# FINANCIAL BOND SIMILARITY SEARCH USING REPRESENTATION LEARNING

**Amin Haeri, Mahdi Ghelichi**
Model Development, Risk Management
TD Bank, Toronto, Canada
{amin.haeri, mahdi.ghelichi}@td.com

**Nishant Agrawal, David Li, Catalina Gomez Sanchez**
Model Development, Risk Management
TD Bank, Toronto, Canada
{nishant.agrawal, zhao.li, catalina.gomezsanchez}@tdsecurities.com

## ABSTRACT

Finding similar bonds remains challenging in fixed-income analytics, as numerical financial attributes often overshadow categorical non-financial ones such as issuer sector and domicile. This paper shows that these categorical attributes dominate the predictability of spread curves and proposes embedding models to capture their semantic similarities, outperforming one-hot and many other baselines. Evaluated via sparse-issuer augmentation, the approach improves risk modeling and curve construction.

## 1 Introduction

In fixed-income markets, bonds are debt instruments through which governments, corporations, and other entities obtain capital from investors. They specify contractual terms such as coupon rate, maturity, and redemption method, which define how interest and principal are repaid. The bond universe spans sovereign, corporate, municipal, and structured issues, each carrying distinct risk and return patterns shaped by credit quality, issuer profile, and market dynamics. Given this diversity, comparing bonds is essential for portfolio construction and risk management.

Similarity search provides a data-driven framework to identify bonds with comparable characteristics across large datasets. Building on decades of research in distance metric learning (Kulis, 2013) and information retrieval (Li, 2014), similarity-based approaches have proven particularly valuable in financial contexts where direct price discovery is limited due to market microstructure constraints. It assesses closeness not only in financial features such as yield, duration, and credit risk but also in contextual dimensions like issuer industry, credit rating, and descriptive information. For instance, an algorithmic comparison should recognize that a bond issued by the Province of Manitoba is more closely related to other Canadian provincial issuers than to a corporate bond issued by Amazon. Such methods improve the relevance of peer selection and enhance the reliability of bond curve construction when direct market data are limited.

In practice, non-financial categorical attributes such as issuer industry, domicile, and market of issuance often explain more of the cross-sectional structure in spreads than marginal differences in numerical terms such as coupon or residual maturity. For example, a 7-year bond issued by the Province of Manitoba tends to behave more like other Canadian provincial issues of similar rating than a corporate bond issued by Amazon with nearly identical maturity and coupon, despite the latter being numerically closer on standard term-structure dimensions. Traditional similarity engines usually start from numerical financial variables and treat categorical descriptors as afterthoughts or filters, which can underweight the role of issuer-level and sectoral context in shaping spreads. The central claim of this paper is that these categorical non-financial attributes should instead be treated as the primary drivers of bond similarity. Recent advances
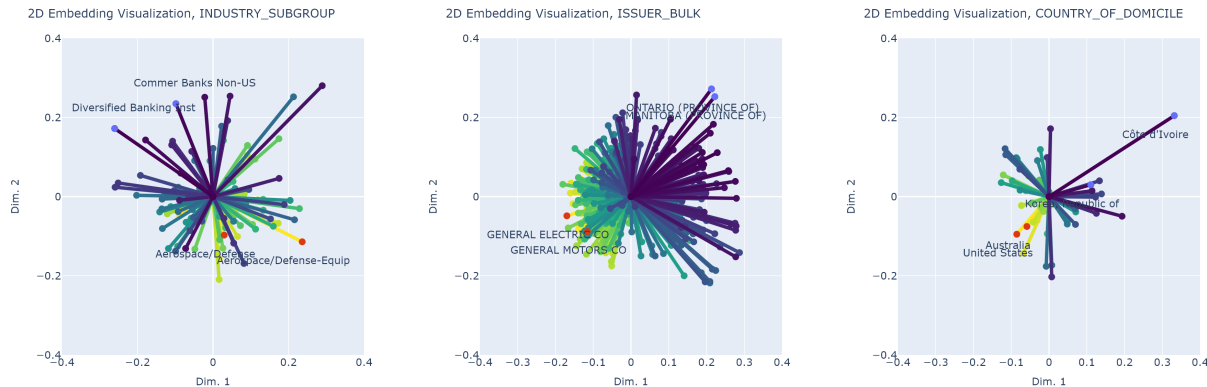
Figure 1: Two-dimensional embedding projections colored by their similarities in the high-dimensional space. From left to right: Industry Subgroup, Issuer Bulk, and Country of Domicile. Each point represents an entity in the learned embedding space, projected onto the first two dimensions, with radial lines indicating displacement from the origin. The visualizations illustrate how semantically related entities cluster and separate according to industry, issuer identity, and geographic domicile, highlighting the embedding's ability to capture structured relationships across different categorical views.

in representation learning (Bengio et al., 2013; Devlin et al., 2019) demonstrate that learned embeddings of categorical features can capture semantic relationships more effectively than explicit one-hot encodings; we extend this finding to the fixed-income domain. Handling categorical information correctly is therefore more important than adding further numerical detail for similarity assessment in bond markets.

Representation learning is a machine learning paradigm focused on automatically discovering useful features or representations of data, rather than relying on handcrafted variables. The goal is to transform raw inputs such as text and structured financial attributes into compact, informative vectors that capture semantic, structural, and statistical relationships. These learned representations enable downstream tasks like classification, clustering, retrieval, or prediction to be performed more effectively and with less domain-specific engineering. Advances in deep learning have made representation learning especially powerful in financial contexts (Goodfellow et al., 2016; Heaton et al., 2017), with models such as autoencoders, word embeddings, and transformer-based architectures that capture rich contextual patterns. For high-cardinality categorical variables (such as bond issuers or industry classifications), learned embeddings overcome the curse of dimensionality inherent in one-hot encoding, producing continuous latent representations where semantic proximity is reflected by vector distance.

