

FinMTEB: Finance Massive Text Embedding Benchmark

Anonymous ACL submission

Abstract

Embedding models play a crucial role in representing and retrieving information across various NLP applications. Recent advances in large language models (LLMs) have further enhanced the performance of embedding models. While these models are often benchmarked on general-purpose datasets, real-world applications demand domain-specific evaluation. In this work, we introduce the **Finance Massive Text Embedding Benchmark** (FinMTEB), a specialized counterpart to MTEB designed for the financial domain. FinMTEB comprises 64 financial domain-specific embedding datasets across 7 tasks that cover diverse textual types in both Chinese and English, such as financial news articles, corporate annual reports, ESG reports, regulatory filings, and earnings call transcripts. We also develop a finance-adapted model, Fin-E5, using a persona-based data synthetic method to cover diverse financial embedding tasks for training. Through extensive evaluation of 15 embedding models, including Fin-E5, we show three key findings: (1) performance on general-purpose benchmarks shows limited correlation with financial domain tasks; (2) domain-adapted models consistently outperform their general-purpose counterparts; and (3) surprisingly, a simple Bag-of-Words (BoW) approach outperforms sophisticated dense embeddings in financial Semantic Textual Similarity (STS) tasks, underscoring current limitations in dense embedding techniques. Our work establishes a robust evaluation framework for financial NLP applications and provides crucial insights for developing domain-specific embedding models.

1 Introduction

Embedding models, which transform text sequences into dense vector representations, serve as fundamental building blocks in natural language processing (NLP) tasks (Mikolov et al., 2013; Pennington et al., 2014; Peters et al., 2018). The quality



Figure 1: Word cloud visualization of Fin-E5’s training data, contain common financial terms.

of text embeddings directly impacts the effectiveness of information retrieval, semantic understanding, and other downstream applications. Although recent large language model (LLM)-based embedding models have shown remarkable performance on general benchmarks (Wang et al., 2023; Li et al., 2023; Meng et al., 2024), their effectiveness in specialized domains, particularly finance, remains understudied. Financial text analysis requires precise handling of domain-specific terminology, temporal sensitivity, and complex numerical relationships (Li et al., 2024; Anderson et al., 2024). This raises two critical questions:

- How effectively do modern embedding models capture domain-specific financial information?
- Can domain adaptation improve LLM-based embeddings for financial applications?

These questions are motivated by three key insights. First, financial semantics often diverge from general language usage. For example, the term "*liability*" inherently conveys negative sentiment in financial contexts due to its association with obligations and risks, whereas in general usage, it neutrally denotes legal responsibility. Such semantic divergence becomes particularly crucial for real-world applications such as Retrieval Augmented

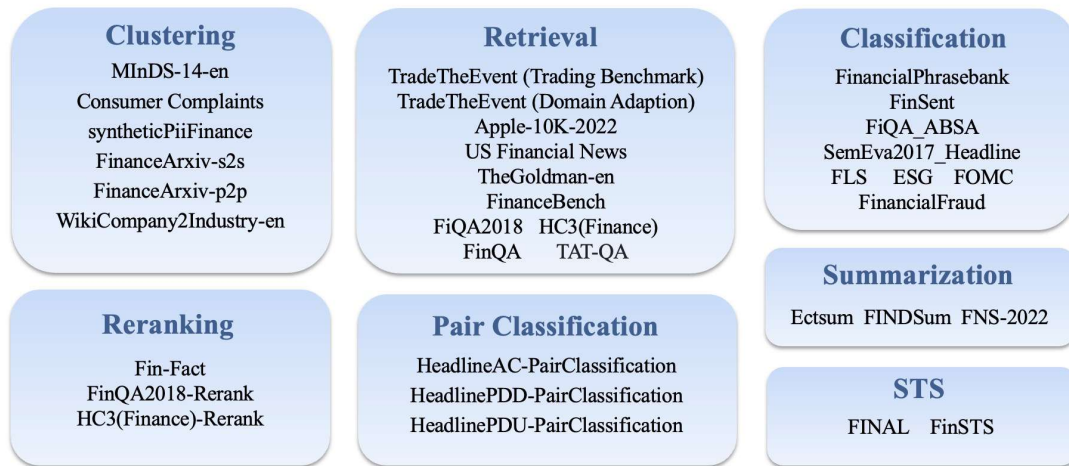


Figure 2: An overview of tasks and datasets used in FinMTEB. All the dataset descriptions and examples are provided in the Appendix A.

Generation (RAG) systems, where accurate document retrieval underpins effective knowledge enhancement. While recent work adapts RAG frameworks for finance (Li et al., 2024; Malandri et al., 2025), the fundamental role of embedding quality in retrieval performance remains overlooked.

Second, empirical evidence increasingly suggests that domain adaptation is crucial for achieving optimal performance in specialized fields (Ling et al., 2023; Gururangan et al., 2020), even with recent advanced LLMs. This necessity for domain specialization has led to the development of field-specific models across various domains: BiMedLM (Bolton et al., 2024) for biomedical texts, SaulLM-7B (Colombo et al., 2024) for legal documents, and BloombergGPT (Wu et al., 2023) for financial applications. This specialization trend extends to embedding models, where domain-specific variants have demonstrated superior performance in capturing specialized vocabulary and semantic relationships. For instance, BioWordVec (Zhang et al., 2019) and BioSentVec (Chen et al., 2019) are used in biomedical text analysis, while FinBERT (Yang et al., 2020) shows promising results in financial applications. However, while financial domain embedding models have shown promising improvement (e.g., BAM (Anderson et al., 2024), a RoBERTa-based (Liu, 2019a) model outperforming the general model in retrieval tasks), they are still based on traditional architectures. Compared to the general domain, there is a gap in the current landscape for finance NLP: **while commercial solutions like voyage-finance-2 (VoyageAI, 2025) exist, there remains a lack of open-source LLM-**

based financial embedding models available to researchers.

Third, financial NLP lacks comprehensive evaluation frameworks for embedding models. Current benchmarks (Islam et al., 2023; Chen et al., 2021) primarily assess text generation rather than embedding quality. Even embedding-specific evaluations (FiQA, 2018; Liu et al., 2024a) focus narrowly on single task types (e.g., classification) or limited text genres (e.g., earnings call transcripts). This gap is deepened by financial texts’ unique characteristics, such as the prevalence of boilerplate language (e.g., *"The company’s performance is subject to various risks..."*) that creates noise in semantic representation. These standardized legal disclaimers appear frequently across documents but offer little information, complicating the models’ ability to differentiate meaningful business insights from routine compliance text. Thus, there is a critical need for comprehensive financial embedding benchmarks.

To bridge this gap, we introduce the **Finance Massive Text Embedding Benchmark (FinMTEB)**, a comprehensive evaluation framework specialized for the financial domain. FinMTEB comprises 64 domain-specific datasets that span both Chinese and English, covering seven distinct tasks: classification, clustering, retrieval, pair classification, reranking, summarization, and semantic textual similarity. We also develop Fin-E5, a finance-adapted version of e5-Mistral-7B-Instruct (Wang et al., 2023), utilizing a persona-based data synthesis method. As shown in Figure 1, our training data encompasses a diverse range of financial topics concepts. Experimental results show that LLM-based

embedding models consistently outperform traditional approaches, while domain adaptation further improves performance. Interestingly, in the STS task, we find that the simple Bag-of-Words (BoW) model outperforms all dense models. This indicates that current embedding models still encounter difficulties in interpreting complex financial texts.

Our main contributions are twofold: First, we propose FinMTEB, the first comprehensive financial domain evaluation benchmark encompassing 64 datasets across seven distinct tasks in both Chinese and English. Second, we develop and release Fin-E5, a finance-adapted embedding model that achieves state-of-the-art performance on FinMTEB. To support future research, we will make both the FinMTEB benchmark and our Fin-E5 model available as open source.

2 Related Work

Recent advances in embedding models have shown remarkable success in general domain tasks, yet their effectiveness in specialized domains remains a critical challenge.

2.1 General-purpose Embedding Models

The evolution of embedding models marks significant progress in natural language processing. Starting with static word representations like Word2Vec (Mikolov et al., 2013) and GloVe (Pennington et al., 2014), the field advanced to contextualized embeddings through transformer-based architectures such as BERT (Devlin et al., 2019) and RoBERTa (Liu, 2019b). A notable advancement came with Sentence-BERT (Reimers and Gurevych, 2019), which introduced Siamese and triplet network architectures to generate meaningful sentence-level representations. Recent developments in large language models have further pushed the boundaries, with models such as e5-mistral-7b-instruct (Wang et al., 2023) and gte-Qwen2-1.5B-instruct (Yang et al., 2024) achieving better performance in various embedding tasks. However, these general-purpose models may not adequately capture the nuanced semantics of specialized domains.

