

DO WE NEED DOMAIN-SPECIFIC EMBEDDING MODELS? AN EMPIRICAL INVESTIGATION ON FINANCE DOMAIN AND A DOMAIN-MTEB BENCHMARK

Anonymous authors

Paper under double-blind review

ABSTRACT

Embedding models play a crucial role in representing and retrieving information across various NLP applications. Recent advancements in Large Language Models (LLMs) have further enhanced the performance of embedding models, which are trained on massive amounts of text covering almost every domain. These models are often benchmarked on general-purpose datasets like Massive Text Embedding Benchmark (MTEB), where they demonstrate superior performance. However, a critical question arises: Is the development of domain-specific embedding models necessary when general-purpose models are trained on vast corpora that already include specialized domain texts? In this paper, we empirically investigate this question, choosing the finance domain as an example. We introduce the Finance Massive Text Embedding Benchmark (FinMTEB), a counterpart to MTEB that consists of financial domain-specific text datasets. We evaluate the performance of seven state-of-the-art embedding models on FinMTEB and observe a significant performance drop compared to their performance on MTEB. To account for the possibility that this drop is driven by FinMTEB’s higher complexity, we propose four measures to quantify dataset complexity and control for this factor in our analysis. Our analysis provides compelling evidence that state-of-the-art embedding models struggle to capture domain-specific linguistic and semantic patterns. Moreover, we find that the performance of general-purpose embedding models on MTEB is not correlated with their performance on FinMTEB, indicating the need for domain-specific embedding benchmarks for domain-specific embedding models. This study sheds light on developing domain-specific embedding models in the LLM era.

1 INTRODUCTION

Embedding models, which transform text sequences into dense vector representations, play a crucial role in various natural language processing (NLP) tasks (Mikolov et al., 2013; Pennington et al., 2014; Peters et al., 2018). The quality of text embeddings directly impacts the effectiveness of information retrieval, semantic understanding, and other downstream applications. Recently, many state-of-the-art embedding models have been developed using large language models (LLMs) as the foundational model (Wang et al., 2023; Li et al., 2023; Meng et al., 2024). Since LLMs are trained on massive text corpora spanning nearly every domain, these LLM-based embedding models have demonstrated superior and robust performance in general-purpose embedding benchmarks such as Massive Text Embedding Benchmark (MTEB) (Muennighoff et al., 2022).

Given that general-purpose embedding models are becoming the backbone of NLP applications, and companies like OpenAI and Cohere are offering general-purpose embeddings that potentially serve a wide range of industry applications, a critical question arises: *Do we still need domain-specific embedding models?* The answer is not immediately clear. On the one hand, as mentioned earlier, state-of-the-art embedding models are primarily built from general-purpose LLMs that have been trained on vast text corpora covering nearly every domain. There is no strong evidence suggesting that these models cannot grasp domain-specific languages or linguistic patterns. On the other hand, while there has been limited development of domain-specific embedding models, researchers have advocated for training domain-specific LLMs (Gururangan et al., 2020) to better

054 capture domain-specific terminology and semantics. For example, domain-specific LLMs such as
055 BioMedLM (Bolton et al., 2024) for the biomedical domain, SaulLM-7B (Colombo et al., 2024) for
056 the legal domain, and BloombergGPT (Wu et al., 2023) for the finance domain are pre-trained on
057 large, domain specialized corpora.

058 To address this question, we empirically evaluate the necessity of domain-specific embedding mod-
059 els, focusing on the finance domain as our research context. We select the finance domain because
060 financial NLP is a critical area within the research community, with a wealth of established fi-
061 nancial NLP datasets (FiQA, 2018; Islam et al., 2023; Liu et al., 2024a; Ju et al., 2023; Sinha &
062 Khandait, 2021; Mukherjee et al., 2022; Malo et al., 2014). The importance of representing finan-
063 cial texts in downstream applications has also driven the development of finance-specific LLMs,
064 such as BloombergGPT (Wu et al., 2023) and InvestLM (Yang et al., 2023b). Moreover, the com-
065 plexity and specificity of the financial domain provide a unique opportunity to assess how effectively
066 general-purpose embedding models can represent specialized texts.

067 We first develop the Finance Massive Text Embedding Benchmark (FinMTEB), a finance-specific
068 counterpart to MTEB. FinMTEB consists of 64 datasets and, like MTEB, covers seven distinct tasks,
069 including classification, clustering, retrieval, pair-classification, reranking, and semantic textual sim-
070 ilarity. Unlike MTEB, all datasets in FinMTEB are based on financial text data, which feature
071 substantially longer text sequences and token lengths. Using seven state-of-the-art embedding mod-
072 els, we observe a significant performance drop on the FinMTEB compared to the general-purpose
073 MTEB. ANOVA analysis further indicates that this average performance drop is primarily driven by
074 the differences between the benchmarks, rather than model-specific factors.

075 While the performance drop on FinMTEB may seem expected given the domain shift, one concern
076 is whether the datasets in FinMTEB are inherently more complex than those in MTEB. Is the re-
077 duced performance a result of the benchmark’s complexity, or do these models lack the necessary
078 understanding of domain-specific context? If the FinMTEB datasets were of equal complexity to
079 MTEB, we might not observe the same performance gap, suggesting that dataset complexity could
080 be contributing to the performance decline.

081 To eliminate the confounding factor of dataset complexity, we propose four different measures to
082 quantify complexity: ChatGPT’s response error rate, dataset readability, information entropy, and
083 text dependency distance. Our analysis shows that even when controlling for complexity, general-
084 purpose embedding models still perform worse on domain-specific texts. Moreover, the more com-
085 plex the domain-specific data, the greater the performance drop—although this trend is less promi-
086 nent in general-purpose tasks on the MTEB. Collectively, this evidence suggests that state-of-the-art,
087 general-purpose embedding models may not fully capture the linguistic nuances and semantic pat-
088 terns unique to a particular domain.

089 Moreover, we observe that the performance of general-purpose embedding models on MTEB does
090 not correlate with their performance on FinMTEB. Models that perform exceptionally well on gen-
091 eral embedding tasks do not necessarily maintain their superiority in the financial domain. This
092 underscores the importance of evaluating embedding models within the specific context in which
093 they will be applied and emphasizes the necessity of domain-specific embedding benchmarks.

094 Our contributions in this paper are threefold:

- 095
096
097 • Our main research contribution is the empirical investigation into the necessity of domain-
098 specific embedding models. To the best of our knowledge, this is one of the first studies
099 to address the critical question of whether domain-specific embeddings are required, espe-
100 cially given the widespread adoption of general-purpose embedding models across various
101 industry applications.
- 102
103 • Our analysis on the necessity of domain-specific embedding models is based on a rigorous
104 evaluation framework. Rather than simply developing a domain benchmark and demon-
105 strating a performance drop, our analysis accounts for dataset complexity, eliminating po-
106 tential confounding factors. This allows us to conclude that the performance gap is due to
107 the models’ inability to encode domain-specific text, rather than inherent dataset complex-
ity.

- The development of the FinMTEB dataset, as a byproduct of our study, may serve as a valuable resource for researchers and practitioners interested in building financial domain-specific embedding models.

2 RELATED WORK

2.1 GENERAL-PURPOSE EMBEDDING MODELS

Embedding models like Word2Vec (Mikolov et al., 2013) and GloVe (Pennington et al., 2014) lay the groundwork by capturing word-level semantics through contextual co-occurrence. The introduction of transformer-based models such as BERT (Devlin et al., 2019) and RoBERTa (Liu, 2019) marks a significant shift by utilizing deep bidirectional encoders, enabling contextualized word embeddings. Building on these, Sentence-BERT (Reimers & Gurevych, 2019) is designed to generate semantically meaningful sentence embeddings using Siamese and triplet networks, improving performance on semantic similarity tasks. Recent advancements in LLMs have driven the development of LLM-based embedding models, such as e5-mistral-7b-instruct (Wang et al., 2023) and gte-Qwen2-1.5B-instruct (Yang et al., 2024), which have achieved state-of-the-art performance across a wide range of NLP tasks.

2.2 DOMAIN-SPECIFIC MODELS

Different domains exhibit distinct linguistic patterns and terminologies, often requiring domain-specific models or adaptations for specialized tasks. Researchers have advocated for training domain-specific models or fine-tuning general models for particular domains (Gururangan et al., 2020). For instance, domain-specific LLMs like BioMedLM (Bolton et al., 2024) for biomedical text, SaulLM-7B (Colombo et al., 2024) for legal documents, and BloombergGPT (Wu et al., 2023) for financial applications are pre-trained on large domain-specific corpora. In addition, instruction-tuned domain-specific models such as InvestLM (Yang et al., 2023b) and FinGPT (Yang et al., 2023a) are fine-tuned for specific downstream tasks in the finance domain.

