**, TabEmbed yields substantial improvements, with the 0.6B model surpassing its Qwen3 backbone by over 35 points in MRR@10. This indicates that our unified contrastive learning paradigm effectively bridges the semantic gap between natural language queries and structured data, addressing a capability largely absent in text embeddings. In **Tabular Classification**, TabEmbed consistently improves accuracy and F1 scores, suggesting that the learned representations capture the fine-grained decision boundaries essential for linear separability. Crucially, our method exhibits remarkable parameter efficiency. TabEmbed-0.6B outperforms all baselines on the aggregate metric, including those in the 7B and 8B regimes. This finding suggests that domain-specific contrastive learning is more critical for tabular understanding than model scaling alone. Nevertheless, scaling TabEmbed from 0.6B to 8B yields consistent performance gains, confirming that our unified paradigm effectively leverages the capacity of larger foundation models to establish a new performance standard for tabular representation.

5.3 Performance on Diverse Backbones

To investigate the universality and robustness of our proposed training paradigm, we extend our evaluation beyond the Qwen3 family to a diverse set of backbone architectures. Specifically, we apply the unified contrastive learning paradigm to eight distinct foundation models, spanning different architectures (e.g., Qwen3, Mistral, and XLM-RoBERT) and parameter scales (ranging from 0.6B to 8B). We compare the performance of these models before and after applying our training framework, utilizing the original performance as baselines.

As illustrated in Figure 4, our approach consistently yields substantial performance improvements across all evaluated backbones, regardless of their architectural design or pre-training objective. Notably, models based on the Qwen3 architecture

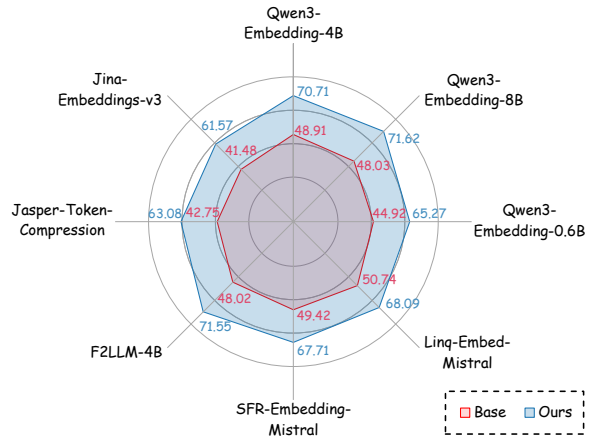


Figure 4: Performance comparison across backbone architectures using our proposed training paradigm.

(e.g., F2LLM-4B) and the Mistral architecture (e.g., Linq-Embed-Mistral) exhibit significant enhancements, with Qwen3-Embedding-4B achieving the most significant improvement, surging from 48.91 to 70.71. Even for Jina-Embeddings-v3, which relies on an encoder-only XLM-RoBERT encoder architecture, our method achieves a remarkable gain of over 20 points (rising from 41.48 to 61.57). These results demonstrate that the improvements stem from the unified contrastive data paradigm rather than model-specific inductive biases, confirming that our paradigm effectively equips diverse text-based foundation models with generalized tabular understanding capabilities.

6 Analysis and Discussion

6.1 Fine-grained Analysis on Retrieval Capabilities

While the aggregate metrics demonstrate the overall superiority of TabEmbed, it is crucial to understand how the model behaves under different semantic modalities and logical complexities. To this end, we conduct a fine-grained breakdown of the retrieval performance on the Qwen3-Embedding-0.6B backbone, categorizing the test queries by type (Numeric, Categorical, and Mixed) and the number of logical constraints (from 1 to 3).

As illustrated in Figure 5, TabEmbed achieves consistent and substantial improvements across all query scenarios, yet the difficulty varies significantly by task type. The dashed lines representing the average performance reveal an inherent hierarchy of difficulty: *Categorical* queries are the most solvable (84.61), followed by *Mixed* (65.96), with *Numeric* queries presenting the greatest challenge

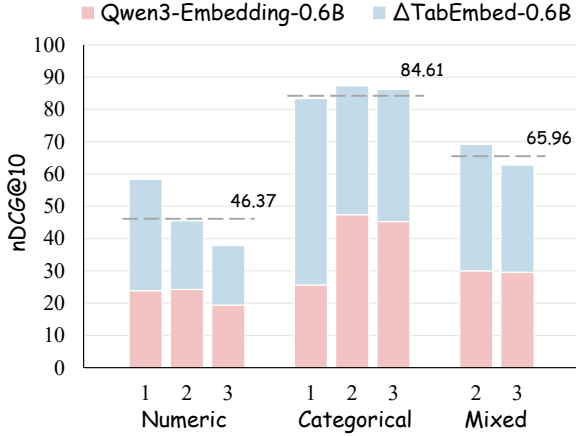


Figure 5: Fine-grained retrieval performance on TabBench (nDCG@10). The dashed lines indicate the average performance of TabEmbed for each query type.

(46.37). Crucially, the baseline model exhibits severe limitations in handling numerical queries, often failing to capture magnitude and range relationships. In contrast, TabEmbed contributes a massive performance gain in the Numeric category, effectively bridging the gap between text-based retrieval and numerical reasoning. Furthermore, regarding logical complexity, we observe that performance generally correlates with the number of constraints. For instance, in the Numeric setting, performance naturally decreases as the number of conditions increases from 1 to 3. Despite this increased difficulty, TabEmbed maintains robust performance, validating its ability to handle complex, multi-condition logical intersections within the embedding space.

6.2 Numerical Sensitivity Analysis

Standard text embedding models often treat numbers as arbitrary tokens, lacking awareness of magnitude and inequality. To investigate whether TabEmbed has acquired genuine numerical reasoning capabilities beyond surface-level token matching, we conduct a **Numerical Sensitivity Test**. Specifically, for a given query containing a numerical constraint (e.g., $q = \text{“Revenue greater than 500”}$), we generate a sequence of candidate values x ranging from small to large. We then compute the Spearman correlation (ρ) between the cosine similarity $\text{sim}(q, x)$ and the ground truth logical satisfaction (i.e., the ideal curve should step up when $x > 500$).