2.2 Current Embedding Evaluation Landscape

To assess embedding quality, several evaluation frameworks have been developed. General-purpose embedding benchmarks, such as the Massive Text Embedding Benchmark (MTEB) (Muennighoff

et al., 2022), provide broad coverage across multiple tasks and languages. Specialized benchmarks like BEIR (Thakur et al., 2021) focus on specific aspects, such as information retrieval. Although they incorporate some domain-specific datasets, such as FiQA (FiQA, 2018), the size of the data and the coverage of the task are limited.

2.3 Domain Adaptation Approaches

Recognizing the limitations of general-purpose models in specialized domains, researchers have pursued two main adaptation strategies. The first approach develops domain-specific models from scratch, exemplified by BioMedLM (Bolton et al., 2024) for biomedicine, SaulLM-7B (Colombo et al., 2024) for legal texts, and BloombergGPT (Wu et al., 2023) for finance. The second strategy fine-tunes existing models for domain-specific tasks, as demonstrated by InvestLM (Yang et al., 2023b) and FinGPT (Yang et al., 2023a). This trend extends to embedding models, with specialized versions such as BioWordVec (Zhang et al., 2019), BioSentVec (Chen et al., 2019), and FinBERT (Yang et al., 2020) showing superior domain-specific performance. However, evaluating these specialized embedding models remains challenging due to the lack of comprehensive domain-specific benchmarks.

2.4 The Gap in Domain-specific Evaluation

While domain-specific language models have stimulated the development of specialized evaluation frameworks across various fields, these benchmarks primarily emphasize generative and reasoning capabilities instead of embedding quality. The financial sector has seen the emergence of frameworks like CFLUE (Zhu et al., 2024), FinEval (Zhang et al., 2023), and FinanceBench (Islam et al., 2023), whereas the legal and medical domains have introduced LawBench (Fei et al., 2023), MedBench (Liu et al., 2024b), and DrBenchmark (Labrak et al., 2024). These benchmarks consistently illustrate that general-purpose models often fall short in specialized areas (Zhu et al., 2024; Fei et al., 2023), highlighting the necessity of domain adaptation (Ling et al., 2023). Despite this acknowledgment, there is still a critical lack of comprehensive evaluation frameworks for domain-specific embeddings that assess performance across essential tasks such as semantic similarity, classification, and retrieval. Even recent financial embedding developments, such as BAM embedding (Anderson

et al., 2024), rely on narrow evaluation frameworks, typically focusing on single-task performance metrics (e.g., FinanceBench (Islam et al., 2023) for retrieval tasks). This limited evaluation may not fully reflect how the models perform in real-world financial applications.

3 The FinMTEB Benchmark

In this section, we introduce the Finance MTEB (FinMTEB) benchmark. As illustrated in Figure 2, FinMTEB encompasses seven embedding tasks, following a structure similar to MTEB (Muenighoff et al., 2022) but with datasets specifically curated for the finance domain.

3.1 FinMTEB Tasks

Semantic Textual Similarity (STS) evaluates the semantic similarity between pairs of financial text. This task is crucial for automated financial analysis and risk management; for example, detecting subtle semantic differences between quarterly earnings statements could reveal important shifts in a company’s financial strategy that impact investment decisions. To ensure comprehensive evaluation, we incorporate diverse financial datasets, including FinSTS (Liu et al., 2024a) and FINAL (Ju et al., 2023) from company annual reports, and BQ-Corpus (Chen et al., 2018) from banking documents. Model performance is quantified using Spearman’s rank correlation, which measures the alignment between predicted cosine similarity scores and human-annotated similarity ratings.

Retrieval evaluates a model’s capability to identify and extract relevant financial information in response to specific queries. Unlike general domain retrieval, financial information retrieval presents unique challenges, requiring precise handling of complex numerical data, temporal dependencies, and regulatory context. For comprehensive evaluation, we leverage established finance QA datasets including FinanceBench (Islam et al., 2023), FiQA2018 (FiQA, 2018), and HPC3 (Guo et al., 2023). To further assess models’ understanding of professional financial terminology, we introduce TheGoldman dataset, constructed from the Goldman Sachs Financial Dictionary. Performance is measured using NDCG@10, a metric that evaluates both the relevance of retrieved information and its ranking position, reflecting the real-world requirement for highly precise top results in financial applications.

Clustering evaluates a model’s ability to automatically group similar financial texts based on their semantic content. To ensure comprehensive evaluation, we developed multiple specialized datasets that capture different aspects of financial text clustering: (1) FinanceArxiv-s2s and FinanceArxiv-p2p, constructed from titles and abstracts of finance-related papers on arXiv, providing rich academic financial content; (2) CompanyWiki2Industry dataset, derived from Wikipedia company descriptions, offering diverse industry categorization scenarios; and (3) complementary resources including consumer complaints from CFPB¹, financial intent detection data (Gerz et al., 2021a; Watson et al., 2024), and other established datasets. Model performance is quantified using the V-measure (Rosenberg and Hirschberg, 2007), a comprehensive metric that evaluates cluster quality through both completeness (all members of a class are assigned to the same cluster) and homogeneity (each cluster contains only members of a single class).

Classification evaluates a model’s ability to categorize financial texts into predefined classes based on their semantic content. This capability is essential for automated financial decision-making; for example, in algorithmic trading, accurately classifying sentiment in earnings calls or news articles can directly influence trading strategies and portfolio adjustments. The classification task encompasses diverse financial scenarios through multiple specialized datasets, including: financial sentiment analysis (Malo et al., 2014; FiQA, 2018; Cortis et al., 2017; Lu et al., 2023), Federal Reserve monetary policy classification (Shah et al., 2023), organization’s strategy classification, and forward-looking statement identification (Yang et al., 2023b). Performance is measured using Mean Average Precision (MAP), which provides a comprehensive assessment of classification accuracy while accounting for ranking quality and confidence scores.

Reranking evaluates the model’s ability to order retrieved documents based on their relevance to financial queries. We utilize financial question-answering datasets such as Fin-Fact and FinQA (Rangapur et al., 2023; Chen et al., 2021) to construct the reranking tasks. Specifically, for each query in these datasets, we retrieve top-k relevant

¹<https://huggingface.co/datasets/CFPB/consumer-finance-complaints>

documents along with the ground truth answers to construct the reranking training and evaluation pairs. The main evaluation metric for reranking in Finance MTEB is Mean Average Precision (MAP).

Pair-Classification evaluates a model’s ability to determine semantic relationships between financial text pairs. This task includes two datasets: (1) the AFQMC dataset² for customer intention, and (2) three financial news headline datasets (Sinha and Khandait, 2021). We use Average Precision (AP) as the evaluation metric to assess model performance across different decision thresholds.

Summarization is evaluated based on the correlation between dense embeddings derived from the summarized texts and those of the original texts, utilizing Spearman’s correlation coefficient as the main metric. The evaluation corpus encompasses a comprehensive range of financial texts, including earnings call transcripts (Mukherjee et al., 2022), financial news articles (Lu et al., 2023), and SEC Form 10-K filings (El-Haj et al., 2022), ensuring robust assessment across diverse financial contexts and writing styles.

3.2 Characteristics of FinMTEB

FinMTEB contains 35 English datasets and 29 Chinese datasets. Detailed information about these datasets is provided in Appendix A.

Linguistic Pattern. Table 9 presents a comparative analysis of linguistic features between MTEB (Muennighoff et al., 2022) and FinMTEB benchmarks, examining aspects such as average sentence length, token length, syllables per token, and dependency distance (Oya, 2011). The results indicate that texts in FinMTEB consistently exhibit longer and more complex sentences than those in MTEB, with an average sentence length of 26.37 tokens compared to MTEB’s 18.2 tokens. This highlights the linguistic differences between financial and general domain texts.

Semantic Diversity. We examine the inter-dataset semantic similarity within FinMTEB. Using the all-MiniLM-L6-v2 model³, we embed 1,000 randomly sampled texts from each dataset, compute their mean embeddings to represent each dataset, and measure inter-dataset similarities using cosine similarity. As shown in Figure 4, most datasets in FinMTEB display inter-dataset similarity scores below 0.6, with a mean cosine similarity

of 0.4, indicating semantic distinctions among various types of financial texts.

4 Fin-E5: Finance-Adapted Text Embedding Model

Data is important for domain adaptation (Ling et al., 2023). However, existing public financial retrieval datasets exhibit significant limitations in their scope and applicability. For example, FiQA (FiQA, 2018), a widely used financial retrieval dataset, primarily focuses on opinion-based content from online platforms, neglecting crucial aspects such as fundamental financial knowledge, technical terminology, and important investment data. This narrow task focus creates a substantial gap in training comprehensive financial embedding models. Therefore, we use persona-based data generation to address this problem and synthesize a diverse range of tasks, as illustrated in Figure 3.

4.1 Data Formation

We aim to construct each training instance as a triplet structure (q, d^+, D^-) , where q represents a financial query, d^+ denotes a relevant document that provides substantive information addressing the query, and D^- comprises carefully selected negative examples that share the financial domain but differ in semantic intent.

4.2 Training Data Construction

To create a comprehensive dataset tailored for financial embedding training, we employ a systematic approach that combines expert-curated seed data with persona-based synthetic data generation.