While domain-specific LLMs have been widely studied and developed, domain-specific embedding models have received relatively less attention. In the biomedical domain, models like BioWordVec (Zhang et al., 2019) and BioSentVec (Chen et al., 2019) generate word and sentence embeddings tailored to biomedical texts. In finance, FinBERT (Yang et al., 2020) is pre-trained on a large corpus of financial texts to enhance text encoding capabilities. However, most domain-specific embedding models are based on smaller language models instead of state-of-the-art LLMs. Since LLMs are trained on extensive data across multiple domains with numerous parameters, it remains unclear whether general-purpose LLM-based embeddings can adequately handle specialized texts. This paper aims to address this research gap.

2.3 EMBEDDING BENCHMARKS

To comprehensively evaluate embedding models, benchmarks like the Massive Text Embedding Benchmark (MTEB) (Muennighoff et al., 2022) have been established. MTEB assesses embedding models across a wide array of tasks using numerous datasets and languages. This extensive evaluation provides insights into a model’s generalizability and effectiveness across different linguistic contexts and task types. Similarly, the BEIR benchmark (Thakur et al., 2021) focuses on the information retrieval task, encompassing 18 diverse datasets. While BEIR includes some domain-specific datasets such as FiQA (FiQA, 2018), it is not tailored for comprehensive domain analysis. The inclusion of a few specialized datasets does not fully address the unique challenges posed by domain-specific language and terminology. There are also scenario-specific RAG evaluation benchmarks like RAGeval (Zhu et al., 2024b). These benchmarks acknowledge the necessity for domain-specific evaluations, particularly highlighting the impact of accurate retrieval in specialized contexts. However, they primarily focus on retrieval tasks and often overlook other crucial embedding tasks such as semantic similarity and clustering.

2.4 DOMAIN-SPECIFIC MODEL BENCHMARKS

Numerous benchmarks tailored to specific domains have been developed with the emergence of domain-specific large language models (LLMs). For example, in the finance domain, benchmarks such as CFLUE (Zhu et al., 2024a), FinEval (Zhang et al., 2023), DocMath-Eval (Zhao et al., 2024), and FinanceBench (Islam et al., 2023) have been introduced to assess the comprehension capabilities of LLMs within financial contexts. Similarly, in the legal domain, LawBench (Fei et al., 2023) has been established to evaluate LLMs across a variety of legal tasks. Besides, MedBench (Liu et al., 2024b), MedEval (He et al., 2023), and DrBenchmark (Labrak et al., 2024) have been developed to test the proficiency in understanding and generating medical information. Most of these benchmarking papers conclude that general-purpose LLMs may fall short on domain tasks (Zhu et al., 2024a; Fei et al., 2023). The importance of domain adaptation has gradually gained attention (Ling et al., 2023). However, to our knowledge, there is little work benchmarking the embedding model’s performance on domain texts.

3 THE FINMTEB BENCHMARK

In this section, we briefly introduce the proposed Finance MTEB (FinMTEB) benchmark, which serves as the foundation for our analysis. The construction of FinMTEB closely resembles the widely used general embedding benchmark, MTEB (Muennighoff et al., 2022).

3.1 FINMTEB TASKS



Figure 1: An overview of tasks and datasets used in FinMTEB. All the dataset descriptions and examples are provided in the Appendix D.

Figure 1 provides an overview of the tasks and datasets included in FinMTEB. Similar to MTEB (Muennighoff et al., 2022), FinMTEB includes seven embedding tasks, but with datasets specifically tailored to the finance domain, as follows.

Semantic Textual Similarity (STS) involves assessing the semantic similarity between two sentences from the financial text. For this task, we employ datasets such as FinSTS (Liu et al., 2024a) and FINAL (Ju et al., 2023) from company annual reports, along with other data types such as BQ-Corpus (Chen et al., 2018) sourced from the banking corpus.

Retrieval focuses on identifying the most relevant evidence in response to a query from a financial corpus. This task utilizes some popular finance QA datasets such as FinanceBench (Islam et al., 2023), FiQA2018 (FiQA, 2018) and HPC3 (Guo et al., 2023). These datasets pair each query with relevant contextual information. Additionally, we also develop specific queries for finance terms from various sources, such as the TradeTheEvent (Zhou et al., 2021), to further enhance the finance domain evaluation.

Classification involves predicting the label of a financial text based on its text embedding. The classification task includes multiple datasets, such as financial sentiment analysis (Malo et al., 2014;

FiQA, 2018; Cortis et al., 2017; Lu et al., 2023), Fed’s monetary policy classification (Shah et al., 2023), and organization’s strategy, as well as forward-looking statement classification (Yang et al., 2023b).

Clustering is the process of grouping sentences based on their embedding similarities. We compile a diverse and comprehensive corpus that includes consumer complaints from CFPB ¹, financial papers from arXiv, company industry descriptions (Qader et al., 2018), financial events and intent detection(Gerz et al., 2021a).

Reranking includes a set of financial datasets that have the ranking of retrieved documents to user queries such as FinQA2018-Rerank (Chen et al., 2021).

Pair-Classification focuses on comparing the class label of two financial text. We use the data from AFQMC ² and finance news headline (Sinha & Khandait, 2021).

Summarization focuses on summarizing the main content of the financial text. The corpus used for this task includes earnings call transcripts (Mukherjee et al., 2022), financial news (Lu et al., 2023), and Form 10-K filings (El-Haj et al., 2022).

In summary, FinMTEB consists of a total of 64 datasets, spanning seven different tasks. The key difference between MTEB and FinMTEB is that all datasets in FinMTEB are finance-domain specific, either previously used in financial NLP research or newly developed by the authors. Semantic similarity between datasets in FinMTEB are shown in Appendix A. The detailed dataset information and descriptions are presented in the Appendix D. The main scoring metric for each task is the same as used with that of the MTEB benchmark, and the details are presented in the Appendix E.

3.2 CHARACTERISTICS OF FINMTEB

Having built the FinMTEB benchmark, we now provide an analysis to understand its characteristics.

Linguistic Pattern. Table 1 presents a comparative analysis of linguistic features such as average sentence length, token length, syllables per token, and dependency distance (Oya, 2011) across two benchmarks. The results indicate that texts in FinMTEB consistently have longer and more complex sentences than those in MTEB, with an average sentence length of 26.37 tokens compared to 18.2 tokens in MTEB. This highlights significant linguistic differences between financial and general texts. However, does this difference warrant the need for a domain-specific embedding model? We will explore this question later.

Table 1: Comparison of Text Characteristics Between FinMTEB and MTEB. The numbers represent the average scores across all samples from all 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

Semantic Diversity. We examine the inter-dataset semantic similarity. Using the all-MiniLM-L6-v2 model³, we embed 1000 samples from each dataset, compute their averages to represent the dataset embedding, and measure inter-dataset similarity using cosine similarity. As shown in Appendix A, most datasets in FinMTEB have an inter-dataset similarity score below 0.6, with a mean cosine similarity of 0.4. Despite being finance-domain specific, this highlights the diverse narratives and contexts present in the financial datasets.

3.3 THE PERFORMANCE OF STATE-OF-THE-ART EMBEDDING MODELS ON FINMTEB

General-purpose Embedding Models. We consider **seven** state-of-the-art, general-purpose embedding models in our experiments. 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

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

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

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

(Yang et al., 2024); bge-large-en-v1.5 (Xiao et al., 2023) and all-MiniLM-L12-v2 (Reimers & 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).

We evaluate the performance of these embedding models on FinMTEB tasks, with the results presented in Table 2, alongside their performance on MTEB for comparison. The results clearly demonstrate a significant performance drop on the FinMTEB benchmarks. For instance, the best-performing model on MTEB, bge-en-icl, achieves an average score of 71.67, while its performance on FinMTEB is notably lower, at 63.09.

Table 2: State-of-the-art Embedding Model Performance on MTEB and FinMTEB. The “MTEB Score” column represents performance on the MTEB benchmark (Muennighoff et al., 2022) as reported on the Hugging Face MTEB Leaderboard ⁴. The “FinMTEB Score” column shows the average performance score evaluated on the proposed FinMTEB benchmark. To ensure a fair comparison, FinMTEB uses the same evaluation metrics as MTEB. More evaluation results on other SoTA models are presented in Appendix B.