Figure 6 visualizes the pairwise comparison of these correlation coefficients across diverse test cases, including inequalities ($>$, $<$), equality ($=$), and range queries (*Between*). The results reveal

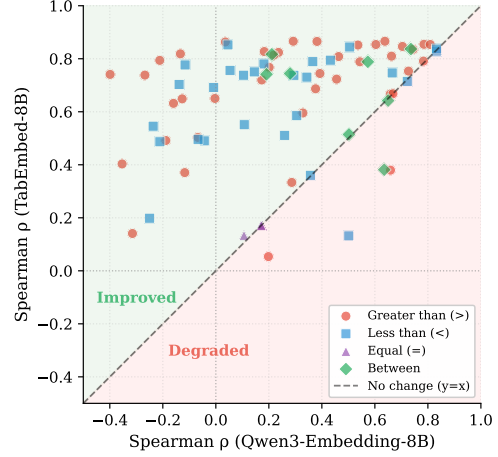


Figure 6: Pairwise comparison of numerical sensitivity between the baseline and TabEmbed. Each point represents a distinct test case, plotted by the Spearman correlation (ρ) between similarity scores and ground truth logic. Points in the green region indicate TabEmbed aligns significantly better with numerical constraints.

a distinct performance gap: the baseline Qwen3-Embedding (X-axis) frequently exhibits near-zero or weakly positive correlations, suggesting it struggles to distinguish between numerically valid and invalid candidates. In contrast, TabEmbed (Y-axis) shifts the majority of test cases into the upper-left “Improved” region, with many cases achieving high correlations ($\rho > 0.8$). This substantial shift indicates that our model has successfully internalized numerical semantics, mapping mathematically close or logically valid values to closer proximity in the vector space. Detailed visualizations of similarity curves are provided in Appendix C.

6.3 Visualization of Embedding Spaces

To provide a qualitative assessment of the learned representations, we project the high-dimensional embeddings into a 2D space using PCA and t-SNE. We visualize the geometric structures for both classification and retrieval tasks, comparing the baseline Qwen3-Embedding-8B against TabEmbed-8B. To quantify the clustering quality, we report the *Cluster Ratio*, defined as the ratio of inter-cluster distance to intra-cluster distance, where a higher ratio indicates better separability.

Figure 7(A) illustrates the feature space for classification. The baseline exhibits a highly entangled distribution (Ratio: 1.04) with significant overlap between classes, suggesting a failure to capture discriminative boundaries. In contrast, TabEmbed effectively disentangles these classes into well-separated clusters (regions A-D), substantially in-

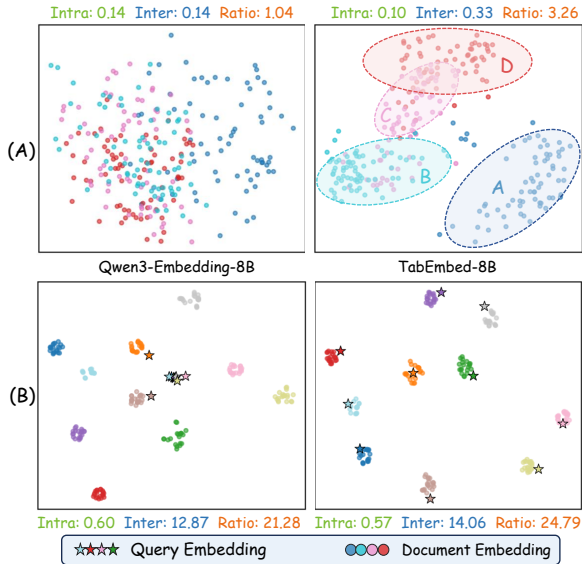


Figure 7: Visualization comparing Qwen3-Embedding-8B (left) and TabEmbed-8B (right) on tabular classification (A) and tabular retrieval (B) tasks. We report the Cluster Ratio to quantify the clustering quality.

creasing the Cluster Ratio to 3.26. This confirms that our contrastive paradigm imparts linear separability to the embedding space, enabling efficient downstream classification.

Figure 7(B) visualizes semantic alignment for retrieval. Although the baseline exhibits partial alignment capabilities, many queries (\star) remain drifting away from their target document clusters. In contrast, TabEmbed consistently anchors queries within their corresponding groups and pulls relevant documents tighter around query centers, producing significantly more compact clusters (Intra: 0.60 \rightarrow 0.57) and higher separability (Ratio: 21.28 \rightarrow 24.79). This demonstrates that TabEmbed learns a precise alignment between natural language constraints and structured tabular data, effectively correcting the misalignment observed in the baseline.

6.4 Robustness to Irrelevant Table Columns

Real-world tabular data is often characterized by high dimensionality, where a user’s query typically targets only a small subset of columns (e.g., filtering by *Price* and *City*) while ignoring numerous irrelevant attributes. Standard text embedding models are susceptible to **semantic dilution**, where irrelevant text diminishes the weight of target information in high-dimensional tables. To evaluate robustness against such structural noise, we incrementally inject up to 30 irrelevant columns into documents initially containing 15 columns and ob-

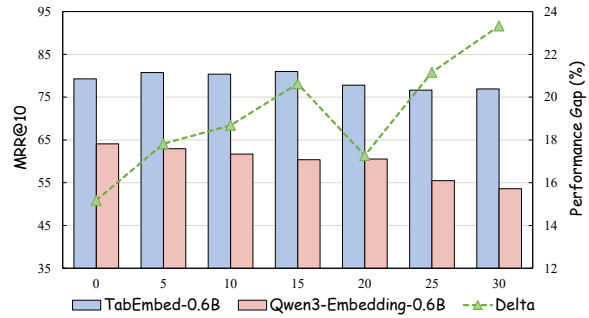


Figure 8: Robustness analysis against irrelevant table columns. We incrementally inject noise columns (0 to 30) into the documents while maintaining fixed queries.

serve the degradation in MRR@10.