Seed Data. Our seed data comes from the finance-specific QA dataset provided by InvestLM (Yang et al., 2023b), which offers expert-validated financial content across various domains, such as market analysis, investment strategies, and corporate finance. To ensure evaluation integrity, we conduct rigorous overlap checks between our training data and the FinMTEB benchmark, guaranteeing no overlap.

Persona-based Data Augmentation. To enhance the diversity of financial task representations, we develop a persona-based data augmentation framework derived from QA data generation (Ge et al., 2024). Our framework employs a three-stage process that specifically targets the expansion of task coverage while preserving domain consistency.

²<https://tianchi.aliyun.com/dataset/106411>

³<https://huggingface.co/sentence-transformers/all-MiniLM-L6-v2>

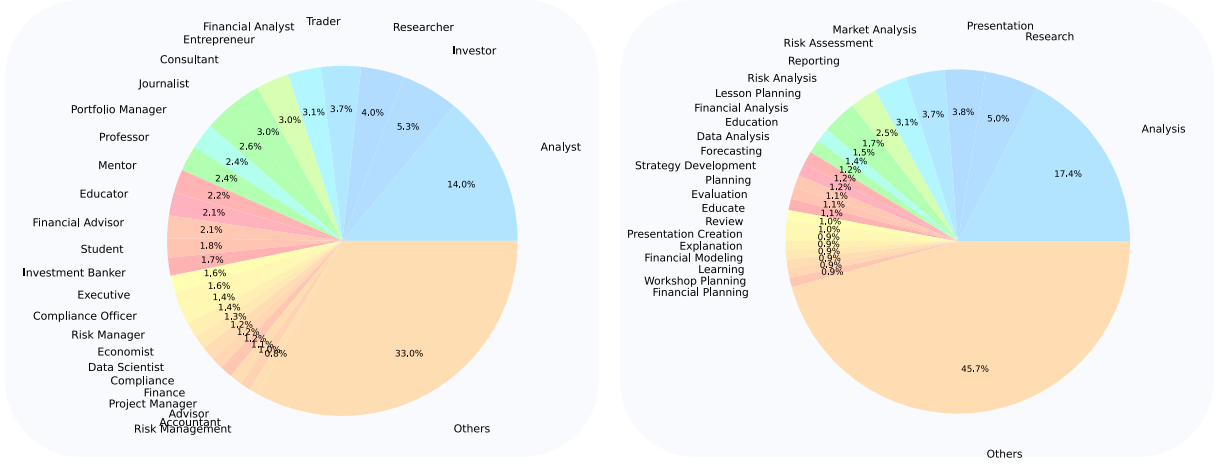


Figure 3: Distribution analysis of 5000 randomly sampled training data showing the breakdown of Tasks and Person Types. Left: Persona distribution. Right: Task distribution.

- **Persona and Task Identification:** We first employ Qwen2.5-72B-Instruct (Team, 2024) to analyze each question-answer pair in the seed data, aiming to identify the persona of the intended user (e.g., equity analyst, risk manager, financial advisor, retail investor) by using the prompt "Who is likely to use this text?" Different personas correspond to various tasks.
- **Contextual Query Expansion:** For each identified persona-task pair, we generate context-specific queries q that reflect the persona's unique information needs and risk preferences, using the prompt "Guess a prompt (i.e., instructions) that the following persona may ask you to do:". For example, a pension fund manager's query might emphasize long-term asset allocation, while a venture capitalist's query would prioritize startup valuation metrics.
- **Synthetic Document Generation:** We used LLMs to synthesize financial documents d^+ tailored to each persona's task, ensuring that the dataset represents diverse perspectives in financial decision-making. This step improves the representativeness of the dataset, ensuring that the embeddings are trained in real-world financial contexts. The prompt is "Please synthesize some real context information, which is related to this question:".

We randomly sample 5,000 data points from the training data, then use GPT-4o (OpenAI, 2024a) to annotate the job-related persona and task for the

query. Visualized in Figure 3, it is clear that our data generation process produces a diverse range of tasks and finance persona.

4.3 Training Pipeline

Following the training recipe of e5-mistral-7b-instruct (Wang et al., 2023), utilizing the last token pooling method, we construct training pairs by selecting queries as anchor points and their corresponding answers as positive samples. To enhance the effectiveness of contrastive learning, we identify challenging negative samples using the all-MiniLM-L12-v2 model (Reimers and Gurevych, 2019). The training process applies the InfoNCE loss (Oord et al., 2018), calculated over in-batch negative samples. The detailed training parameter is illustrated in Appendix C.

5 Experiment

In this section, we benchmark several existing models on FinMTEB, and then provide an in-depth analysis. Since most models are trained on English corpora, we only evaluate their performance on English datasets.

5.1 Models

In addition to Fin-E5, we also evaluate four categories of embedding models on the FinMTEB benchmark in Table 1. The benchmark time is reported in Appendix D.

Bag-of-Words (BOW). As a simple baseline, we implement the traditional BOW approach that represents text as sparse vectors based on word

Model	Size	Tasks							Avg.
		STS	Retrieval	Class.	Cluster.	Rerank.	PairClass.	Summ.	
		2	10	8	6	3	3	3	
BOW	-	0.4845	0.2084	0.4696	0.2547	0.7628	0.7143	0.0542	0.4212
Encoder based Models									
BERT	110M	0.3789	0.0207	0.5496	0.1744	0.3930	0.7111	0.0452	0.3247
FinBERT	110M	0.4198	0.1102	0.5923	0.2833	0.6404	0.6967	0.0417	0.3978
instructor-base	110M	0.3732	0.5772	0.6208	0.5300	0.9734	0.6138	0.1465	0.5479
bge-large-en-v1.5	335M	0.3396	0.6463	0.6436	0.5725	0.9825	0.7400	0.2019	0.5895
AngIE-BERT	335M	0.3080	0.5730	0.6439	0.5774	0.9650	0.6891	0.5049	0.6088
LLM-based Models									
gte-Qwen1.5-7B-instruct	7B	0.3758	0.6697	0.6438	0.5854	0.9890	0.6998	0.2350	0.5998
Echo	7B	<u>0.4380</u>	0.6443	0.6525	0.5776	0.9765	0.6261	0.4722	0.6267
bge-en-icl	7B	0.3233	0.6789	0.6569	0.5742	0.9898	0.6738	0.5197	0.6309
NV-Embed v2	7B	0.3739	0.7061	0.6393	0.6096	0.9822	0.6043	0.5103	0.6322
e5-mistral-7b-instruct	7B	0.3800	0.6749	0.6449	0.5783	0.9875	<u>0.7394</u>	0.5275	0.6475
Commercial Models									
text-embedding-3-small	-	0.3254	0.6641	0.6387	0.5802	0.9825	0.5957	0.5085	0.6136
text-embedding-3-large	-	0.3615	<u>0.7112</u>	0.6596	0.6081	<u>0.9910</u>	0.7309	<u>0.5671</u>	0.6613
voyage-3-large	-	0.4145	0.7463	<u>0.6861</u>	<u>0.5944</u>	0.9938	0.6519	0.6484	<u>0.6765</u>
Finance Adapted LLM-based Models									
Fin-E5	7B	0.4342	0.7105	0.7565	0.5650	0.9896	0.8014	0.4797	0.6767

Table 1: Performance comparison across different embedding models on FinMTEB benchmark. The evaluation metrics include semantic textual similarity (STS), retrieval, classification (Class.), clustering (Cluster.), reranking (Rerank.), pair classification (PairClass.), and summarization (Summ.). **Best** results are in bold. The underline represents the second-best performance.

frequencies, providing a reference point for comparing more sophisticated methods.

Encoder-based Models. We evaluate various transformer encoder architectures, including: (1) classical models like BERT (CLS pooling) (Devlin et al., 2019) and domain-specific FinBERT (Yang et al., 2020); (2) optimized models such as msmarco-bert-base-dot-v5 and all-MiniLM-L12-v2 (Reimers and Gurevych, 2019); and (3) advanced architectures including bge-large-en-v1.5 (Xiao et al., 2023), AngIE-BERT (Li and Li, 2023) and instructor-base (Su et al., 2022).

LLM-based Models. We investigate several state-of-the-art decoder-based embedding models: (1) Mistral-7B-based models including bge-en-icl (Xiao et al., 2023), e5-mistral-7b-instruct (Wang et al., 2023) and Echo (Springer et al., 2024); (2) NV-Embed v2 (Lee et al., 2024); and (3) gte-Qwen1.5-7B-instruct (Li et al., 2023) built on the Qwen2 (Yang et al., 2024) architecture. These models utilize the powerful representation capabilities of LLMs to generate high-quality embeddings.

Commercial Models. To provide a comprehensive comparison with commercial solutions, we include industry-leading closed-source models,

specifically OpenAI’s text-embedding-3-large, text-embedding-3-small (OpenAI, 2024b) and voyage-3-large (VoyageAI, 2025)⁴.