Embedding Model	Base Model	Embedding Dimensions	MTEB Score	FinMTEB Score
bge-en-icl	Mistral	4096	71.67	63.09 (↓ -8.58)
gte-Qwen2-1.5B-instruct	Qwen2	1536	67.16	59.98 (↓ -7.18)
e5-mistral-7b-instruct	Mistral	4096	66.63	64.75 (↓ -1.88)
bge-large-en-v1.5	Bert	1024	64.23	58.95 (↓ -5.28)
text-embedding-3-small	-	1536	62.26	61.36 (↓ -0.90)
instructor-base	T5Encoder	768	59.54	54.79 (↓ -4.75)
all-MiniLM-L12-v2	Bert	384	56.53	54.31 (↓ -2.22)

3.3.1 WHAT DRIVES THE PERFORMANCE GAP?

Having observed a performance discrepancy between general-purpose embedding models across the two benchmarks, we aim to investigate what drives this difference. Specifically, we consider two possible factors: (1) the **model** used (i.e., which embedding model is applied), and (2) the **domain** (i.e., whether it’s the general benchmark, MTEB, or the domain-specific benchmark, FinMTEB). Since we evaluate seven embedding models across two domains, this results in 14 unique model-domain combinations.

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

Table 3: ANOVA analysis results. The reported numbers represent the partial eta squared (effect size) for each factor (Model or Domain). Asterisks indicate statistical significance levels, with ** denoting p -value < 0.05.

	STS	Classification	Retrieval	Reranking	Clustering	Pair-classification	Summarization	Average
Model	0.79**	0.02**	0.89**	0.30**	0.04**	0.00	0.11**	0.00
Domain	0.11**	0.23**	0.31**	0.05**	0.82**	0.63**	0.12**	1.00**

We present the ANOVA analysis results in Appendix C. First, the results indicate that the choice of embedding model (Model factor) significantly impacts performance in most tasks, such as STS (0.79), Retrieval (0.89), and Reranking (0.30), with the exception of Pair-classification (0.00), where model choice has no significant impact. Second, the Domain factor also shows significant effects across all embedding tasks. Interestingly, the average scores reveal that, from an overall perspective, the Model factor has little impact on performance, with an effect size of 0.00 and an insignificant p -value. This suggests that while individual models may excel at specific tasks, their performance discrepancies balance out when averaged. However, the Domain factor (1.00) demonstrates a much

⁴<https://huggingface.co/spaces/mteb/leaderboard>

324 more prominent influence, underscoring the necessity for domain-specific models or fine-tuning
325 when addressing specialized tasks like those in finance.

326 **Research Question.** While the performance drop on FinMTEB and the subsequent ANOVA analy-
327 sis suggests that domain-specific embedding tasks may pose greater challenges for general-purpose
328 embedding models, does this necessarily indicate a need for domain-specific models? Not neces-
329 sarily. The difference in datasets between FinMTEB and MTEB could contribute to the observed
330 performance drop. For instance, FinMTEB datasets might be inherently more difficult or linguisti-
331 cally complex compared to those in MTEB as illustrated in Table 1. If both benchmarks contained
332 datasets of equivalent complexity, general-purpose models might even perform better on FinMTEB
333 tasks. Therefore, the performance drop does not necessarily imply that the models fail to under-
334 stand domain-specific language or concepts. To draw meaningful conclusions about the necessity
335 of domain-specific models, we must first control for differences in dataset difficulty. In the next
336 section, we will analyze model performance while considering these inherent differences between
337 the FinMTEB and MTEB datasets.

338 339 4 PERFORMANCE ANALYSIS AFTER CONTROLLING FOR DATASET 340 COMPLEXITY 341

342
343 To answer the above research question, we conduct a detailed analysis of the embedding models’
344 performance, while accounting for dataset complexity.

345 346 4.1 QUANTIFYING DATASET COMPLEXITY 347

348 We propose four different measures to quantify a dataset’s complexity.

349 **ChatGPT Error Rate.** The first measure quantifies how challenging it is for ChatGPT to answer a
350 dataset’s questions. Specifically, for each example in the dataset across different tasks, we reformat
351 the example into a question-and-answer pair, as shown in Appendix G, and prompt GPT-4o mini
352 (OpenAI, 2024a). The rationale is that if ChatGPT fails to answer a question correctly, it indicates
353 the difficulty level of the question. Additionally, since state-of-the-art LLM-based embedding mod-
354 els present each query with an instruction in a question-answer format, we use the ChatGPT error
355 rate as an indicator of dataset complexity.

356 **Information Theory.** We borrow the concept of information entropy from information theory to
357 measure the complexity of a text sequence. Information entropy is calculated as the average amount
358 of information produced by a stochastic source of data. Higher Information Entropy indicates text
359 that contains more information or is less predictable, potentially implying greater complexity.

360 **Readability.** We also use readability to measure dataset complexity, specifically applying the Gun-
361 ning Fog Index (Gunning, 1952), which factors in sentence length and the number of complex words.
362 The index estimates the years of formal education required to understand a text on the first reading.
363 A higher Gunning Fog Index score indicates more complex sentences.

364 **Mean Dependency Distance.** Finally, we measure linguistic complexity using the dependency dis-
365 tance between two syntactically related words in a sentence (Oya, 2011). A longer dependency
366 distance indicates that more context is needed for comprehension, reflecting greater sentence com-
367 plexity.

368
369 For all of these four complexity measures, a higher score indicates higher dataset complexity. Details
370 on the measures and their calculations are provided in Appendix F.

371 372 4.2 SUBGROUP ANALYSIS: EMBEDDING PERFORMANCE VS. DATASET COMPLEXITY 373

374 We conduct a subgroup analysis to examine the impact of the domain on embedding model per-
375 formance. First, we calculate dataset complexity using one of the four complexity measures and
376 categorize the datasets into three subgroups: low, medium, and high complexity. This ensures that
377 METB and FinMTEB datasets within each subgroup have the same level of complexity. We then
calculate the average performance score of seven LLM-based embedding models across datasets

378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431

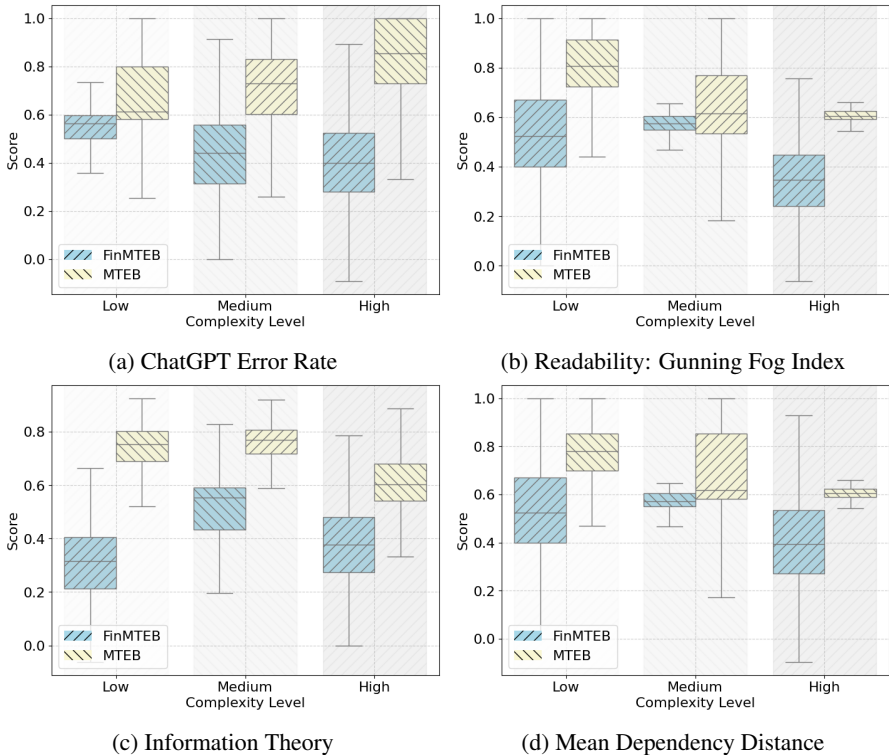


Figure 2: Subgroup analysis results. The x-axis represents the three dataset complexity levels: Low, Medium, and High. The y-axis reports the average score for each dataset across all benchmark tasks.

within each subgroup. The results of this analysis appear in Figure 4.2. The subgroup analysis reveals two key findings.

First, embedding models perform substantially worse on FinMTEB datasets compared to MTEB datasets, even after accounting for dataset complexity. The performance of embedding models on FinMTEB datasets is consistently lower than on MTEB datasets within the same group. This performance gap persists regardless of the complexity measure used. Given that datasets in the same subgroup have similar levels of complexity (e.g., comparable ChatGPT error rates), the lower performance on FinMTEB tasks suggests that state-of-the-art LLM-based embedding models struggle with encoding domain-specific terminologies and semantic patterns.