As shown in Figure 8, the baseline Qwen3-Embedding exhibits a marked sensitivity to noise. Its performance declines steadily as the number of irrelevant columns increases, dropping from $\sim 64\%$ to below 55%. This confirms that without structural awareness, the model struggles to attend to the relevant signal amidst a growing volume of noise tokens. In contrast, TabEmbed demonstrates exceptional stability, consistently maintaining an MRR@10 above 75% even when 30 irrelevant columns are added. Crucially, the green dashed line highlights that the performance gap (Δ) between the two models widens monotonically from $\sim 15\%$ at the noise-free baseline to over 23% at the maximum noise level. This result suggests that TabEmbed has effectively learned an implicit structural attention mechanism, enabling it to selectively align query constraints with matching columns while filtering out unrelated tabular context.

7 Conclusion

We introduced **TabEmbed**, a unified embedding model that bridges the gap between tabular classification and retrieval. Supported by our proposed **TabBench** benchmark, we demonstrated that standard text embeddings struggle with tabular structure and numerical semantics. TabEmbed addresses these challenges through a unified contrastive learning paradigm, utilizing task-adaptive query generation and hard negative mining to learn discriminative representations. Our experiments reveal that TabEmbed achieves state-of-the-art performance, with the 0.6B model surpassing significantly larger baselines. This work establishes a baseline for generalist tabular embeddings, demonstrating that rigorous domain alignment is a more effective path to tabular intelligence than parameter scaling alone.

578 Limitations

579 Despite the promising results of TabEmbed on
580 the proposed benchmark, there are several lim-
581 itations to our current study. First, due to the
582 substantial scale of TabBench (comprising over
583 300 datasets) and budgetary constraints, we did
584 not include commercial closed-source embedding
585 APIs (e.g., Google Gemini Embedding (Lee et al.,
586 2025)) in our evaluation. While our comparison
587 covers a wide range of state-of-the-art open-source
588 models, a comprehensive benchmarking against
589 these commercial systems remains a direction for
590 future research. Second, our method relies on
591 serializing tabular data into natural language se-
592 quences. For extremely wide tables with hundreds
593 of columns, the serialized text may exceed the max-
594 imum context window of the backbone models, po-
595 tentially leading to information truncation. Devel-
596 oping more token-efficient serialization strategies
597 or employing long-context architectures to handle
598 ultra-wide tables is an avenue we plan to explore
599 in future work.

600 References

601 Bernd Bischl, Giuseppe Casalicchio, Matthias Feurer,
602 Pieter Gijsbers, Frank Hutter, Michel Lang, Rafael G
603 Mantovani, Jan N van Rijn, and Joaquin Vanschoren.
604 2017. Openml benchmarking suites. *arXiv preprint*
605 *arXiv:1708.03731*.

606 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and
607 Kristina Toutanova. 2019. Bert: Pre-training of deep
608 bidirectional transformers for language understand-
609 ing. In *Proceedings of the 2019 conference of the*
610 *North American chapter of the association for com-*
611 *putational linguistics: human language technologies,*
612 *volume 1 (long and short papers)*, pages 4171–4186.

613 Gus Eggert, Kevin Huo, Mike Biven, and Justin Waugh.
614 2023. Tablib: A dataset of 627m tables with context.
615 *arXiv preprint arXiv:2310.07875*.

616 Sebastian Felix Fischer, Matthias Feurer, and Bernd
617 Bischl. 2023. Openml-ctr23—a curated tabular re-
618 gression benchmarking suite. In *AutoML Conference*
619 *2023 (Workshop)*.

620 Tianyu Gao, Xingcheng Yao, and Danqi Chen. 2021.
621 Simcse: Simple contrastive learning of sentence em-
622 beddings. *arXiv preprint arXiv:2104.08821*.

623 Josh Gardner, Juan C Perdomo, and Ludwig Schmidt.
624 2024. Large scale transfer learning for tabular data
625 via language modeling. *Advances in Neural Informa-*
626 *tion Processing Systems*, 37:45155–45205.

627 Léo Grinsztajn, Edouard Oyallon, and Gaël Varoquaux.
628 2022. Why do tree-based models still outperform

deep learning on typical tabular data? *Advances in*
neural information processing systems, 35:507–520. 629 630

Sylvain Gugger, Lysandre Debut, Thomas Wolf, Philipp
Schmid, Zachary Mueller, Sourab Mangrulkar, Marc
Sun, and Benjamin Bossan. 2022. Accelerate: Train-
ing and inference at scale made simple, efficient and
adaptable. [https://github.com/huggingface/](https://github.com/huggingface/accelerate)
[accelerate](https://github.com/huggingface/accelerate). 631 632 633 634 635 636

Stefan Hegselmann, Alejandro Buendia, Hunter Lang,
Monica Agrawal, Xiaoyi Jiang, and David Sontag.
2023. Tabllm: Few-shot classification of tabular
data with large language models. In *International*
conference on artificial intelligence and statistics,
pages 5549–5581. PMLR. 637 638 639 640 641 642

Jihoon Kwon Sangmo Gu Yejin Kim Minkyung Cho
Jy-yong Sohn Chanyeol Choi Junseong Kim, Seol-
hwa Lee. 2024. Linq-embed-mistral: elevating text
retrieval with improved gpt data through task-specific
control and quality refinement. Linq AI Research
Blog. 643 644 645 646 647 648

Chankyu Lee, Rajarshi Roy, Mengyao Xu, Jonathan
Raiman, Mohammad Shoeybi, Bryan Catanzaro, and
Wei Ping. 2024. Nv-embed: Improved techniques for
training llms as generalist embedding models. *arXiv*
preprint arXiv:2405.17428. 649 650 651 652 653

Jinhyuk Lee, Feiyang Chen, Sahil Dua, Daniel
Cer, Madhuri Shanbhogue, Iftekhar Naim, Gus-
tavo Hernández Ábrego, Zhe Li, Kaifeng Chen, Hen-
rique Schechter Vera, and 1 others. 2025. Gemini
embedding: Generalizable embeddings from gemini.
arXiv preprint arXiv:2503.07891. 654 655 656 657 658 659

Zehan Li, Xin Zhang, Yanzhao Zhang, Dingkun Long,
Pengjun Xie, and Meishan Zhang. 2023. Towards
general text embeddings with multi-stage contrastive
learning. *arXiv preprint arXiv:2308.03281*. 660 661 662 663