5.2 Analysis

Based on the results presented in Table 1, our analysis focuses on three key findings.

5.2.1 Impact of Domain Adaptation

As illustrated in Table 1, domain specialization considerably boosts performance: FinBERT outperforms BERT by 15.6% (0.6721 vs. 0.5812), while Fin-E5 exceeds its general-domain counterpart e5-mistral-7b-instruct by 4.5% (0.6767 vs. 0.6475), particularly excelling in classification (0.842 vs. 0.807) and semantic textual similarity (0.721 vs. 0.685). The finance-adapted Fin-E5 also achieves state-of-the-art performance (0.6767 average score) on the FinMTEB benchmark, exceeding both general-purpose and commercial models. Notably, this peak performance is achieved with just 100 training steps, showcasing a cost-effective adaptation without the risk of data leakage.

⁴We thank VoyageAI for supporting us in conducting the evaluation.

5.2.2 The Role of Model Architecture and Size

Our experiments reveal three distinct performance tiers across architectural paradigms (Table 1). Traditional bag-of-words (BOW) models achieve baseline performance (STS: 0.4845) and show notable limitations in retrieval tasks. Encoder-based architectures, such as bge-large-en-v1.5, demonstrate significant improvements, increasing retrieval performance by 107% (0.6463) and STS by 38% (0.6692) over BOW. A paradigm shift occurs with LLM-based models; e5-mistral-7b-instruct sets new standards with an average score of 0.6475. This progression from BOW (lexical) to LLM-based (contextual) architectures reveals a 52% overall performance improvement, suggesting that model capacity plays a critical role in capturing financial semantics.

5.2.3 Limitations of Current Models in Financial STS Tasks

The STS results reveal a counterintuitive finding: BOW models (0.4845) outperform all dense architectures (maximum 0.4342) in terms of financial document similarity. This reversal of typical NLP performance hierarchies arises from two characteristics of the corpus: (1) Extensive boilerplate content in annual reports introduces noise for contextual embeddings, and (2) Specialized terminology (27% unique financial terms per document) decreases lexical overlap. BOW benefits from exact term matches in standardized disclosures; the best dense model only captures 64% of human-annotated similarity relationships, revealing fundamental limitations in current strategies for financial documents.

6 Domain-specific Embedding Benchmark is needed

This section addresses another research question. *To what extent do general-purpose embedding evaluations appropriately capture domain-specific performance?* To solve this question, we run a quantitative comparison between MTEB (Muennighoff et al., 2022) and FinMTEB.

Models. We evaluate seven state-of-the-art general-purpose embedding model. Specifically, we consider the following models: bge-en-icl (Xiao et al., 2023) and e5-mistral-7b-instruct (Wang et al., 2023), which are developed from Mistral-7B-v0.1 (Jiang et al., 2023); gte-Qwen2-1.5B-instruct (Li et al., 2023), developed from Qwen2 (Yang et al.,

2024); bge-large-en-v1.5 (Xiao et al., 2023) and all-MiniLM-L12-v2 (Reimers and Gurevych, 2019), both developed from BERT (Devlin et al., 2019); instructor-base (Su et al., 2022) from T5Encoder (Raffel et al., 2020); and OpenAI’s text-embedding-3-small (OpenAI, 2024b).

Method. To ensure robust statistical analysis, we use bootstrapping methods to generate a large sample dataset. For each task in both MTEB and FinMTEB, we aggregate the datasets associated with the task into a task pool. From each task pool, we randomly select 50 examples to create a bootstrap sample and evaluate the embedding model’s performance on this bootstrap. We repeat this process 500 times, resulting in 500 bootstraps for each combination. Thus, we have 14 unique combinations (model and domain), each with 500 bootstraps and their corresponding performance scores.

Analysis of Variance. We conduct an Analysis of Variance (ANOVA) that examines the effects of both the model and the domain. The results reveal that the Domain Factor demonstrates statistical significance across all tasks ($p < 0.001$), with notably large F statistics in classification ($F = 2086.30$), clustering ($F = 32161.37$), and STS ($F = 25761.71$). Furthermore, the Domain Factor generally accounts for a greater share of the variance than the Model Factor, as indicated by the Sum of Squares (e.g., in Classification: Domain = 56.82 vs. Model = 4.17). These findings suggest that domain-specific characteristics significantly impact model performance, reinforcing the importance of specialized evaluation frameworks such as FinMTEB for financial applications.

7 Conclusion

This paper introduces FinMTEB, the first comprehensive benchmark for evaluating embedding models in the financial domain. Our main contributions include establishing a large-scale evaluation framework with 64 datasets across seven tasks in Chinese and English, and developing Fin-E5, a finance-adapted embedding model demonstrating competitive performance through persona-based data augmentation. Our empirical results highlight the importance of domain-specific adaptation and reveal current limitations in financial text embeddings. We believe FinMTEB will serve as a valuable resource for both researchers and practitioners in advancing financial language models.

8 Limitation

This work has two primary limitations. First, it relies on several existing financial datasets that could potentially overlap with the training data of contemporary embedding models. This overlap may introduce contamination, making it difficult to ensure completely fair comparisons between different models. Second, our adapted model and evaluation methods are currently limited to the English language, which restricts their applicability to non-English financial texts.

References

- Peter Anderson, Mano Vikash Janardhanan, Jason He, Wei Cheng, and Charlie Flanagan. 2024. *Greenback bears and fiscal hawks: Finance is a jungle and text embeddings must adapt*. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing: Industry Track*, pages 362–370, Miami, Florida, US. Association for Computational Linguistics.
- Elliot Bolton, Abhinav Venigalla, Michihiro Yasunaga, David Hall, Betty Xiong, Tony Lee, Roxana Daneshjou, Jonathan Frankle, Percy Liang, Michael Carbin, et al. 2024. Biomedlm: A 2.7 b parameter language model trained on biomedical text. *arXiv preprint arXiv:2403.18421*.
- CCKS. 2022. Ccks2022: Few-shot event extraction for the financial sector. https://www.biendata.xyz/competition/ccks2022_eventext/.
- CFPB. 2024. Consumer finance complaints. <https://huggingface.co/datasets/CFPB/consumer-finance-complaints>.
- Jing Chen, Qingcai Chen, Xin Liu, Haijun Yang, Daohe Lu, and Buzhou Tang. 2018. *The BQ corpus: A large-scale domain-specific Chinese corpus for sentence semantic equivalence identification*. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 4946–4951, Brussels, Belgium. Association for Computational Linguistics.
- Qingyu Chen, Yifan Peng, and Zhiyong Lu. 2019. Biosentvec: creating sentence embeddings for biomedical texts. In *2019 IEEE International Conference on Healthcare Informatics (ICHI)*, pages 1–5. IEEE.
- Wei Chen, Qiushi Wang, Zefei Long, Xianyin Zhang, Zhongtian Lu, Bingxuan Li, Siyuan Wang, Jiarong Xu, Xiang Bai, Xuanjing Huang, and Zhongyu Wei. 2023. Disc-finllm: A chinese financial large language model based on multiple experts fine-tuning. *arXiv preprint arXiv:2310.15205*.

- Zhiyu Chen, Wenhui Chen, Charese Smiley, Sameena Shah, Iana Borova, Dylan Langdon, Reema Moussa, Matt Beane, Ting-Hao Huang, Bryan Routledge, and William Yang Wang. 2021. *FinQA: A dataset of numerical reasoning over financial data*. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 3697–3711, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Pierre Colombo, Telmo Pessoa Pires, Malik Boudiaf, Dominic Culver, Rui Melo, Caio Corro, Andre FT Martins, Fabrizio Esposito, Vera Lúcia Raposo, Sofia Morgado, et al. 2024. Saullm-7b: A pioneering large language model for law. *arXiv preprint arXiv:2403.03883*.
- Keith Cortis, André Freitas, Tobias Daudert, Manuela Huerlimann, Manel Zarrouk, Siegfried Handschuh, and Brian Davis. 2017. *SemEval-2017 task 5: Fine-grained sentiment analysis on financial microblogs and news*. In *Proceedings of the 11th International Workshop on Semantic Evaluation (SemEval-2017)*, pages 519–535, Vancouver, Canada. Association for Computational Linguistics.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. *BERT: Pre-training of deep bidirectional transformers for language understanding*. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Mahmoud El-Haj, Nadhem Zmandar, Paul Rayson, Ahmed AbuRa’ed, Marina Litvak, Nikiforos Pittaras, George Giannakopoulos, Aris Kosmopoulos, Blanca Carbajo-Coronado, and Antonio Moreno-Sandoval. 2022. The financial narrative summarisation shared task (FNS 2022). In *Proceedings of the 4th Financial Narrative Processing Workshop @LREC2022*, pages 43–52, Marseille, France. European Language Resources Association.
- Zhiwei Fei, Xiaoyu Shen, Dawei Zhu, Fengzhe Zhou, Zhuo Han, Songyang Zhang, Kai Chen, Zongwen Shen, and Jidong Ge. 2023. Lawbench: Benchmarking legal knowledge of large language models. *arXiv preprint arXiv:2309.16289*.
- FiQA. 2018. Financial question answering. <https://sites.google.com/view/fiqa>.
- Tao Ge, Xin Chan, Xiaoyang Wang, Dian Yu, Haitao Mi, and Dong Yu. 2024. Scaling synthetic data creation with 1,000,000,000 personas. *arXiv preprint arXiv:2406.20094*.
- Daniela Gerz, Pei-Hao Su, Razvan Kusztos, Avishek Mondal, Michał Lis, Eshan Singhal, Nikola Mrkšić, Tsung-Hsien Wen, and Ivan Vulić. 2021a. *Multilingual and cross-lingual intent detection from spoken data*. In *Proceedings of the 2021 Conference on*