Second, embedding models perform worst on FinMTEB datasets with the highest complexity levels. For example, using readability as a complexity measure, the average performance of embedding models on high-complexity FinMTEB datasets is around 0.3, significantly lower than their performance on low-complexity datasets (around 0.5) and medium-complexity datasets (around 0.6). This further highlights the challenge embedding models face in capturing complex, domain-specific language and semantics.

4.3 REGRESSION ANALYSIS: EMBEDDING PERFORMANCE VS. DATASET COMPLEXITY

Table 4: Regression analysis results. The numbers represent the estimated coefficient values, with standard errors in parentheses. ** denotes significance at the $p < 0.05$ level.

Domain	ChatGPT Error Rate		Readability		Information Theory		Mean Dependency Distance	
	MTEB	FinMTEB	MTEB	FinMTEB	MTEB	FinMTEB	MTEB	FinMTEB
Intercept	0.6619** (0.003)	0.6097** (0.003)	0.8739** (0.002)	0.6976** (0.005)	0.5905** (0.002)	0.4766** (0.002)	1.2341** (0.007)	0.9357** (0.008)
Complexity	-0.2193** (0.009)	-0.4637** (0.010)	-0.0278** (0.000)	-0.0202** (0.000)	0.0168** (0.001)	-0.0092** (0.001)	-0.2703** (0.003)	-0.1810** (0.003)

Furthermore, we conduct a regression analysis to examine the relationship between dataset complexity and embedding model performance. Specifically, for each dataset complexity measure, we run two ordinary least squares (OLS) regression models—one for the MTEB datasets and one for the FinMTEB datasets. We normalize the variables Performance and Complexity to a range between 0 and 1 using min-max normalization. The regression specification is as follows:

$$\text{Performance} = \text{Intercept} + \beta \times \text{Complexity} + \epsilon \quad (1)$$

where β is the coefficients of the covariate (i.e., dataset complexity), and ϵ is the error term.

The regression results are presented in Table 4 and are largely consistent with the findings from the subgroup analysis. First, there is a negative relationship between dataset complexity and embedding model performance, indicating that models significantly struggle with domain-specific texts exhibiting higher linguistic complexity. Second, the intercept for the MTEB datasets is consistently higher than that for FinMTEB. Given that the Complexity variable is normalized between 0 and 1, these results suggest a significant performance gap between embedding models on MTEB and FinMTEB, even when controlling for the same level of dataset complexity.

Overall, both the subgroup and regression analyses demonstrate that the performance drop reported in Table 2 is not driven by differences in dataset complexity between MTEB and FinMTEB benchmarks. Rather, it suggests that state-of-the-art, general-purpose embedding models may not fully capture the linguistic nuances and semantic patterns specific to a given domain.

5 DOMAIN-SPECIFIC EMBEDDING BENCHMARK IS NEEDED

Another key consideration when discussing domain-specific embeddings is whether we need a domain-specific embedding benchmark. While it may seem intuitive to say yes, there is little empirical evidence supporting this assumption. To explore this question, 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 5, 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. This suggests the necessity of developing domain-specific embedding benchmarks for evaluating domain-specific embeddings. Therefore, the development of FinMTEB constitutes a significant contribution to benchmarking embedding models specifically for the financial domain.

Table 5: 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.

	STS	Classification	Retrieval	Reranking	Clustering	Pair-classification	Summarization
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

6 CONCLUSION

In this study, we empirically investigate a seemingly intuitive yet practically important question: do we need domain-specific embedding models? To rigorously address this, we use the finance domain as an example and develop the FinMTEB benchmark, which comprises a large variety of domain-specific (i.e., finance) embedding tasks. Additionally, we propose four methods to quantify dataset complexity. Our comparative analysis reveals that state-of-the-art LLM-based embedding models exhibit a substantial performance gap between general (MTEB) and domain-specific (FinMTEB) benchmarks. More importantly, this gap persists even when accounting for dataset complexity. The empirical results provide strong evidence that, despite being trained on vast amounts of data that likely include various domains, LLM-based embedding models still fall short in capturing domain-specific linguistic and semantic patterns. Given the widespread use of embedding models in information retrieval and semantic search applications, this highlights the need for further adaptation of these models to specific domains in order to improve their utility. Moreover, the development of the

486 FinMTEB benchmark can serve as a valuable resource for researchers and practitioners interested
487 in financial-specific embedding models.

488 While this study presents compelling evidence for the necessity of domain-specific embedding mod-
489 els, the challenge of how to train these models remains. Should we adapt domain-specific embed-
490 dings from a domain-specific LLM, or should we develop domain-specific datasets and fine-tune
491 a general-purpose LLM? We hope to see more research in this direction to further advance AI’s
492 capabilities in handling domain-specific tasks effectively.

493 REFERENCES

494 Elliot Bolton, Abhinav Venigalla, Michihiro Yasunaga, David Hall, Betty Xiong, Tony Lee, Roxana
495 Daneshjou, Jonathan Frankle, Percy Liang, Michael Carbin, et al. Biomedlm: A 2.7 b parameter
496 language model trained on biomedical text. *arXiv preprint arXiv:2403.18421*, 2024.

497 CCKS. Ccks2022: Few-shot event extraction for the financial sector, 2022. [https://www.
498 biendata.xyz/competition/ccks2022_eventext/](https://www.biendata.xyz/competition/ccks2022_eventext/).

499 CFPB. Consumer finance complaints, 2024. [https://huggingface.co/datasets/
500 CFPB/consumer-finance-complaints](https://huggingface.co/datasets/CFPB/consumer-finance-complaints).

501 Jing Chen, Qingcai Chen, Xin Liu, Haijun Yang, Daohe Lu, and Buzhou Tang. The BQ corpus:
502 A large-scale domain-specific Chinese corpus for sentence semantic equivalence identification.
503 In Ellen Riloff, David Chiang, Julia Hockenmaier, and Jun’ichi Tsujii (eds.), *Proceedings of
504 the 2018 Conference on Empirical Methods in Natural Language Processing*, pp. 4946–4951,
505 Brussels, Belgium, October-November 2018. Association for Computational Linguistics. doi:
506 10.18653/v1/D18-1536.

507 Qingyu Chen, Yifan Peng, and Zhiyong Lu. Biosentvec: creating sentence embeddings for biomed-
508 ical texts. In *2019 IEEE International Conference on Healthcare Informatics (ICHI)*, pp. 1–5.
509 IEEE, 2019.

510 Wei Chen, Qiushi Wang, Zefei Long, Xianyin Zhang, Zhongtian Lu, Bingxuan Li, Siyuan Wang,
511 Jiarong Xu, Xiang Bai, Xuanjing Huang, and Zhongyu Wei. Disc-finllm: A chinese financial
512 large language model based on multiple experts fine-tuning. *arXiv preprint arXiv:2310.15205*,
513 2023.

514 Zhiyu Chen, Wenhui Chen, Charese Smiley, Sameena Shah, Iana Borova, Dylan Langdon, Reema
515 Moussa, Matt Beane, Ting-Hao Huang, Bryan Routledge, and William Yang Wang. FinQA: A
516 dataset of numerical reasoning over financial data. In Marie-Francine Moens, Xuanjing Huang,
517 Lucia Specia, and Scott Wen-tau Yih (eds.), *Proceedings of the 2021 Conference on Empirical
518 Methods in Natural Language Processing*, pp. 3697–3711, Online and Punta Cana, Dominican
519 Republic, November 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.
520 emnlp-main.300.

521 Pierre Colombo, Telmo Pessoa Pires, Malik Boudiaf, Dominic Culver, Rui Melo, Caio Corro, An-
522 dre FT Martins, Fabrizio Esposito, Vera Lúcia Raposo, Sofia Morgado, et al. Saullm-7b: A
523 pioneering large language model for law. *arXiv preprint arXiv:2403.03883*, 2024.

524 Keith Cortis, André Freitas, Tobias Daudert, Manuela Huerlimann, Manel Zarrouk, Siegfried
525 Handschuh, and Brian Davis. SemEval-2017 task 5: Fine-grained sentiment analysis on fi-
526 nancial microblogs and news. In Steven Bethard, Marine Carpuat, Marianna Apidianaki,
527 Saif M. Mohammad, Daniel Cer, and David Jurgens (eds.), *Proceedings of the 11th Interna-
528 tional Workshop on Semantic Evaluation (SemEval-2017)*, pp. 519–535, Vancouver, Canada,
529 August 2017. Association for Computational Linguistics. doi: 10.18653/v1/S17-2089. URL
530 <https://aclanthology.org/S17-2089>.