Duncan McElfresh, Sujay Khandagale, Jonathan
Valverde, Vishak Prasad C, Ganesh Ramakrishnan,
Micah Goldblum, and Colin White. 2023. When do
neural nets outperform boosted trees on tabular data?
Advances in Neural Information Processing Systems,
36:76336–76369. 664 665 666 667 668 669

Andreas C Mueller, Carlo A Curino, and Raghu Ramakr-
ishnan. 2025. **Mothernet: Fast training and inference**
via hyper-network transformers. In *The Thirteenth*
International Conference on Learning Representa-
tions. 670 671 672 673 674

Niklas Muennighoff. 2022. Sgpt: Gpt sentence
embeddings for semantic search. *arXiv preprint*
arXiv:2202.08904. 675 676 677

Niklas Muennighoff, Nouamane Tazi, Loïc Magne, and
Nils Reimers. 2023. Mteb: Massive text embedding
benchmark. In *Proceedings of the 17th Conference*
of the European Chapter of the Association for Com-
putational Linguistics, pages 2014–2037. 678 679 680 681 682

Table 2: Architectural specifications of the TabEmbed model family. All variants are initialized from the corresponding Qwen3-Embedding checkpoints and inherit their structural configurations.

Model	Layers	Hidden Dim	Max Context
TabEmbed-0.6B	28	1024	32K
TabEmbed-4B	36	2560	32K
TabEmbed-8B	36	4096	32K

long-range tabular dependencies, we set the maximum sequence length to 1024 tokens. The global batch size is set to 256. For the dataset composition, we sample 500,000 retrieval triplets and 100,000 classification triplets from the processed T4 corpus. To ensure training stability, particularly for the 8B parameter models, we employ BFloat16 (BF16) mixed-precision training.

A.2 Model Architecture

TabEmbed is built upon the dense decoder-only architecture of the Qwen3-Embedding family. We release TabEmbed in three sizes (0.6B, 4B, and 8B) to cater to diverse computational constraints. While our fine-tuning protocol utilizes a context length of 1,024 tokens to optimize training throughput, the underlying architecture supports distinctively long contexts (up to 32K tokens) and variable embedding dimensions. The detailed architectural specifications for each model variant are summarized in Table 2.

A.3 Evaluation Protocols

To ensure rigorous and reproducible evaluation, we fix the random seed to 42 across all experiments. Table 3 summarizes the scale and composition of our evaluation benchmarks (TabBench). The specific evaluation protocols for the two tasks are as follows:

Tabular Classification (Linear Probing) We assess the linear separability of the embedding space by training a lightweight classifier on top of fixed representations. Specifically, we freeze the parameters of the embedding model and encode all training and testing samples into dense vectors. We then train a Logistic Regression classifier using the `scikit-learn` library. The classifier is configured with a maximum of 1,000 iterations (`max_iter=1000`) to ensure convergence. We report **Accuracy** and **Macro-F1 Score** to account for potential class imbalances in the source datasets.

Table 3: Statistics of the Tabular Embedding Benchmark (TabBench). The benchmark aggregates datasets from four diverse high-quality sources for classification and constructs a large-scale corpus for retrieval tasks.

Category	Count	# Samples / Corpus
CLASSIFICATION BENCHMARKS		
Grinsztajn	56	521,889
OpenML-CC18	66	249,939
OpenML-CTR23	34	210,026
UniPredict	155	386,618
<i>Classification Total</i>	<i>311</i>	<i>1,368,472</i>
RETRIEVAL BENCHMARKS		
Corpus	—	1,394,247
Numeric Queries	10,000	—
Categorical Queries	10,000	—
Mixed Queries	10,000	—
<i>Retrieval Total</i>	<i>30,000</i>	<i>1,394,247</i>

Tabular Retrieval (Dense Retrieval) For the retrieval task, we utilize the `Faiss` library for efficient vector similarity search. We employ an exact search strategy using the `IndexFlatIP` index (Inner Product), which corresponds to Cosine Similarity as all embeddings are L_2 -normalized prior to indexing. For each query, we retrieve the top- k most similar documents from the corpus. Performance is measured using **MRR@10** (Mean Reciprocal Rank) and **nDCG@10** (Normalized Discounted Cumulative Gain), which evaluate the ranking quality of the relevant ground-truth rows. While our main results report metrics at $k = 10$, we also compute Recall and Precision at various cutoffs ($k \in \{1, 5, 10, 20, 50, 100\}$) for comprehensive analysis.

A.4 Hardware and Infrastructure

All experiments are conducted on a high-performance computing cluster equipped with 16 PPU-810E accelerators, each possessing 96GB of high-bandwidth memory. To efficiently fine-tune the large-scale models (up to 8 billion parameters), we implement a composite optimization strategy. This includes **DeepSpeed ZeRO Stage 2** for optimizer state partitioning and gradient checkpointing to reduce memory fragmentation. The multi-GPU training is orchestrated via distributed data parallelism (DDP), ensuring linear scaling of the effective batch size.