748	<i>Empirical Methods in Natural Language Processing</i> ,	Xiang Li, Zhenyu Li, Chen Shi, Yong Xu, Qing	803
749	pages 7468–7475, Online and Punta Cana, Domini-	Du, Mingkui Tan, Jun Huang, and Wei Lin. 2024.	804
750	can Republic. Association for Computational Lin-	Alphafin: Benchmarking financial analysis with	805
751	guistics.	retrieval-augmented stock-chain framework. <i>arXiv</i>	806
		<i>preprint arXiv:2403.12582</i> .	807
752	Daniela Gerz, Pei-Hao Su, Razvan Kusztoş, Avishek		
753	Mondal, Michał Lis, Eshan Singhal, Nikola Mrkšić,	Xianming Li and Jing Li. 2023. Angle-optimized text	808
754	Tsung-Hsien Wen, and Ivan Vulić. 2021b. Multilin-	embeddings. <i>arXiv preprint arXiv:2309.12871</i> .	809
755	gual and cross-lingual intent detection from spoken		
756	data. <i>arXiv preprint arXiv:2104.08524</i> .	Zehan Li, Xin Zhang, Yanzhao Zhang, Dingkun Long,	810
		Pengjun Xie, and Meishan Zhang. 2023. Towards	811
757	Biyang Guo, Xin Zhang, Ziyuan Wang, Minqi Jiang,	general text embeddings with multi-stage contrastive	812
758	Jinran Nie, Yuxuan Ding, Jianwei Yue, and Yupeng	learning. <i>arXiv preprint arXiv:2308.03281</i> .	813
759	Wu. 2023. How close is chatgpt to human experts?		
760	comparison corpus, evaluation, and detection. <i>arXiv</i>	Chen Ling, Xujiang Zhao, Jiaying Lu, Chengyuan Deng,	814
761	<i>preprint arXiv:2301.07597</i> .	Can Zheng, Junxiang Wang, Tanmoy Chowdhury,	815
		Yun Li, Hejie Cui, Xuchao Zhang, et al. 2023. Do-	816
762	Suchin Gururangan, Ana Marasović, Swabha	main specialization as the key to make large language	817
763	Swayamdipta, Kyle Lo, Iz Beltagy, Doug Downey,	models disruptive: A comprehensive survey. <i>arXiv</i>	818
764	and Noah A. Smith. 2020. Don't stop pretraining:	<i>preprint arXiv:2305.18703</i> .	819
765	Adapt language models to domains and tasks . In		
766	<i>Proceedings of the 58th Annual Meeting of the</i>	Jiaxin Liu, Yi Yang, and Kar Yan Tam. 2024a. Beyond	820
767	<i>Association for Computational Linguistics</i> , pages	surface similarity: Detecting subtle semantic shifts	821
768	8342–8360, Online. Association for Computational	in financial narratives . In <i>Findings of the Association</i>	822
769	Linguistics.	<i>for Computational Linguistics: NAACL 2024</i> , pages	823
		2641–2652, Mexico City, Mexico. Association for	824
770	Pranab Islam, Anand Kannappan, Douwe Kiela, Re-	Computational Linguistics.	825
771	becca Qian, Nino Scherrer, and Bertie Vidgen. 2023.		
772	Financebench: A new benchmark for financial ques-	Mianxin Liu, Jinru Ding, Jie Xu, Weiguo Hu, Xiaoyang	826
773	tion answering. <i>arXiv preprint arXiv:2311.11944</i> .	Li, Lifeng Zhu, Zhian Bai, Xiaoming Shi, Benyou	827
		Wang, Haitao Song, et al. 2024b. Medbench: A com-	828
774	Albert Q Jiang, Alexandre Sablayrolles, Arthur Men-	prehensive, standardized, and reliable benchmarking	829
775	sch, Chris Bamford, Devendra Singh Chaplot, Diego	system for evaluating chinese medical large language	830
776	de las Casas, Florian Bressand, Gianna Lengyel, Guil-	models. <i>arXiv preprint arXiv:2407.10990</i> .	831
777	laume Lample, Lucile Saulnier, et al. 2023. Mistral		
778	7b. <i>arXiv preprint arXiv:2310.06825</i> .	Shuaiqi Liu, Jiannong Cao, Ruosong Yang, and Zhiyuan	832
		Wen. 2022. Long text and multi-table summarization:	833
779	Jia-Huei Ju, Yu-Shiang Huang, Cheng-Wei Lin, Che Lin,	Dataset and method . In <i>Findings of the Association</i>	834
780	and Chuan-Ju Wang. 2023. A compare-and-contrast	<i>for Computational Linguistics: EMNLP 2022</i> , pages	835
781	multistage pipeline for uncovering financial signals	1995–2010, Abu Dhabi, United Arab Emirates. Asso-	836
782	in financial reports . In <i>Proceedings of the 61st An-</i>	ciation for Computational Linguistics.	837
783	<i>annual Meeting of the Association for Computational</i>		
784	<i>Linguistics (Volume 1: Long Papers)</i> , pages 14307–	Yinhan Liu. 2019a. Roberta: A robustly opti-	838
785	14321, Toronto, Canada. Association for Computa-	mized bert pretraining approach. <i>arXiv preprint</i>	839
786	tional Linguistics.	<i>arXiv:1907.11692</i> , 364.	840
787	Yanis Labrak, Adrien Bazoge, Oumaima El Khettari,	Yinhan Liu. 2019b. Roberta: A robustly opti-	841
788	Mickaël Rouvier, Natalia Grabar, Beatrice Daille,	mized bert pretraining approach. <i>arXiv preprint</i>	842
789	Solen Quiniou, Emmanuel Morin, Pierre-Antoine	<i>arXiv:1907.11692</i> .	843
790	Gourraud, Richard Dufour, et al. 2024. Drbench-		
791	mark: A large language understanding evaluation	Dakuan Lu, Hengkui Wu, Jiaqing Liang, Yipei Xu,	844
792	benchmark for french biomedical domain. <i>arXiv</i>	Qianyu He, Yipeng Geng, Mengkun Han, Yingsi	845
793	<i>preprint arXiv:2402.13432</i> .	Xin, and Yanghua Xiao. 2023. Bbt-fin: Compre-	846
		hensive construction of chinese financial domain	847
794	Yinyu Lan, Yanru Wu, Wang Xu, Weiqiang Feng, and	pre-trained language model, corpus and benchmark.	848
795	Youhao Zhang. 2023. Chinese fine-grained financial	<i>arXiv preprint arXiv:2302.09432</i> .	849
796	sentiment analysis with large language models. <i>arXiv</i>		
797	<i>preprint arXiv:2306.14096</i> .	Lorenzo Malandri, Fabio Mercorio, Mario Mezzan-	850
		zanica, and Filippo Pallucchini. 2025. RE-FIN:	851
798	Chankyu Lee, Rajarshi Roy, Mengyao Xu, Jonathan	Retrieval-based enrichment for financial data . In	852
799	Raiman, Mohammad Shoeibi, Bryan Catanzaro, and	<i>Proceedings of the 31st International Conference on</i>	853
800	Wei Ping. 2024. Nv-embed: Improved techniques for	<i>Computational Linguistics: Industry Track</i> , pages	854
801	training llms as generalist embedding models. <i>arXiv</i>	751–759, Abu Dhabi, UAE. Association for Compu-	855
802	<i>preprint arXiv:2405.17428</i> .	tational Linguistics.	856