531 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of
532 deep bidirectional transformers for language understanding. In Jill Burstein, Christy Doran, and
533 Thamar Solorio (eds.), *Proceedings of the 2019 Conference of the North American Chapter of
534 the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long
535 and Short Papers)*, pp. 4171–4186, Minneapolis, Minnesota, June 2019. Association for Compu-
536 tational Linguistics. doi: 10.18653/v1/N19-1423.

- 540 Mahmoud El-Haj, Nadhem Zmandar, Paul Rayson, Ahmed AbuRa'ed, Marina Litvak, Nikiforos
541 Pittaras, George Giannakopoulos, Aris Kosmopoulos, Blanca Carbajo-Coronado, and Antonio
542 Moreno-Sandoval. The financial narrative summarisation shared task (FNS 2022). In Mahmoud
543 El-Haj, Paul Rayson, and Nadhem Zmandar (eds.), *Proceedings of the 4th Financial Narrative
544 Processing Workshop @LREC2022*, pp. 43–52, Marseille, France, June 2022. European Lan-
545 guage Resources Association.
- 546 Zhiwei Fei, Xiaoyu Shen, Dawei Zhu, Fengzhe Zhou, Zhuo Han, Songyang Zhang, Kai Chen,
547 Zongwen Shen, and Jidong Ge. Lawbench: Benchmarking legal knowledge of large language
548 models. *arXiv preprint arXiv:2309.16289*, 2023.
- 549
550 FiQA. Financial question answering., 2018. <https://sites.google.com/view/fiqa>.
- 551 Daniela Gerz, Pei-Hao Su, Razvan Kusztos, Avishek Mondal, Michał Lis, Eshan Singhal, Nikola
552 Mrkšić, Tsung-Hsien Wen, and Ivan Vulić. Multilingual and cross-lingual intent detection from
553 spoken data. In Marie-Francine Moens, Xuanjing Huang, Lucia Specia, and Scott Wen-tau Yih
554 (eds.), *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Process-
555 ing*, pp. 7468–7475, Online and Punta Cana, Dominican Republic, November 2021a. Association
556 for Computational Linguistics. doi: 10.18653/v1/2021.emnlp-main.591.
- 557
558 Daniela Gerz, Pei-Hao Su, Razvan Kusztos, Avishek Mondal, Michał Lis, Eshan Singhal, Nikola
559 Mrkšić, Tsung-Hsien Wen, and Ivan Vulić. Multilingual and cross-lingual intent detection from
560 spoken data. *arXiv preprint arXiv:2104.08524*, 2021b.
- 561 Robert Gunning. The technique of clear writing, 1952.
- 562
563 Biyang Guo, Xin Zhang, Ziyuan Wang, Minqi Jiang, Jinran Nie, Yuxuan Ding, Jianwei Yue, and Yu-
564 peng Wu. How close is chatgpt to human experts? comparison corpus, evaluation, and detection.
565 *arXiv preprint arXiv:2301.07597*, 2023.
- 566
567 Suchin Gururangan, Ana Marasović, Swabha Swayamdipta, Kyle Lo, Iz Beltagy, Doug Downey,
568 and Noah A. Smith. Don’t stop pretraining: Adapt language models to domains and tasks. In Dan
569 Jurafsky, Joyce Chai, Natalie Schluter, and Joel Tetreault (eds.), *Proceedings of the 58th Annual
570 Meeting of the Association for Computational Linguistics*, pp. 8342–8360, Online, July 2020.
571 Association for Computational Linguistics. doi: 10.18653/v1/2020.acl-main.740.
- 572
573 Zexue He, Yu Wang, An Yan, Yao Liu, Eric Chang, Amilcare Gentili, Julian McAuley, and Chun-
574 Nan Hsu. MedEval: A multi-level, multi-task, and multi-domain medical benchmark for language
575 model evaluation. In Houda Bouamor, Juan Pino, and Kalika Bali (eds.), *Proceedings of the 2023
576 Conference on Empirical Methods in Natural Language Processing*, pp. 8725–8744, Singapore,
577 December 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-main.
540.
- 578
579 Pranab Islam, Anand Kannappan, Douwe Kiela, Rebecca Qian, Nino Scherrer, and Bertie Vid-
580 gen. Financebench: A new benchmark for financial question answering. *arXiv preprint
581 arXiv:2311.11944*, 2023.
- 582
583 Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot,
584 Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al.
585 Mistral 7b. *arXiv preprint arXiv:2310.06825*, 2023.
- 586
587 Jia-Huei Ju, Yu-Shiang Huang, Cheng-Wei Lin, Che Lin, and Chuan-Ju Wang. A compare-and-
588 contrast multistage pipeline for uncovering financial signals in financial reports. In Anna Rogers,
589 Jordan Boyd-Graber, and Naoaki Okazaki (eds.), *Proceedings of the 61st Annual Meeting of the
590 Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 14307–14321, Toronto,
591 Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.acl-long.
592 800.
- 593
594 Yanis Labrak, Adrien Bazoge, Oumaima El Khettari, Mickaël Rouvier, Natalia Grabar, Beatrice
595 Daille, Solen Quiniou, Emmanuel Morin, Pierre-Antoine Gourraud, Richard Dufour, et al. Dr-
596 benchmark: A large language understanding evaluation benchmark for french biomedical domain.
597 *arXiv preprint arXiv:2402.13432*, 2024.

- 594 Yinyu Lan, Yanru Wu, Wang Xu, Weiqiang Feng, and Youhao Zhang. Chinese fine-grained financial
595 sentiment analysis with large language models. *arXiv preprint arXiv:2306.14096*, 2023.
596
- 597 Chankyu Lee, Rajarshi Roy, Mengyao Xu, Jonathan Raiman, Mohammad Shoeybi, Bryan Catan-
598 zaro, and Wei Ping. Nv-embed: Improved techniques for training llms as generalist embedding
599 models. *arXiv preprint arXiv:2405.17428*, 2024.
- 600 Xiang Li, Zhenyu Li, Chen Shi, Yong Xu, Qing Du, Mingkui Tan, Jun Huang, and Wei Lin. Al-
601 phafin: Benchmarking financial analysis with retrieval-augmented stock-chain framework. *arXiv*
602 *preprint arXiv:2403.12582*, 2024.
603
- 604 Xianming Li and Jing Li. Angle-optimized text embeddings. *arXiv preprint arXiv:2309.12871*,
605 2023.
606
- 607 Xianming Li and Jing Li. Bellm: Backward dependency enhanced large language model for sen-
608 tence embeddings. In *Proceedings of the 2024 Conference of the North American Chapter of*
609 *the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long*
610 *Papers)*, pp. 792–804, 2024.
- 611 Zehan Li, Xin Zhang, Yanzhao Zhang, Dingkun Long, Pengjun Xie, and Meishan Zhang. Towards
612 general text embeddings with multi-stage contrastive learning. *arXiv preprint arXiv:2308.03281*,
613 2023.
- 614 Chen Ling, Xujiang Zhao, Jiaying Lu, Chengyuan Deng, Can Zheng, Junxiang Wang, Tanmoy
615 Chowdhury, Yun Li, Hejie Cui, Xuchao Zhang, et al. Domain specialization as the key to make
616 large language models disruptive: A comprehensive survey. *arXiv preprint arXiv:2305.18703*,
617 2023.
618
- 619 Jiaxin Liu, Yi Yang, and Kar Yan Tam. Beyond surface similarity: Detecting subtle semantic shifts
620 in financial narratives. In Kevin Duh, Helena Gomez, and Steven Bethard (eds.), *Findings of the*
621 *Association for Computational Linguistics: NAACL 2024*, pp. 2641–2652, Mexico City, Mexico,
622 June 2024a. Association for Computational Linguistics. doi: 10.18653/v1/2024.findings-naacl.
623 168.
- 624 Mianxin Liu, Jinru Ding, Jie Xu, Weiguo Hu, Xiaoyang Li, Lifeng Zhu, Zhian Bai, Xiaoming
625 Shi, Benyou Wang, Haitao Song, et al. Medbench: A comprehensive, standardized, and reli-
626 able benchmarking system for evaluating chinese medical large language models. *arXiv preprint*
627 *arXiv:2407.10990*, 2024b.
- 628 Shuaiqi Liu, Jiannong Cao, Ruosong Yang, and Zhiyuan Wen. Long text and multi-table sum-
629 marization: Dataset and method. In Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang (eds.),
630 *Findings of the Association for Computational Linguistics: EMNLP 2022*, pp. 1995–2010, Abu
631 Dhabi, United Arab Emirates, December 2022. Association for Computational Linguistics. doi:
632 10.18653/v1/2022.findings-emnlp.145.
633
- 634 Yinhan Liu. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint*
635 *arXiv:1907.11692*, 2019.
- 636 Dakuan Lu, Hengkui Wu, Jiaqing Liang, Yipei Xu, Qianyu He, Yipeng Geng, Mengkun Han, Yingsi
637 Xin, and Yanghua Xiao. Bbt-fin: Comprehensive construction of chinese financial domain pre-
638 trained language model, corpus and benchmark. *arXiv preprint arXiv:2302.09432*, 2023.
639
- 640 Pekka Malo, Ankur Sinha, Pekka Korhonen, Jyrki Wallenius, and Pyry Takala. Good debt or bad
641 debt: Detecting semantic orientations in economic texts. *Journal of the Association for Informa-*
642 *tion Science and Technology*, 65(4):782–796, 2014.
- 643 Rui Meng, Ye Liu, Shafiq Rayhan Joty, Caiming Xiong, Yingbo Zhou, and Semih Yavuz. Sfr-
644 embedding-2: Advanced text embedding with multi-stage training, 2024. URL [https://](https://huggingface.co/Salesforce/SFR-Embedding-2_R)
645 huggingface.co/Salesforce/SFR-Embedding-2_R.
646
- 647 Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word represen-
tations in vector space. *arXiv preprint arXiv:1301.3781*, 2013.