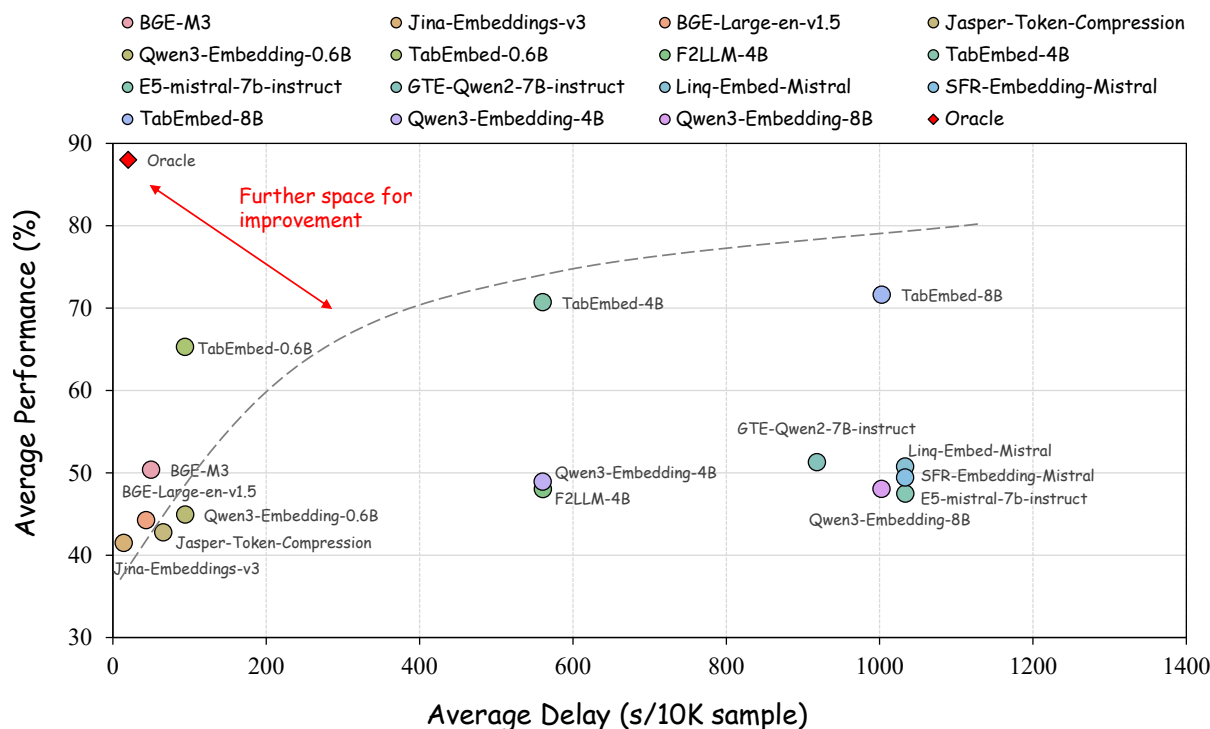


Figure 9: Performance vs. Latency trade-off. The Y-axis represents the overall average performance on TabBench, while the X-axis denotes the average inference delay (seconds per 10,000 samples).

B Inference Efficiency Analysis

To assess the practical viability of TabEmbed for real-world deployment, particularly in resource-constrained environments, we conduct a comprehensive analysis of the trade-off between model performance and inference latency. Figure 9 illustrates the relationship between the aggregate performance on TabBench (Y-axis) and the computational cost (X-axis) across different model scales.

We measured inference latency using a standardized benchmarking protocol on a single PPU-810E accelerator. To simulate realistic input distributions comparable to those found in TabBench, we constructed a synthetic dataset comprising serialized tabular rows with lengths varying uniformly between 50 and 200 words. All models were evaluated under identical conditions: a batch size of 64 and a maximum sequence length of 1024 tokens. To ensure statistical stability, we performed a warm-up phase followed by three independent experimental runs, reporting the average latency normalized per 10,000 samples.

The results reveal distinct performance-efficiency clusters corresponding to parameter scales. In the low-latency regime, standard text embedding models such as Jina-Embeddings-v3 and Qwen3-Embedding-0.6B offer high throughput but

demonstrate limited capability in capturing tabular semantics, with performance scores hovering around 45%. **TabEmbed-0.6B** significantly disrupts this trend, achieving a performance score of 65.27% while maintaining a highly efficient latency profile (≈ 94 seconds per 10k samples). This indicates that domain-specific contrastive learning can unlock tabular reasoning capabilities in lightweight architectures without incurring additional inference costs.

In the high-capacity regime (4B and 8B parameters), TabEmbed continues to push the performance boundary, reaching up to 71.62% with the 8B variant. However, this performance gain comes with a considerable increase in computational cost, with latency exceeding 1,000 seconds per 10k samples. TabEmbed-4B offers a compelling middle ground, delivering near-peak performance at approximately half the inference cost of the 8B model. The plot also includes a theoretical "Oracle" point, highlighting the gap that remains between current state-of-the-art models and an ideal system with minimal delay and maximum accuracy. This suggests that future research directions should focus on knowledge distillation or quantization techniques to retain the structural reasoning capabilities of TabEmbed-8B within the latency budget of smaller models.

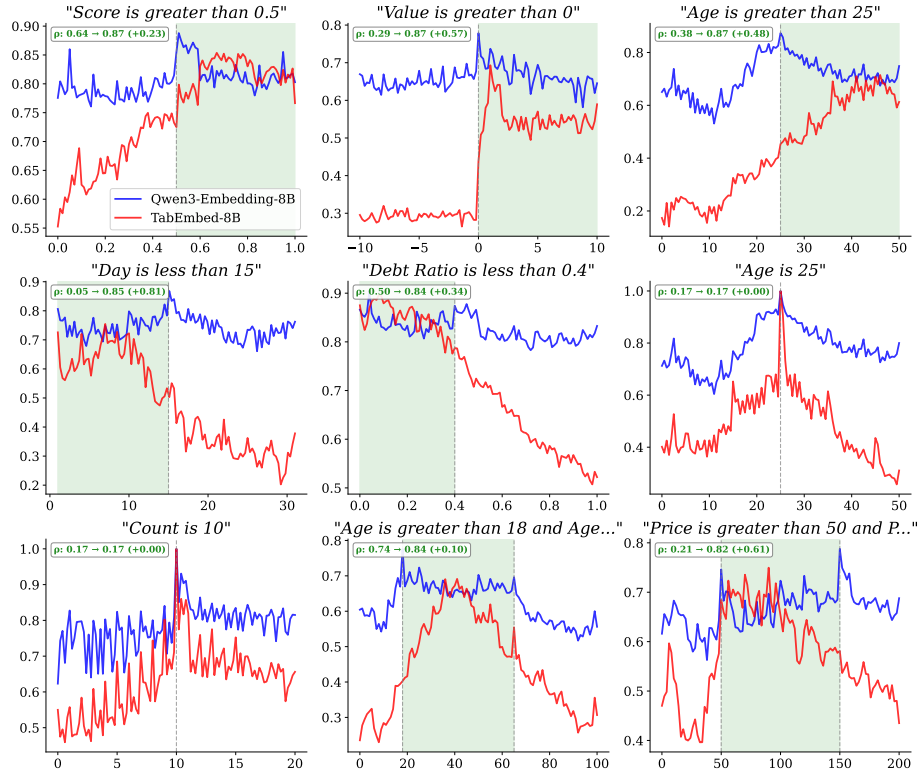


Figure 10: Similarity curves for 9 representative numerical reasoning tasks. **X-axis:** candidate value; **Y-axis:** cosine similarity with the query. **Green shaded regions** indicate valid ranges where conditions are satisfied. **Blue lines:** baseline Qwen3-Embedding-8B; **Red lines:** TabEmbed-8B. Spearman correlation (ρ) improvements are annotated in each subplot.