857	Pekka Malo, Ankur Sinha, Pekka Korhonen, Jyrki Wal-	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine	911
858	lenius, and Pyry Takala. 2014. Good debt or bad	Lee, Sharan Narang, Michael Matena, Yanqi Zhou,	912
859	debt: Detecting semantic orientations in economic	Wei Li, and Peter J Liu. 2020. Exploring the lim-	913
860	texts. <i>Journal of the Association for Information</i>	its of transfer learning with a unified text-to-text	914
861	<i>Science and Technology</i> , 65(4):782–796.	transformer. <i>Journal of machine learning research</i> ,	915
		21(140):1–67.	916
862	Rui Meng, Ye Liu, Shafiq Rayhan Joty, Caiming	Aman Rangapur, Haoran Wang, Ling Jian, and Kai Shu.	917
863	Xiong, Yingbo Zhou, and Semih Yavuz. 2024. <i>Sfr-</i>	2023. Fin-fact: A benchmark dataset for multimodal	918
864	<i>embedding-2: Advanced text embedding with multi-</i>	financial fact checking and explanation generation.	919
865	<i>stage training</i> .	<i>arXiv preprint arXiv:2309.08793</i> .	920
866	Tomas Mikolov, Kai Chen, Greg Corrado, and Jef-	Nils Reimers and Iryna Gurevych. 2019. <i>Sentence-bert:</i>	921
867	frey Dean. 2013. <i>Efficient estimation of word</i>	<i>Sentence embeddings using siamese bert-networks</i> .	922
868	<i>representations in vector space</i> . <i>arXiv preprint</i>	In <i>Proceedings of the 2019 Conference on Empirical</i>	923
869	<i>arXiv:1301.3781</i> .	<i>Methods in Natural Language Processing</i> . Associa-	924
		tion for Computational Linguistics.	925
870	Niklas Muennighoff, Nouamane Tazi, Loïc Magne, and	Andrew Rosenberg and Julia Hirschberg. 2007. V-	926
871	Nils Reimers. 2022. Mteb: Massive text embedding	measure: A conditional entropy-based external clus-	927
872	benchmark. <i>arXiv preprint arXiv:2210.07316</i> .	ter evaluation measure. In <i>Proceedings of the 2007</i>	928
873	Rajdeep Mukherjee, Abhinav Bohra, Akash Banerjee,	<i>Joint Conference on Empirical Methods in Natural</i>	929
874	Soumya Sharma, Manjunath Hegde, Afreen Shaikh,	<i>Language Processing and Computational Natural</i>	930
875	Shivani Shrivastava, Koustuv Dasgupta, Niloy Gan-	<i>Language Learning (EMNLP-CoNLL)</i> , pages 410–	931
876	guly, Saptarshi Ghosh, et al. 2022. Ectsum: A new	420, Prague, Czech Republic. Association for Com-	932
877	benchmark dataset for bullet point summarization	putational Linguistics.	933
878	of long earnings call transcripts. <i>arXiv preprint</i>		
879	<i>arXiv:2210.12467</i> .	Agam Shah, Suvan Paturi, and Sudheer Chava. 2023.	934
880	Qiong Nan, Juan Cao, Yongchun Zhu, Yanyan Wang,	<i>Trillion dollar words: A new financial dataset, task &</i>	935
881	and Jintao Li. 2021. Mdfend: Multi-domain fake	<i>market analysis</i> . In <i>Proceedings of the 61st Annual</i>	936
882	news detection. In <i>Proceedings of the 30th ACM In-</i>	<i>Meeting of the Association for Computational Lin-</i>	937
883	<i>ternational Conference on Information & Knowledge</i>	<i>guistics (Volume 1: Long Papers)</i> , pages 6664–6679,	938
884	<i>Management</i> , pages 3343–3347.	Toronto, Canada. Association for Computational Lin-	939
885	Aaron van den Oord, Yazhe Li, and Oriol Vinyals. 2018.	guistics.	940
886	Representation learning with contrastive predictive	Ankur Sinha and Tanmay Khandait. 2021. Impact of	941
887	coding. <i>arXiv preprint arXiv:1807.03748</i> .	news on the commodity market: Dataset and results.	942
888	OpenAI. 2024a. Openai (august 24 version). https:	In <i>Advances in Information and Communication:</i>	943
889	https://api.openai.com/v1/chat .	<i>Proceedings of the 2021 Future of Information and</i>	944
890	OpenAI. 2024b. Openai (august 24 version). https:	<i>Communication Conference (FICC)</i> , Volume 2, pages	945
891	https://api.openai.com/v1/embeddings .	589–601. Springer.	946
892	Masanori Oya. 2011. Syntactic dependency distance	Jacob Mitchell Springer, Suhas Kotha, Daniel Fried,	947
893	as sentence complexity measure. In <i>Proceedings</i>	Graham Neubig, and Aditi Raghunathan. 2024. Rep-	948
894	<i>of the 16th International Conference of Pan-Pacific</i>	etition improves language model embeddings. <i>arXiv</i>	949
895	<i>Association of Applied Linguistics</i> , volume 1.	<i>preprint arXiv:2402.15449</i> .	950
896	Jeffrey Pennington, Richard Socher, and Christopher	Hongjin Su, Weijia Shi, Jungo Kasai, Yizhong Wang,	951
897	Manning. 2014. <i>GloVe: Global vectors for word</i>	Yushi Hu, Mari Ostendorf, Wen-tau Yih, Noah A	952
898	<i>representation</i> . In <i>Proceedings of the 2014 Confer-</i>	Smith, Luke Zettlemoyer, and Tao Yu. 2022. One	953
899	<i>ence on Empirical Methods in Natural Language Pro-</i>	embedder, any task: Instruction-finetuned text em-	954
900	<i>cessing (EMNLP)</i> , pages 1532–1543, Doha, Qatar.	beddings. <i>arXiv preprint arXiv:2212.09741</i> .	955
901	Association for Computational Linguistics.	Maosong Sun, Jingyang Li, Zhipeng Guo, Yu Zhao,	956
902	Matthew E. Peters, Mark Neumann, Mohit Iyyer, Matt	Yabin Zheng, Xiance Si, and Zhiyuan Liu. 2016.	957
903	Gardner, Christopher Clark, Kenton Lee, and Luke	Thuctc: An efficient chinese text classifier. http:	958
904	Zettlemoyer. 2018. <i>Deep contextualized word repre-</i>	http://thuctc.thunlp.org/ .	959
905	<i>sentations</i> . In <i>Proceedings of the 2018 Conference of</i>	Qwen Team. 2024. <i>Qwen2.5: A party of foundation</i>	960
906	<i>the North American Chapter of the Association for</i>	<i>models</i> .	961
907	<i>Computational Linguistics: Human Language Tech-</i>	Nandan Thakur, Nils Reimers, Andreas Rücklé, Ab-	962
908	<i>nologies, Volume 1 (Long Papers)</i> , pages 2227–2237,	hishek Srivastava, and Iryna Gurevych. 2021. BEIR:	963
909	New Orleans, Louisiana. Association for Computa-	A heterogeneous benchmark for zero-shot evaluation	964
910	tional Linguistics.	of information retrieval models. In <i>Thirty-fifth Con-</i>	965
		<i>ference on Neural Information Processing Systems</i>	966
		<i>Datasets and Benchmarks Track (Round 2)</i> .	967