- 648 Niklas Muennighoff, Nouamane Tazi, Loïc Magne, and Nils Reimers. Mteb: Massive text embed-
649 ding benchmark. *arXiv preprint arXiv:2210.07316*, 2022.
- 650
- 651 Rajdeep Mukherjee, Abhinav Bohra, Akash Banerjee, Soumya Sharma, Manjunath Hegde, Afreen
652 Shaikh, Shivani Shrivastava, Koustuv Dasgupta, Niloy Ganguly, Saptarshi Ghosh, et al. Ectsum:
653 A new benchmark dataset for bullet point summarization of long earnings call transcripts. *arXiv*
654 *preprint arXiv:2210.12467*, 2022.
- 655 Qiong Nan, Juan Cao, Yongchun Zhu, Yanyan Wang, and Jintao Li. Mdfend: Multi-domain fake
656 news detection. In *Proceedings of the 30th ACM International Conference on Information &*
657 *Knowledge Management*, pp. 3343–3347, 2021.
- 658 OpenAI. Openai (august 24 version), 2024a. <https://api.openai.com/v1/chat>.
- 659
- 660 OpenAI. Openai (august 24 version), 2024b. <https://api.openai.com/v1/embeddings>.
- 661 Masanori Oya. Syntactic dependency distance as sentence complexity measure. In *Proceedings of*
662 *the 16th International Conference of Pan-Pacific Association of Applied Linguistics*, volume 1,
663 2011.
- 664
- 665 Jeffrey Pennington, Richard Socher, and Christopher Manning. GloVe: Global vectors for word
666 representation. In Alessandro Moschitti, Bo Pang, and Walter Daelemans (eds.), *Proceedings*
667 *of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp.
668 1532–1543, Doha, Qatar, October 2014. Association for Computational Linguistics. doi: 10.
669 3115/v1/D14-1162. URL <https://aclanthology.org/D14-1162>.
- 670 Matthew E. Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee,
671 and Luke Zettlemoyer. Deep contextualized word representations. In Marilyn Walker, Heng Ji,
672 and Amanda Stent (eds.), *Proceedings of the 2018 Conference of the North American Chapter of*
673 *the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long*
674 *Papers)*, pp. 2227–2237, New Orleans, Louisiana, June 2018. Association for Computational
675 Linguistics. doi: 10.18653/v1/N18-1202.
- 676 Raheel Qader, Khoder Jneid, François Portet, and Cyril Labbé. Generation of company descriptions
677 using concept-to-text and text-to-text deep models: dataset collection and systems evaluation. In
678 Emiel Kraemer, Albert Gatt, and Martijn Goudbeek (eds.), *Proceedings of the 11th International*
679 *Conference on Natural Language Generation*, pp. 254–263, Tilburg University, The Netherlands,
680 November 2018. Association for Computational Linguistics. doi: 10.18653/v1/W18-6532. URL
681 <https://aclanthology.org/W18-6532>.
- 682 Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi
683 Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text
684 transformer. *Journal of machine learning research*, 21(140):1–67, 2020.
- 685 Nils Reimers and Iryna Gurevych. Sentence-bert: Sentence embeddings using siamese bert-
686 networks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language*
687 *Processing*. Association for Computational Linguistics, 11 2019. URL [https://arxiv.](https://arxiv.org/abs/1908.10084)
688 [org/abs/1908.10084](https://arxiv.org/abs/1908.10084).
- 689
- 690 Andrew Rosenberg and Julia Hirschberg. V-measure: A conditional entropy-based external cluster
691 evaluation measure. In Jason Eisner (ed.), *Proceedings of the 2007 Joint Conference on Empir-*
692 *ical Methods in Natural Language Processing and Computational Natural Language Learning*
693 *(EMNLP-CoNLL)*, pp. 410–420, Prague, Czech Republic, June 2007. Association for Computa-
694 tional Linguistics.
- 695 Agam Shah, Suvan Paturi, and Sudheer Chava. Trillion dollar words: A new financial dataset, task &
696 market analysis. In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (eds.), *Proceedings*
697 *of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long*
698 *Papers)*, pp. 6664–6679, Toronto, Canada, July 2023. Association for Computational Linguistics.
699 doi: 10.18653/v1/2023.acl-long.368.
- 700 Ankur Sinha and Tanmay Khandait. Impact of news on the commodity market: Dataset and results.
701 In *Advances in Information and Communication: Proceedings of the 2021 Future of Information*
and Communication Conference (FICC), Volume 2, pp. 589–601. Springer, 2021.

- 702 Jacob Mitchell Springer, Suhas Kotha, Daniel Fried, Graham Neubig, and Aditi Raghunathan. Rep-
703 etition improves language model embeddings. *arXiv preprint arXiv:2402.15449*, 2024.
704
- 705 Hongjin Su, Weijia Shi, Jungo Kasai, Yizhong Wang, Yushi Hu, Mari Ostendorf, Wen-tau Yih,
706 Noah A Smith, Luke Zettlemoyer, and Tao Yu. One embedder, any task: Instruction-finetuned
707 text embeddings. *arXiv preprint arXiv:2212.09741*, 2022.
- 708 Maosong Sun, Jingyang Li, Zhipeng Guo, Yu Zhao, Yabin Zheng, Xiance Si, and Zhiyuan Liu.
709 Thuctc: An efficient chinese text classifier, 2016. <http://thuctc.thunlp.org/>.
710
- 711 Nandan Thakur, Nils Reimers, Andreas Rücklé, Abhishek Srivastava, and Iryna Gurevych. BEIR: A
712 heterogeneous benchmark for zero-shot evaluation of information retrieval models. In *Thirty-fifth
713 Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round
714 2)*, 2021.
- 715 Liang Wang, Nan Yang, Xiaolong Huang, Linjun Yang, Rangan Majumder, and Furu Wei. Improv-
716 ing text embeddings with large language models. *arXiv preprint arXiv:2401.00368*, 2023.
717
- 718 Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A Smith, Daniel Khachabi, and
719 Hannaneh Hajishirzi. Self-instruct: Aligning language models with self-generated instructions.
720 *arXiv preprint arXiv:2212.10560*, 2022.
- 721 Alex Watson, Yev Meyer, Maarten Van Segbroeck, Matthew Grossman, Sami Torbey, Piotr
722 Mlocek, and Johnny Greco. Synthetic-PII-Financial-Documents-North-America: A synthetic
723 dataset for training language models to label and detect pii in domain specific formats, June
724 2024. URL [https://huggingface.co/datasets/gretelai/synthetic_pii_](https://huggingface.co/datasets/gretelai/synthetic_pii_finance_multilingual)
725 [finance_multilingual](https://huggingface.co/datasets/gretelai/synthetic_pii_finance_multilingual).
- 726 Shijie Wu, Ozan Irsoy, Steven Lu, Vadim Dabravolski, Mark Dredze, Sebastian Gehrmann, Prab-
727 hanjan Kambadur, David Rosenberg, and Gideon Mann. Bloomberggpt: A large language model
728 for finance. *arXiv preprint arXiv:2303.17564*, 2023.
729
- 730 Shitao Xiao, Zheng Liu, Peitian Zhang, and Niklas Muennighoff. C-pack: Packaged resources to
731 advance general chinese embedding. *arXiv preprint arXiv:2309.07597*, 2023.
732
- 733 Ziyue Xu, Peilin Zhou, Xinyu Shi, Jiageng Wu, Yikang Jiang, Bin Ke, and Jie Yang. Fintruthqa: A
734 benchmark dataset for evaluating the quality of financial information disclosure. *arXiv preprint
735 arXiv:2406.12009*, 2024.
- 736 An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li,
737 Chengyuan Li, Dayiheng Liu, Fei Huang, et al. Qwen2 technical report. *arXiv preprint
738 arXiv:2407.10671*, 2024.
- 739 Hongyang Yang, Xiao-Yang Liu, and Christina Dan Wang. Fingpt: Open-source financial large
740 language models. *arXiv preprint arXiv:2306.06031*, 2023a.
741
- 742 Yi Yang, Mark Christopher Siy Uy, and Allen Huang. Finbert: A pretrained language model for
743 financial communications. *arXiv preprint arXiv:2006.08097*, 2020.
- 744 Yi Yang, Yixuan Tang, and Kar Yan Tam. Investlm: A large language model for investment using
745 financial domain instruction tuning. *arXiv preprint arXiv:2309.13064*, 2023b.
746
- 747 Liwen Zhang, Weige Cai, Zhaowei Liu, Zhi Yang, Wei Dai, Yujie Liao, Qianru Qin, Yifei Li, Xingyu
748 Liu, Zhiqiang Liu, et al. Fineval: A chinese financial domain knowledge evaluation benchmark
749 for large language models. *arXiv preprint arXiv:2308.09975*, 2023.
- 750 Yijia Zhang, Qingyu Chen, Zhihao Yang, Hongfei Lin, and Zhiyong Lu. Biowordvec, improving
751 biomedical word embeddings with subword information and mesh. *Scientific data*, 6(1):52, 2019.
752
- 753 Yilun Zhao, Yitao Long, Hongjun Liu, Ryo Kamoi, Linyong Nan, Lyuhao Chen, Yixin Liu, Xiangru
754 Tang, Rui Zhang, and Arman Cohan. Docmath-eval: Evaluating math reasoning capabilities of
755 llms in understanding financial documents. In *Proceedings of the 62nd Annual Meeting of the
Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 16103–16120, 2024.