C Numeric Sensitivity Curves

To provide a granular view of the numerical reasoning capabilities discussed in Section 6.2, Figure 10 displays the detailed cosine similarity trajectories for nine representative test cases. For each test case, we define a query q containing a specific numerical constraint (e.g., “Age is greater than 25”) and generate a sequence of 101 candidate documents $d(x)$ with values linearly spaced across a relevant range. We then compute the cosine similarity score regarding the logical truth value: ideal embeddings should yield high similarity only when the condition is met (indicated by the green shaded regions).

As illustrated by the blue lines in Figure 10, the baseline Qwen3-Embedding-8B typically exhibits random fluctuations or weak correlations. For instance, in “Score is greater than 0.5”, the baseline’s similarity scores remain relatively flat or erratic regardless of x , confirming that standard text embeddings treat numbers primarily as independent tokens without inherent ordinal semantics.

In contrast, TabEmbed-8B (red lines) demonstrates distinct, logic-aware behaviors tailored to the specific operator types:

- Inequalities ($>$, $<$):** The model approximates a step function with sharp transitions at the decision boundary. For example, in “Age is greater than 25”, the similarity rises abruptly as x approaches the threshold and sustains a high plateau within the valid range, whereas for “Day is less than 15”, it drops significantly once the threshold is exceeded.
- Equalities ($=$):** For exact matching tasks like “Age is 25” or “Count is 10”, TabEmbed produces a sharp peak centered exactly at the target value, mimicking a Dirac delta function to distinguish the target from numerically adjacent distractors.
- Composite Ranges (Between):** For queries involving logical conjunctions (e.g., “Age is greater than 18...”), the model accurately delineates the intersection interval, maintaining high similarity only where both conditions hold true.

The substantial improvements in Spearman correlation (ρ) annotated in each subplot (e.g., $0.64 \rightarrow 0.87$) quantitatively verify that TabEmbed has successfully aligned its embedding space with the underlying mathematical logic.

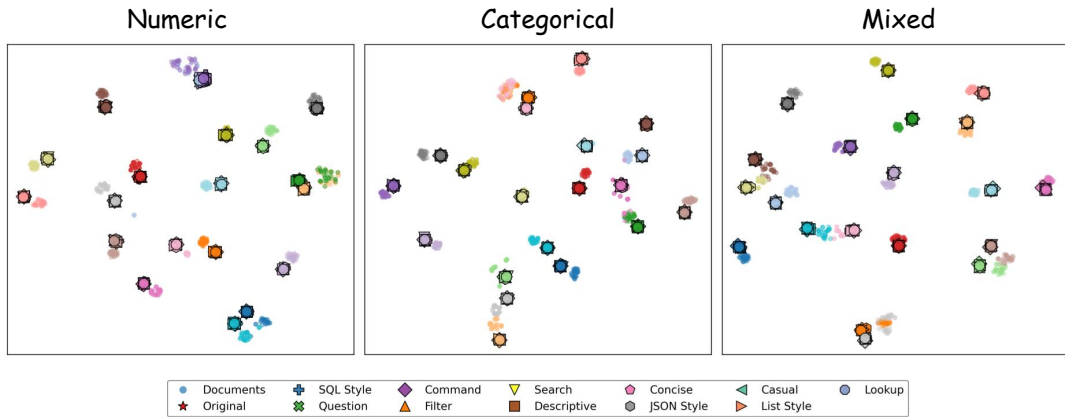


Figure 11: t-SNE visualization of query template robustness across Numeric, Categorical, and Mixed tasks. For each semantic intent, we generate 12 distinct query variations (e.g., SQL-style, JSON-style, Casual) as defined in Table 4. The visualization demonstrates that despite significant syntactic differences, all query variations (represented by distinct markers) cluster tightly around the same relevant documents (colored circles), indicating that TabEmbed learns a syntax-agnostic representation of tabular constraints.

D Robustness to Query Template Variations

A critical requirement for a generalist tabular embedding model is the ability to understand the underlying user intent regardless of the input format. Users may express the same retrieval constraint through diverse modalities, ranging from formal syntaxes (e.g., SQL, JSON) to unstructured natural language (e.g., questions, commands). To evaluate whether TabEmbed has learned a **syntax-agnostic representation** of tabular constraints, we conducted a qualitative visualization experiment. We selected representative samples from the Numeric, Categorical, and Mixed retrieval tasks. For each sample, we generated 12 distinct variations using the templates listed in Table 4, effectively creating a set of semantic equivalence classes where queries differ in surface form but share identical logical constraints.

Figure 11 presents the t-SNE projection of these embeddings. The visualization reveals a striking geometric pattern: for every semantic intent, the diverse query variations (represented by distinct markers such as stars, crosses, and triangles) form tight, cohesive clusters surrounding their corresponding ground-truth documents (colored circles). Notably, this alignment persists across extreme syntactic disparities. For instance, highly structured formats like **JSON Style** (T9: "age": 25...) and **SQL Style** (T2: SELECT * FROM...) are mapped to the immediate vicinity of unstructured natural language queries like **Casual** (T10) and **Question** (T3). This observation confirms that TabEmbed

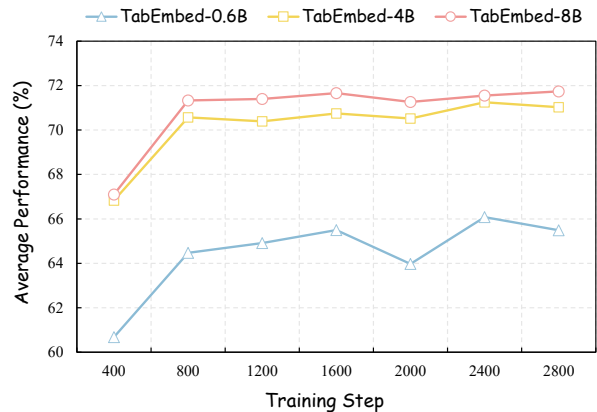


Figure 12: The impact of training steps on the average performance across TabBench. We track the evaluation metrics of TabEmbed at different parameter scales (0.6B, 4B, and 8B) from 400 to 2800 training steps.

does not merely rely on keyword matching or rigid template overfitting. Instead, it has successfully learned to extract the **invariant logical semantics** (e.g., numerical magnitude and equality constraints) from the input, projecting semantically equivalent queries to the same point on the manifold regardless of their linguistic style. This capability ensures that TabEmbed can generalize to diverse real-world search scenarios where user querying habits may vary significantly.