968	VoyageAI. 2025. Voyageai (jan 25 version). https://api.voyageai.com/v1/embeddings .	for Computational Linguistics: ACL-IJCNLP 2021, pages 2114–2124, Online. Association for Computational Linguistics.	1022 1023 1024
970	Liang Wang, Nan Yang, Xiaolong Huang, Linjun Yang, Rangan Majumder, and Furu Wei. 2023. Improving text embeddings with large language models. <i>arXiv preprint arXiv:2401.00368</i> .	Fengbin Zhu, Wenqiang Lei, Youcheng Huang, Chao Wang, Shuo Zhang, Jiancheng Lv, Fuli Feng, and Tat-Seng Chua. 2021. Tat-qa: A question answering benchmark on a hybrid of tabular and textual content in finance. <i>arXiv preprint arXiv:2105.07624</i> .	1025 1026 1027 1028 1029
974	Alex Watson, Yev Meyer, Maarten Van Segbroeck, Matthew Grossman, Sami Torbey, Piotr Mlocek, and Johnny Greco. 2024. Synthetic-P11-Financial-Documents-North-America: A synthetic dataset for training language models to label and detect pii in domain specific formats .	Jie Zhu, Junhui Li, Yalong Wen, and Lifan Guo. 2024. Benchmarking large language models on cflue—a chinese financial language understanding evaluation dataset. <i>arXiv preprint arXiv:2405.10542</i> .	1030 1031 1032 1033
980	Shijie Wu, Ozan Irsoy, Steven Lu, Vadim Dabrovolski, Mark Dredze, Sebastian Gehrmann, Prabhajan Kam-badur, David Rosenberg, and Gideon Mann. 2023. Bloomberggpt: A large language model for finance. <i>arXiv preprint arXiv:2303.17564</i> .	A Datasets	1034
981		The detailed description of each dataset used in this work is listed in the Table tables 2 to 8.	1035 1036
982		B Dataset Characteristic	1037
983		Figure 4 presents the semantic similarity across all datasets in the FinMTEB benchmark. The semantic similarity is calculated by cosine similarity. Table 9 presents a comparative analysis of linguistic features between MTEB (Muennighoff et al., 2022) and FinMTEB benchmarks, examining aspects such as average sentence length, token length, syllables per token, and dependency distance (Oya, 2011).	1038 1039 1040 1041 1042 1043 1044 1045 1046 1047 1048
984		C Training Details For Fin-E5	1049
985	Shitao Xiao, Zheng Liu, Peitian Zhang, and Niklas Muennighoff. 2023. C-pack: Packaged resources to advance general chinese embedding. <i>arXiv preprint arXiv:2309.07597</i> .	The training dataset size is 19,467. The model is trained for 100 steps using the augmented dataset with a batch size of 128. For optimization, we use the AdamW optimizer with a learning rate of 1e-5 and implement a linear warmup schedule. For a given data (q, d^+, D^-) , we adopt an instruction-based methodology for embedding training. The instruction template is as follows:	1050 1051 1052 1053 1054 1055 1056 1057
986		$q_{\text{inst}} = \text{Instruct: } \{task_definition\} \backslash n \{q\} \quad (1)$	1058
987		where $\{task_definition\}$ represents a concise single-sentence description of the embedding task.	1059 1060
988		D Benchmarking Time Reporting.	1061
989	Ziyue Xu, Peilin Zhou, Xinyu Shi, Jiageng Wu, Yikang Jiang, Bin Ke, and Jie Yang. 2024. Fintruthqa: A benchmark dataset for evaluating the quality of financial information disclosure. <i>arXiv preprint arXiv:2406.12009</i> .	The benchmarking was conducted on the NVIDIA H800 GPU using a batch size of 512. Echo Embedding (Springer et al., 2024) required the longest processing time at 12 hours, followed by BeLLM (Li	1062 1063 1064 1065
994	An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, et al. 2024. Qwen2 technical report. <i>arXiv preprint arXiv:2407.10671</i> .		
995			
996			
997			
998	Hongyang Yang, Xiao-Yang Liu, and Christina Dan Wang. 2023a. Fingpt: Open-source financial large language models. <i>arXiv preprint arXiv:2306.06031</i> .		
999			
1000			
1001	Yi Yang, Yixuan Tang, and Kar Yan Tam. 2023b. Investlm: A large language model for investment using financial domain instruction tuning. <i>arXiv preprint arXiv:2309.13064</i> .		
1002			
1003			
1004			
1005	Yi Yang, Mark Christopher Siy Uy, and Allen Huang. 2020. Finbert: A pretrained language model for financial communications. <i>arXiv preprint arXiv:2006.08097</i> .		
1006			
1007			
1008			
1009	Liwen Zhang, Weige Cai, Zhaowei Liu, Zhi Yang, Wei Dai, Yujie Liao, Qianru Qin, Yifei Li, Xingyu Liu, Zhiqiang Liu, et al. 2023. Fineval: A chinese financial domain knowledge evaluation benchmark for large language models. <i>arXiv preprint arXiv:2308.09975</i> .		
1010			
1011			
1012			
1013			
1014			
1015	Yijia Zhang, Qingyu Chen, Zhihao Yang, Hongfei Lin, and Zhiyong Lu. 2019. Biowordvec, improving biomedical word embeddings with subword information and mesh. <i>Scientific data</i> , 6(1):52.		
1016			
1017			
1018			
1019	Zhihan Zhou, Liqian Ma, and Han Liu. 2021. Trade the event: Corporate events detection for news-based event-driven trading. In <i>Findings of the Association</i>		
1020			
1021			

Dataset Name	Language	Description
FINAL (Ju et al., 2023)	English	A dataset designed for discovering financial signals in narrative financial reports.
FinSTS (Liu et al., 2024a)	English	A dataset focused on detecting subtle semantic shifts in financial narratives.
AFQMC ⁵	Chinese	A Chinese dataset for customer service question matching in the financial domain.
BQ-Corpus (Chen et al., 2018)	Chinese	A large-scale Chinese corpus for sentence semantic equivalence identification (SSEI) in the banking domain.

Table 2: Summary of STS Datasets

and Li, 2023) at 11.98 hours. AnglE-BERT (Li and Li, 2023) completed the evaluation in 8 hours, while NV-Embed v2 (Lee et al., 2024) demonstrated the highest efficiency, completing all tasks in just 5.6 hours.

E Spearman’s Correlation of Embedding Models’ Performance

We evaluate the performance ranking of embedding models on both the general MTEB and FinMTEB datasets, calculating Spearman’s rank correlation between the two. The results, shown in Table 10, indicate that the ranking correlation is not statistically significant (p-values all greater than 0.05). In other words, a general-purpose embedding model performing well on MTEB does not necessarily perform well on domain-specific tasks.

F Analysis of Variance (ANOVA)

Table 11 illustrates the full results of ANOVA analysis.

⁷<https://tianchi.aliyun.com/dataset/106411>

⁸<https://lighthouz.ai/blog/rag-benchmark-finance-apple-10K-2022/>

⁹<https://www.kaggle.com/datasets/jeet2016/us-financial-news-articles>

¹⁰https://github.com/alipay/financial_evaluation_dataset/tree/main

¹¹<https://github.com/smoothnlp/SmoothNLP>

¹²https://github.com/alipay/financial_evaluation_dataset/tree/main

¹³<https://github.com/amitkedia007/Financial-Fraud-Detection-Using-LLMs/tree/main>

¹⁴<https://github.com/open-compass/OpenFinData?tab=readme-ov-file>

Dataset Name	Language	Description
FiQA2018 (FiQA, 2018)	English	Financial opinion mining and question answering dataset.
FinanceBench (Islam et al., 2023)	English	Open book financial question answering dataset.
HC3(Finance) (Guo et al., 2023)	English	A human-ChatGPT comparison corpus in the finance domain.
Apple-10K-2022 ⁶	English	A retrieval-augmented generation (RAG) benchmark for finance applications.
FinQA (Chen et al., 2021)	English	Financial numerical reasoning dataset with structured and unstructured evidence.
TAT-QA (Zhu et al., 2021)	English	Question answering benchmark combining tabular and textual content in finance.
US Financial News ⁷	English	Finance news articles paired with headlines and stock ticker symbols.
TradeTheEvent (Trading Benchmark) (Zhou et al., 2021)	English	Finance news articles paired with headlines and stock ticker symbols.
TradeTheEvent (Domain Adaption) (Zhou et al., 2021)	English	Financial terms and explanations dataset.
TheGoldman-en	English	English version of the Goldman Sachs Financial Dictionary.
FinTruthQA (Xu et al., 2024)	Chinese	Dataset for evaluating the quality of financial information disclosure.
Fin-Eva (Retrieval task) ⁸	Chinese	Financial scenario QA dataset focusing on retrieval tasks.
AlphaFin (Li et al., 2024)	Chinese	Comprehensive financial dataset including NLI, QA, and stock trend predictions.
DISC-FinLLM (Retrieval Part Data) (Chen et al., 2023)	Chinese	Financial scenario QA dataset.
FinQA (from DuEE-fin) (Lu et al., 2023)	Chinese	Financial news bulletin event quiz dataset.
DISC-FinLLM (Computing) (Chen et al., 2023)	Chinese	Financial scenario QA dataset focusing on numerical tasks.
SmoothNLP ⁹	Chinese	Chinese finance news dataset.
THUCNews (Sun et al., 2016)	Chinese	Chinese finance news dataset.
Fin-Eva (Terminology) ¹⁰	Chinese	Financial terminology dataset used in the industry.
TheGoldman-cn	Chinese	Chinese version of the Goldman Sachs Financial Dictionary.

Table 3: Summary of Retrieval Datasets

Dataset Name	Language	Description
FinancialPhrasebank (Malo et al., 2014)	English	Polar sentiment dataset of sentences from financial news, categorized by sentiment into positive, negative, or neutral.
FinSent (Yang et al., 2023b)	English	Polar sentiment dataset of sentences from the financial domain, categorized by sentiment into positive, negative, or neutral.
FiQA_ABSA (FiQA, 2018)	English	Polar sentiment dataset of sentences from the financial domain, categorized by sentiment into positive, negative, or neutral.
SemEva2017_Headline (Cortis et al., 2017)	English	Polar sentiment dataset of sentences from the financial domain, categorized by sentiment into positive, negative, or neutral.
FLS (Yang et al., 2023b)	English	A finance dataset detects whether the sentence is a forward-looking statement.
ESG (Yang et al., 2023b)	English	A finance dataset performs sentence classification under the environmental, social, and corporate governance (ESG) framework.
FOMC (Shah et al., 2023)	English	A task of hawkish-dovish classification in finance domain.
Financial-Fraud ¹¹	English	This dataset was used for research in detecting financial fraud.
FinNSP (Lu et al., 2023)	Chinese	Financial negative news and its subject determination dataset.
FinChina (Lan et al., 2023)	Chinese	Polar sentiment dataset of sentences from the financial domain, categorized by sentiment into positive, negative, or neutral.
FinFE (Lu et al., 2023)	Chinese	Financial social media text sentiment categorization dataset.
OpenFinData ¹²	Chinese	Financial scenario QA dataset including sentiment task.
MDFEND-Weibo2 (finance) (Nan et al., 2021)	Chinese	Fake news detection in the finance domain.