756 Zhihan Zhou, Liqian Ma, and Han Liu. Trade the event: Corporate events detection for news-
 757 based event-driven trading. In Chengqing Zong, Fei Xia, Wenjie Li, and Roberto Navigli
 758 (eds.), *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pp. 2114–
 759 2124, Online, August 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.
 760 findings-acl.186.

761 Fengbin Zhu, Wenqiang Lei, Youcheng Huang, Chao Wang, Shuo Zhang, Jiancheng Lv, Fuli Feng,
 762 and Tat-Seng Chua. Tat-qa: A question answering benchmark on a hybrid of tabular and textual
 763 content in finance. *arXiv preprint arXiv:2105.07624*, 2021.

764 Jie Zhu, Junhui Li, Yalong Wen, and Lifan Guo. Benchmarking large language models on cflue—a
 765 chinese financial language understanding evaluation dataset. *arXiv preprint arXiv:2405.10542*,
 766 2024a.

767 Kunlun Zhu, Yifan Luo, Dingling Xu, Ruobing Wang, Shi Yu, Shuo Wang, Yukun Yan, Zhenghao
 768 Liu, Xu Han, Zhiyuan Liu, et al. Rageval: Scenario specific rag evaluation dataset generation
 769 framework. *arXiv preprint arXiv:2408.01262*, 2024b.

772 A DATASET SEMANTIC SIMILARITY

773 Figure 3 presents the semantic similarity across all datasets in the Finance MTEB benchmark. The
 774 semantic similarity is calculated by cosine similarity.

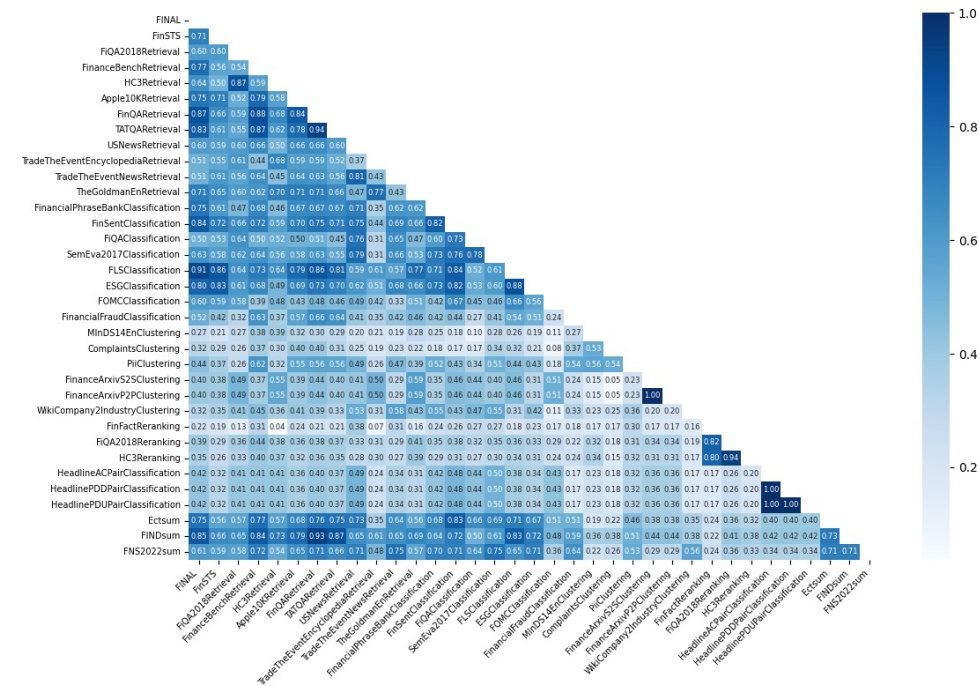


Figure 3: Semantic similarity across all the datasets in FinMTEB benchmark.

B ADDITIONAL STATE-OF-THE-ART EMBEDDING MODEL PERFORMANCE

C ANOVA DATA

Table 7 illustrates the full results of ANOVA analysis.

Table 6: Performance metrics of more state-of-the-art (SoTA) models on FinMTEB and their testing times (batch size = 512)

Model	STS	Retr.	Rerank.	Class.	Summ.	PairClass.	Clus.	Avg.	Duration (hr)
Echo Embedding (Springer et al., 2024)	0.4380	0.6443	0.9765	0.6525	0.4722	0.6261	0.5776	0.6267	12.00
Angle-BERT (Li & Li, 2023)	0.3080	0.5730	0.9650	0.6439	0.5049	0.6891	0.5774	0.6088	8.00
NV-Embed v2 (Lee et al., 2024)	0.3739	0.7061	0.9822	0.6393	0.5103	0.6043	0.6096	0.5739	5.60
BeLLM (Li & Li, 2024)	0.3919	0.0169	0.5661	0.1168	0.3906	0.5682	0.0685	0.3027	11.98

Table 7: **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.

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

D DATASETS

The detailed description of each dataset used in this work is listed in the Table tables 8 to 14.

Table 8: Summary of STS Datasets

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.

⁵<https://tianchi.aliyun.com/dataset/106411>

864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917

Table 9: Summary of Retrieval Datasets

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 10: Summary of Classification 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 11: Summary of Clustering 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 12: Summary of Summarization 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)	Chinese	A financial summarization dataset.
Fin-Eva (Abstract)	Chinese	A financial summarization dataset.

⁶<https://lighthouse.ai/blog/rag-benchmark-finance-apple-10K-2022/>⁷<https://www.kaggle.com/datasets/jeet2016/us-financial-news-articles>⁸https://github.com/alipay/financial_evaluation_dataset/tree/main⁹<https://github.com/smoothnlp/SmoothNLP>¹⁰https://github.com/alipay/financial_evaluation_dataset/tree/main¹¹<https://github.com/amitkedia007/Financial-Fraud-Detection-Using-LLMs/tree/main>¹²<https://github.com/open-compass/OpenFinData?tab=readme-ov-file>

Table 13: Summary of Reranking Datasets

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

Table 14: Summary of PairClassification Datasets

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

E MAIN METRIC

Semantic Textual Similarity (STS): The main metric used to measure performance in this task is Spearman’s rank correlation of predicted cosine similarity scores with the true similarity score.

Classification: The main metric for evaluating is accuracy, ensuring that the model’s assessment is based on different types of financial texts and frameworks.

Clustering: The main evaluation metric for this task is the V-measure (Rosenberg & Hirschberg, 2007), which assesses the quality of the clusters by examining both the completeness and the homogeneity of the data within each group.

Rerank: The main evaluation metric for reranking in Finance MTEB is Mean Average Precision (MAP).

Pair-Classification: The main evaluation metric for Pair-Classification is Average Precision (AP), which measures the model’s accuracy across various decision thresholds.

Summarization: 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.

Retrieval: The main evaluation metric employed in this task is NDCG@10, which assesses the quality of the results based on their relevance and position in the list returned.