E Training Convergence Analysis

To determine the optimal training duration and investigate the convergence behavior of our unified contrastive learning paradigm, we monitor the average performance of TabEmbed across three param-

Table 4: Query template variations used to evaluate semantic robustness of embedding models. All templates express the same underlying constraints but differ in linguistic style and format.

ID	Template Name	Example Query
T1	Original	Find records where age is 25 and salary greater than 50000
T2	SQL Style	SELECT * FROM table WHERE age = 25 AND salary > 50000
T3	Question	Which records have age equal to 25 and salary greater than 50000?
T4	Command	Get all entries with age of 25 and salary above 50000
T5	Filter	Filter: age==25, salary>50000
T6	Search	Search for records: age:25 salary:>50000
T7	Descriptive	I need data where the age is 25 and the salary is more than 50000
T8	Concise	age == 25 salary > 50000
T9	JSON Style	{"age": 25, "salary": {"\$gt": 50000}}
T10	Casual	Show me rows that have age as 25 and salary over 50000
T11	List Style	Conditions: 1. age equals 25; 2. salary > 50000
T12	Lookup	Look up records matching: age=25, salary>50000

eter scales (0.6B, 4B, and 8B) at regular intervals, ranging from 400 to 2800 training steps.

As illustrated in Figure 12, we observe several key findings regarding training efficiency and model scaling. First, Rapid Convergence. All models exhibit a steep performance trajectory in the initial phase. Notably, a significant portion of the performance gain is realized within the first 800 steps. For instance, the 4B model improves from roughly 67% to over 70% in this short period. The performance curves generally plateau after approximately 1600 steps, suggesting that our framework is highly data-efficient and does not require excessively prolonged training to learn robust tabular representations. Second, Impact of Model Scale. Consistent with neural scaling laws, larger models consistently achieve higher performance ceilings. TabEmbed-8B (pink line) maintains a clear superiority over the 4B and 0.6B variants throughout the training process. Furthermore, model scale correlates positively with training stability. The 4B and 8B models display smooth and monotonic improvements, whereas the 0.6B model (blue line) exhibits noticeable volatility, particularly around step 2000. This suggests that larger foundation models possess a more robust latent space, making them less susceptible to batch noise during contrastive optimization.

F Target Column Selection Criteria

To transform the unannotated tables from the T4 corpus into high-quality supervised classification

tasks, we employ a dynamic target identification pipeline. For a given table \mathbf{T} , we first consider all columns as potential target candidates and then apply a rigorous filtering protocol to exclude non-informative or trivial prediction targets. Specifically, a column c is excluded from the candidate pool if it satisfies any of the following rejection criteria:

- **Low Informativeness:** The column contains only a single unique value (constant columns), providing no discriminative signal.
- **High Cardinality / Identifiers:** The column possesses a unique value for every row (e.g., UUIDs, Row IDs), or the number of unique classes exceeds 50. Such columns typically lead to trivial memorization rather than semantic generalization.
- **Data Type Constraints:** The column is identified as a date/timestamp, or contains textual values exceeding 256 characters, which are unsuitable for standard classification objectives.
- **Column name constraint:** The column name contains "Unnamed:" (pandas' default marker for unnamed columns).

It is important to note that columns failing these criteria are only excluded from being selected as the *prediction target* y . They remain part of the input feature set \mathbf{x}_{-y} to provide context, unless