Table 4: Summary of Classification Datasets

Dataset Name	Language	Description
MInDS-14-en (Gerz et al., 2021b)	English	MINDS-14 is a dataset for intent detection in e-banking, covering 14 intents across 14 languages.
Consumer Complaints (CFPB, 2024)	English	The Consumer Complaint Database is a collection of complaints about consumer financial products and services that sent to companies for response.
Synthetic PII finance (Watson et al., 2024)	English	Synthetic financial documents containing Personally Identifiable Information (PII).
FinanceArxiv-s2s ¹³	English	Clustering of titles from arxiv (q-fin).
FinanceArxiv-p2p	English	Clustering of abstract from arxiv (q-fin).
WikiCompany2Industry-en ¹⁴	English	Clustering the related industry domain according to the company description.
MInDS-14-zh (Gerz et al., 2021b)	Chinese	MINDS-14 is a dataset for intent detection in e-banking, covering 14 intents across 14 languages.
FinNL (Lu et al., 2023)	Chinese	Financial news categorization dataset.
CCKS2022 (CCKS, 2022)	Chinese	Clustering of financial events.
CCKS2020	Chinese	Clustering of financial events.
CCKS2019	Chinese	Clustering of financial events.

Table 5: Summary of Clustering Datasets

Dataset Name	Language	Description
Ectsum (Mukherjee et al., 2022)	English	A Dataset For Bullet Point Summarization of Long Earnings Call Transcripts.
FINDSum (Liu et al., 2022)	English	A Large-Scale Dataset for Long Text and Multi-Table Summarization.
FNS-2022 (El-Haj et al., 2022)	English	Financial Narrative Summarisation for 10K.
FiNNA (Lu et al., 2023)	Chinese	A financial news summarization dataset.
Fin-Eva (Headline) (Zhang et al., 2023)	Chinese	A financial summarization dataset.
Fin-Eva (Abstract) (Zhang et al., 2023)	Chinese	A financial summarization dataset.

Table 6: Summary of Summarization Datasets

Dataset Name	Language	Description
Fin-Fact (Rangapur et al., 2023)	English	A Benchmark Dataset for Financial Fact Checking and Explanation Generation.
FiQA2018 (FiQA, 2018)	English	Financial opinion mining and question answering.
HC3(Finance) (Guo et al., 2023)	English	A human-ChatGPT comparison finance corpus.
Fin-Eva (Retrieval task) (Zhang et al., 2023)	Chinese	Financial scenario QA dataset including retrieval task.
DISC-FinLLM (Retrieval Part Data) (Chen et al., 2023)	Chinese	Financial scenario QA dataset.

Table 7: Summary of Reranking Datasets

Dataset Name	Language	Description
HeadlineAC-PairClassification (Sinha and Khandait, 2021)	English	Financial text sentiment categorization dataset.
HeadlinePDD-PairClassification (Sinha and Khandait, 2021)	English	Financial text sentiment categorization dataset.
HeadlinePDU-PairClassification (Sinha and Khandait, 2021)	English	Financial text sentiment categorization dataset.
AFQMC	Chinese	Ant Financial Question Matching Corpus.

Table 8: Summary of PairClassification Datasets

Benchmark	Sentence Length	Token Length	Syllables Per Token	Dependency Distance
MTEB	18.20	4.89	1.49	2.49
FinMTEB	26.37	5.12	1.52	2.85

Table 9: Comparison of Text Characteristics Between FinMTEB and MTEB. The numbers represent the average scores across all samples from all datasets.

	STS	Class.	Ret.	Rerank.	Clust.	PairClass.	Summ.
Correlation	0.30	-0.80	0.30	-0.10	-0.70	-0.30	0.60
p-value	0.62	0.10	0.62	0.87	0.18	0.62	0.28

Table 10: Spearman’s correlation of embedding models’ performance on MTEB and FinMTEB across different tasks. The p-value indicates that all correlations are statistically insignificant, suggesting a lack of evidence for a relationship between embedding model performance on the two benchmarks.

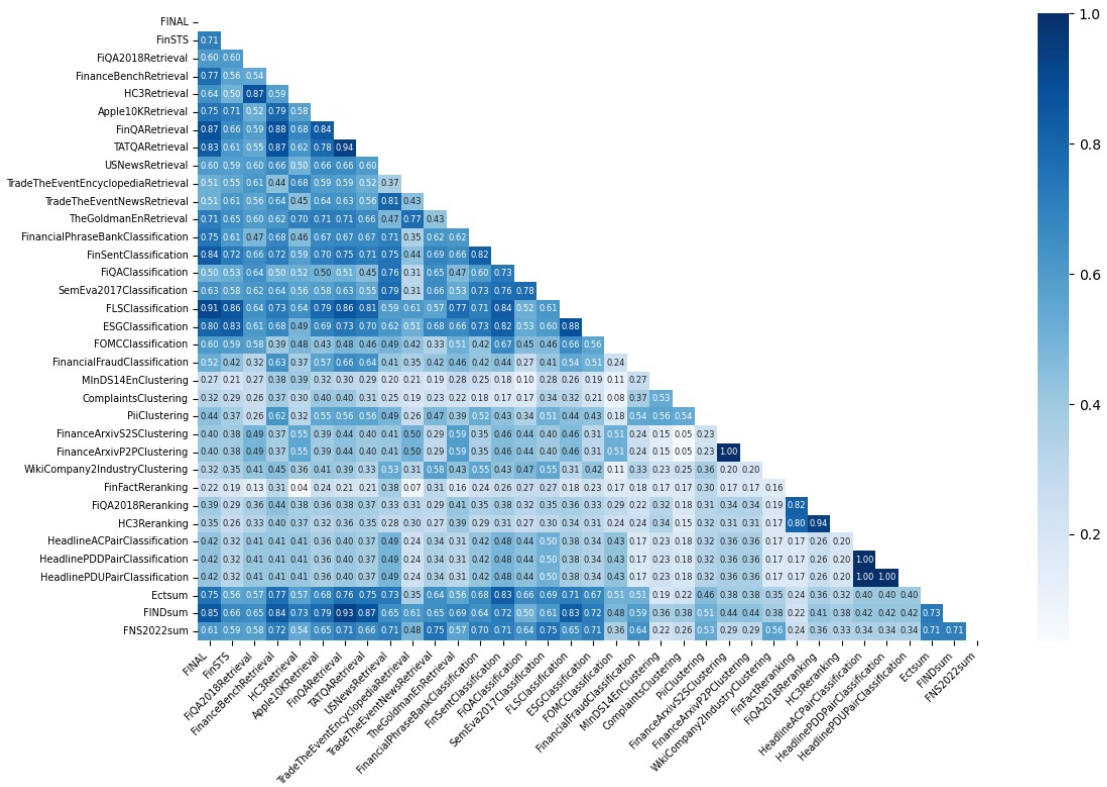


Figure 4: Semantic similarity across all the datasets in FinMTEB benchmark.

Task	Factor	Sum of Squares	Degrees of Freedom	F-Statistic	p-value
Classification	Model Factor	4.17	6.00	25.55	3.41×10^{-30}
	Domain Factor	56.82	1.00	2086.30	≈ 0
	Residual	190.42	6992.00	NA	NA
Retrieval	Model Factor	104.25	6.00	9052.57	≈ 0
	Domain Factor	6.16	1.00	3207.72	≈ 0
	Residual	13.42	6992.00	NA	NA
STS	Model Factor	10.55	6.00	149.00	1.64×10^{-178}
	Domain Factor	304.09	1.00	25761.71	≈ 0
	Residual	82.53	6992.00	NA	NA
Clustering	Model Factor	0.29	6.00	47.60	1.59×10^{-57}
	Domain Factor	32.25	1.00	32161.37	≈ 0
	Residual	7.01	6992.00	NA	NA
Summarization	Model Factor	12.98	6.00	145.31	2.90×10^{-174}
	Domain Factor	14.49	1.00	973.32	3.60×10^{-200}
	Residual	104.07	6992.00	NA	NA
Reranking	Model Factor	5.38	6.00	489.05	≈ 0
	Domain Factor	0.64	1.00	346.78	1.39×10^{-75}
	Residual	12.84	7002.00	NA	NA
Pair Classification	Model Factor	0.25	6.00	1.97	0.07
	Domain Factor	249.19	1.00	11989.92	≈ 0
	Residual	145.31	6992.00	NA	NA
Average	Model Factor	0.00	6.00	1.34	0.37
	Domain Factor	0.08	1.00	253.87	≈ 0
	Residual	0.00	6.00	NA	NA

Table 11: **Analysis of Variance (ANOVA) Results Across Tasks and Factors.** *Factor* represents the independent variables analyzed: **Model Factor** pertains to variations attributed to different models, and **Domain Factor** pertains to variations due to different domains (MTEB or FinMTEB). **Residual** refers to the unexplained variance. The **Sum of Squares**, **Degrees of Freedom**, **F-Statistic**, and **p-value** are presented for each factor within each task. Asterisks denote significance levels, with lower p-values indicating higher statistical significance. The Domain Factor consistently shows high significance across all tasks.