F COMPLEXITY SCORE CALCULATION

F.1 SHANNON ENTROPY IN INFORMATION THEORY

The Shannon entropy is a measure from information theory that quantifies the average level of uncertainty or information content inherent in a set of possible outcomes. A higher Shannon entropy means higher uncertainty. To calculate the Shannon entropy H of a text, we follow these steps:

1. Count Tokens: Identify all unique tokens w_i in the text and count their occurrences n_i .
2. Calculate Probabilities: Compute the probability of each token $P(w_i)$ by dividing its count by the total number of tokens N :

$$P(w_i) = \frac{n_i}{N}, \quad \text{where} \quad N = \sum_{i=1}^M n_i$$

Here, M is the total number of unique tokens.

- 1026 3. Compute Shannon Entropy: Use the probabilities to calculate the entropy, and sum up all
1027 unique tokens in the text.

$$1028 H = - \sum_{i=1}^M P(w_i) \log_2 P(w_i)$$

1032 F.2 READABILITY: GUNNING FOG INDEX

1033 The Gunning Fog Index (Gunning, 1952) is a readability metric that estimates the years of formal
1034 education needed to understand a text upon first reading, it evaluates the complexity of English prose
1035 based on sentence length and the frequency of complex words. To calculate the Gunning Fog Index
1036 (GFI), follow these steps:

- 1037 1. Select a Representative Passage
1038 Choose at least 1000 words from the text that represents the overall writing style.
1039 2. Calculate the Average Sentence Length (ASL)

$$1040 ASL = \frac{\text{Total Number of Words}}{\text{Total Number of Sentences}} \quad (2)$$

- 1041 3. Identify Complex Words
1042 • **Complex words** are words with **three or more syllables**.
1043 • Exclude proper nouns, familiar jargon, compound words, and verbs with common
1044 suffixes (e.g., “-es”, “-ed”, “-ing”).
1045 4. Calculate the Percentage of Complex Words (PCW)

$$1046 PCW = \left(\frac{\text{Number of Complex Words}}{\text{Total Number of Words}} \right) \times 100 \quad (3)$$

- 1047 5. Compute the Gunning Fog Index

$$1048 GFI = 0.4 \times (ASL + PCW) \quad (4)$$

1058 F.3 MEAN DEPENDENCY DISTANCE

1059 The mean dependency distance (MDD) (Oya, 2011) is introduced as a metric to quantify the syntactic
1060 complexity of sentences based on dependency parsing. A higher mean dependency distance
1061 indicates longer dependencies and potentially more complex syntactic structures. For each dependency
1062 relation between a word (a head) and its dependent in a sentence d is calculated as the absolute
1063 difference of their positions in the sentence:
1064

$$1065 d = |\text{Position}_{\text{head}} - \text{Position}_{\text{dependent}}|$$

1066 Here:

- 1067 • $\text{Position}_{\text{head}}$ is the position index of the head word.
1068 • $\text{Position}_{\text{dependent}}$ is the position index of the dependent word.

1069 The sentence-level MDD is computed by averaging the dependency distances of all its N dependency
1070 relations:
1071

$$1072 MDD_{\text{sentence}} = \frac{1}{N} \sum_{i=1}^N d_i = \frac{1}{N} \sum_{i=1}^N |\text{Position}_{\text{head}_i} - \text{Position}_{\text{dependent}_i}|$$

1073 Same with sentence-level, the document-level MDD averages the sentence-level mean dependency
1074 distances across all M sentences in the document:
1075

1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1110
1111
1112
1113
1114
1115
1116
1117
1118
1119
1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133

$$\text{MDD}_{\text{document}} = \frac{1}{M} \sum_{j=1}^M \text{MDD}_{\text{sentence}_j}$$

In our experiments, we calculate the document-level MDD for a test sample by averaging the MDD of all its text. For example, to compute the MDD for a pair-classification data point, we average the MDD of sentence 1 and sentence 2.

G PROMPT FOR CHATGPT ERROR RATE

The detailed example prompt of each task for the ChatGPT Error Rate is listed in Table tables 15 to 21.

Table 15: Prompt for the ChatGPT Error Rate on Semantic Textual Similarity (STS).

Determine whether the following two sentences are similar and answer yes or no.
Sentence1: Excluding the impact of merger-related costs, NSTAR Electric s earnings increased \$67.4 million in 2013, as compared to 2012, due primarily to lower overall operations and maintenance costs and higher retail electric sales due primarily to colder weather in the first and fourth quarters in 2013.
Sentence2: NSTAR Electric’s earnings increased in 2014, as compared to 2013, due primarily to lower operations and maintenance costs primarily attributable to lower employee-related costs and higher transmission earnings, partially offset by higher interest expense, higher depreciation expense, higher property tax expense and the after-tax reserve recorded for the 2014 FERC ROE orders as compared to the reserve recorded in 2013 for the FERC ALJ initial decision in the FERC base ROE complaints.

Table 16: Prompt for the ChatGPT Error Rate on Classification.

Classify the sentiment of a given finance text as either positive, negative, or neutral.
Text: Glencore shares hit 3-month high after refinancing key credit line

Table 17: Prompt for the ChatGPT Error Rate on Retrieval.

Given a financial question, retrieve user replies that best answer the question. Return the index.
Query: What is 'Obligor' ?
Corpus: {A List of 20 different context}

H SUPPLEMENTARY EXPERIMENT: FINETUNE SOTA EMBEDDING USING DOMAIN CORPUS

We fine-tuned e5-mistral-7b-instruct (Wang et al., 2023) using a syntenic finance QA dataset generated through Self-instruct(Wang et al., 2022) with GPT-4o mini (OpenAI, 2024a) and a manually labeled seed finance dataset. The results demonstrate clear domain adaptation effects:

On FinMTEB (Table 22), performance improved from 0.6475 to 0.6735, showing the benefits of finance-specific training. While general MTEB scores (Table 23) slightly decreased from 0.6463 to 0.6320, the model maintained competitive performance on broader tasks.

These results highlight how domain adaptation can enhance specialized task performance while preserving general capabilities.

1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187

Table 18: Prompt for the ChatGPT Error Rate on Clustering.

Identify industries from company descriptions.
Choices: Banking;Retail;Automotive;Aerospace;Financial services
Text: Cobridge Communications was a cable television, high-speed internet, and digital telephone service provider.

Table 19: Prompt for the ChatGPT Error Rate on Reranking.

Given a financial question, retrieve user replies that best answer the question. Return the index.
Query: How to tell if you can trust a loan company?
Corpus: {A List of 20 different context}

Table 20: Prompt for the ChatGPT Error Rate on Pair-Classification.

Classify the sentiment of a given finance text and determine whether label of two sentences are similar. Answer yes or no.
Sentence1: gold falls as dollar strengthens, etf holdings decline
Sentence2: Gold futures succumb to profit-booking, global cues

Table 21: Prompt for the ChatGPT Error Rate on Summarization.

Determine whether the following are documents and summary. Answer yes or no.
Text: deposits grew \$ 167.8 million , or 7 % , to \$ 2.504 billion at december 31 , 2020 , compared to \$ 2.336 billion at december 31 , 2019. non-interest bearing deposits grew by \$ 225.8 million in 2020 , or 20 % , and made up 54 % of total deposits at year-end . cost of deposits remained low at 0.11 % in 2020 , compared to 0.20 % in 2019. net interest income totaled \$ 96.7 million and \$ 95.7 million in 2020 and 2019 , resp....
Summary: cash and cash equivalents : our cash and cash equivalents , which include federal funds sold and short-term investments , were \$ 181.5 million at december 31

1188
1189
1190
1191
1192
1193
1194
1195
1196
1197
1198
1199
1200

Table 22: Overall Performance on FinMTEB

Model	STS	Ret.	Rerank.	Class.	Summ.	PairClass.	Clust.	Avg.
e5-mistral-7b-instruct	0.380	0.675	0.988	0.645	0.528	0.739	0.578	0.648
+ domain adaption	0.428	0.699	0.990	0.757	0.480	0.801	0.560	0.674

1204
1205
1206
1207
1208
1209
1210
1211
1212
1213
1214
1215
1216
1217
1218
1219
1220
1221
1222
1223
1224
1225
1226

Table 23: Overall Performance on MTEB

Model	STS	Ret.	Rerank.	Class.	Summ.	PairClass.	Clust.	Avg.
e5-mistral-7b-instruct	0.859	0.588	0.602	0.770	0.314	0.883	0.508	0.646
+ domain adaption	0.858	0.491	0.603	0.776	0.306	0.875	0.517	0.632

1231
1232
1233
1234
1235
1236
1237
1238
1239
1240
1241