1066	they are removed by standard feature selection processes.	1115
1067		1116
1068	Once the set of valid candidate columns is established, we select a single target y for each training instance. To balance the distribution of task types, we employ a weighted sampling strategy. Based on our qualitative observation that categorical columns often yield higher-quality decision boundaries than arbitrary numerical regression targets, we prioritize classification tasks. Specifically, if both continuous and categorical candidates are present, we sample a categorical target with probability $p = 0.5$ and a continuous target with probability $1 - p = 0.5$.	1117
1069		
1070		
1071		
1072		
1073		
1074		
1075		
1076		
1077		
1078		
1079	In cases where a continuous column is selected as the target, we transform the regression problem into a classification problem via dynamic discretization. We divide the continuous range into quantile-based bins (defaulting to 4 buckets) to ensure class balance. The resulting targets are serialized into natural language class descriptors, such as “less than 15.5”, “between 15.5 and 40.2”, or “greater than 40.2”. This unified serialization allows TabEmbed to handle both original categorical labels and discretized numerical bins within the same semantic embedding space.	
1080		
1081		
1082		
1083		
1084		
1085		
1086		
1087		
1088		
1089		
1090		
1091	G Applications	
1092	The unified representation capability of TabEmbed and the comprehensive evaluation framework of TabBench open up several promising avenues for real-world applications, particularly in scenarios where traditional schema-bound methods fall short.	
1093		
1094		
1095		
1096		
1097	G.1 Semantic Interface for NoSQL and Vector Databases	
1098		
1099	Traditional approaches to querying structured data typically rely on Text-to-SQL parsing, which is notoriously brittle and depends heavily on rigid schemas. This limitation is particularly pronounced in NoSQL databases (e.g., MongoDB) and modern Vector Databases (e.g., Milvus, Pinecone), which often lack full SQL support or store semi-structured data (JSON) where writing explicit queries is complex. TabEmbed offers a paradigm shift by enabling a “ Semantic SQL ” interface. Since TabEmbed maps natural language constraints (e.g., “ <i>High-value users from the tech sector</i> ”) and structured rows into a shared vector space, it allows these databases to perform row-level retrieval directly via vector similarity search. This effectively bypasses the need for intermediate logical	
1100		
1101		
1102		
1103		
1104		
1105		
1106		
1107		
1108		
1109		
1110		
1111		
1112		
1113		
1114		
	forms, enabling robust Retrieval-Augmented Generation (RAG) pipelines over tabular data without the fragility of code generation.	1118
	G.2 Enterprise Data Discovery and Data Lakes	1119
	Large enterprises often maintain massive data lakes containing thousands of unorganized spreadsheets and CSV files (the “dark data” problem). Finding relevant datasets for a specific analytical task is a major challenge, as file names are often non-descriptive. TabEmbed can serve as a core component for Semantic Data Discovery . By embedding sample rows or summarized schemas from thousands of tables into a unified index, users can search for datasets using vague intent queries (e.g., “ <i>I need sales data regarding Q3 revenue in Southeast Asia</i> ”). Unlike keyword-based search, TabEmbed understands the numerical and categorical semantics within the table content, retrieving relevant tables even if the column headers do not explicitly match the query keywords.	1120
		1121
		1122
		1123
		1124
		1125
		1126
		1127
		1128
		1129
		1130
		1131
		1132
		1133
		1134
		1135
	G.3 Zero-Shot and Cold-Start Tabular Prediction	1136
		1137
	In many dynamic industrial applications (e.g., fraud detection in new markets, user churn prediction for new products), historical training data is scarce or unavailable, rendering traditional supervised models (like XGBoost) ineffective. As demonstrated by our classification experiments, TabEmbed possesses strong zero-shot transfer capabilities. It can function as a generic feature extractor or a nearest-neighbor classifier for Cold-Start Scenarios . Practitioners can simply convert a handful of labeled examples (support set) into vectors and classify new incoming rows based on embedding similarity, enabling immediate predictive capabilities without the need for time-consuming feature engineering or model training.	1138
		1139
		1140
		1141
		1142
		1143
		1144
		1145
		1146
		1147
		1148
		1149
		1150
		1151
		1152
	H Baselines	1153
	We compare TabEmbed against a diverse set of state-of-the-art generalist text embedding models, ranging from lightweight encoders to large-scale LLM-based embeddings.	1154
		1155
		1156
		1157
	0.6B Parameter Scale	1158
	• Jina-Embeddings-v3 (Sturua et al., 2024): A multilingual, multi-task embedding model based on the Jina-XLM-RoBERTa architecture. It incorporates Rotary Position Embeddings (RoPE)	1159
		1160
		1161
		1162

1163 to support extended context windows up to 8192
1164 tokens. A key feature of this model is the in-
1165 tegration of five task-specific LoRA adapters,
1166 allowing for efficient generation of embeddings
1167 tailored to specific downstream applications.

- 1168 • **Jasper-Token-Compression** (Zhang et al.,
1169 2025a): A 600M parameter model from the
1170 Jasper series that introduces dynamic text token
1171 compression, inspired by DeepSeek-OCR strate-
1172 gies. By combining vector distillation with con-
1173 trastive learning, it achieves high performance
1174 while compressing textual information by ap-
1175 proximately 10x, offering a unique approach to
1176 efficient representation.
- 1177 • **Qwen3-Embedding-0.6B** (Zhang et al., 2025b):
1178 The lightweight variant of the latest proprietary
1179 embedding series from the Qwen family. Build-
1180 ing upon the dense Qwen3 foundational architec-
1181 ture, it inherits strong multilingual capabilities
1182 and reasoning skills, optimized specifically for
1183 retrieval and ranking tasks.

1184 4B Parameter Scale

- 1185 • **F2LLM-4B** (Zhang et al., 2025c): Standing for
1186 “Foundation to Feature Large Language Mod-
1187 els,” F2LLM is fine-tuned on a curated corpus
1188 of 6 million high-quality query-document pairs
1189 sourced exclusively from open-source datasets.
1190 It employs a single-stage training process with
1191 homogeneous macro batches, eschewing com-
1192 plex multi-stage pipelines while covering diverse
1193 retrieval and clustering tasks.
- 1194 • **Qwen3-Embedding-4B** (Zhang et al., 2025b):
1195 A mid-sized model in the Qwen3 embedding se-
1196 ries. It balances computational efficiency with
1197 the advanced long-text understanding capabili-
1198 ties of the Qwen3 foundation, serving as a strong
1199 baseline for mid-scale generalist text embed-
1200 dings.

1201 7B-8B Parameter Scale

- 1202 • **SFR-Embedding-Mistral** (Rui Meng, 2024):
1203 Developed by Salesforce Research, this model
1204 is initialized from E5-Mistral-7b-instruct and
1205 Mistral-7B-v0.1. It represents a robust baseline
1206 for instruction-tuned embeddings derived from
1207 decoder-only architectures.
- 1208 • **Linq-Embed-Mistral** (Junseong Kim, 2024):
1209 Also built upon the E5-Mistral and Mistral-7B

foundations, Linq-Embed-Mistral focuses on en- 1210
hancing retrieval performance through advanced 1211
data refinement. Its training pipeline empha- 1212
sizes sophisticated data crafting, rigorous filter- 1213
ing, and hard-negative mining guided by teacher 1214
models to improve the quality of synthetic train- 1215
ing triplets. 1216

- 1217 • **GTE-Qwen2-7B-Instruct** (Li et al., 2023): The 1218
latest addition to the General Text Embedding 1219
(GTE) family, built on the Qwen2-7B LLM. It 1220
leverages the same training data and strategies 1221
as its predecessor (GTE-Qwen1.5) but benefits 1222
from the architectural upgrades of the Qwen2 1223
base model. It is a leading performer on the 1224
MTEB benchmark, particularly in multilingual 1225
evaluation scenarios.
- 1226 • **Qwen3-Embedding-8B** (Zhang et al., 2025b): 1227
The largest model in our comparison suite and 1228
the direct backbone for TabEmbed-8B. It repre- 1229
sents the state-of-the-art in the Qwen family for 1230
dense retrieval, bitext mining, and classification, 1231
providing a rigorous baseline to measure the im- 1232
pact of our domain-specific contrastive learning 1233
paradigm.