Conventional approaches to bond similarity search rely heavily on handcrafted features, rule-based filters, and simple distance metrics applied to a limited set of attributes, such as coupon, maturity, or credit rating. While straightforward, these methods often fail to capture the full complexity of bond structures and issuer-specific nuances, especially when unstructured information from prospectuses or covenant terms is involved. More fundamentally, for high-cardinality categorical variables (e.g., issuer identity, industry classification), one-hot encoding treats all categories as orthogonal and prevents the model from learning meaningful similarities between categories. These conventional approaches can be brittle when handling missing data, issuer sparsity, or evolving market conditions, leading to inaccurate or incomplete similarity assessments. They lack the flexibility to adapt to new issuers or market regimes, restricting their utility in dynamic fixed-income markets.

In this paper, rather than positioning our work as another supervised or unsupervised metric-learning approach, we focus on the distinction between financial and non-financial attributes and on how to represent high-cardinality categorical information. The proposed framework uses embedding models to map these categorical attributes into a latent space where cosine similarity reflects economically meaningful relationships, and evaluates their usefulness via CDS spread-curve estimation from sparsified issuer catalogs.

## 2   Prior Work

The majority of prior work on bond similarity and term-structure modeling has focused on numerical financial variables (i.e., yields, spreads, durations, and volatilities) combined with distance metrics or predictive models, while treating issuer-level and sectoral descriptors as coarse controls or exclusions. Unsupervised clustering and supervised metric

learning have both been explored in this numerical regime, but the relative importance of non-financial categorical attributes and the question of how to represent them remain largely unaddressed.

Early research predominantly focused on unsupervised clustering methods such as K-means, hierarchical clustering, self-organizing maps, and Gaussian mixture models. These methods group bonds based on attributes like coupon rate, yield-to-maturity, duration, and credit rating. They typically operate on standardized numerical features; industry or rating dummies, if used at all, are encoded via sparse one-hot vectors that do not express similarity between categories. While useful for broad classification tasks, such as distinguishing between investment-grade and high-yield bonds, unsupervised methods suffer from several limitations. They rely on predefined distance metrics (e.g., Euclidean or Minkowski), which may not align with the intrinsic geometry of financial data.

More recent work has shifted toward supervised similarity learning, where similarity is defined relative to a prediction target (e.g., yield or spread). In this paradigm, machine learning models are trained to predict bond characteristics, and proximity measures are extracted from the resulting model structure. Random forests, in particular, have been shown to be effective distance-metric learners, with theoretical justification grounded in the adaptive nearest-neighbors interpretation (Jeyapaulraj et al., 2022; Scornet et al., 2015). Proximity between two securities is defined by the proportion of trees in which they co-occur in the same terminal node. Variants such as out-of-bag (OOB) proximity (Breiman, 2004) and geometry-and-accuracy-preserving (GAP) proximities (Mentch & Hooker, 2016) have been proposed to improve robustness and interpretability. In financial applications, Random Forest-based similarities have proven scalable to high-dimensional datasets and naturally handle mixed categorical and numerical features (Desai & Mehta, 2021).

Building on these advances, quantum cognition machine learning (QCML) has recently been introduced as a novel paradigm for supervised distance metric learning (Rosaler et al., 2025). QCML leverages the mathematical formalism of quantum theory and foundational ideas from quantum cognition (Busemeyer & Bruza, 2012) by representing data points as quantum states in a Hilbert space, with features and targets encoded as quantum observables. Similarity is defined via quantum fidelity, a measure of overlap between two quantum states derived from quantum information theory, specifically the absolute value of the inner product between two quantum states. A key theoretical advantage of QCML is its logarithmic scaling of parameters with respect to feature dimensionality (the parameter count scales as $O(N^2)$), where $N$ is the Hilbert space dimension), in contrast to random forests, where tree depth typically scales linearly with the number of features, resulting in exponential growth in parameters and leaves. This logarithmic economy of representation endows QCML with superior generalization in high-dimensional, sparse data regimes. Empirical results demonstrate that QCML outperforms random forests in predicting yields and deriving similarity metrics for high-yield bonds, while performing comparably in investment-grade markets. This suggests that QCML may provide a robust alternative to tree-based methods for illiquid and noisy financial datasets, particularly when the bond universe exhibits high feature dimensionality (due to one-hot encoding of categorical variables) and a high concentration of outliers (as in high-yield markets).

While supervised similarity metrics based on tree ensembles and quantum models excel in capturing complex nonlinear relationships through prediction tasks, they remain bound by the representations inherent in their training data. An orthogonal but complementary line of inquiry is representation learning for categorical variables themselves. Transformer-based embeddings have demonstrated remarkable ability to capture semantic relationships in high-cardinality categorical spaces (Devlin et al., 2019; Vaswani et al., 2023). However, their application to financial similarity, particularly in the context of bond markets where categorical non-financial attributes may dominate numerical features, remains underexplored. The present work bridges that gap by demonstrating empirically that embedding-based representations of categorical bond attributes outperform classical approaches (one-hot and numerical baselines) in sparse-issuer regimes and achieve comparable or superior performance to supervised metric learners such as random forests, particularly when data sparsity is pronounced.

## 3   Representation Learning

The core challenge in financial bond similarity search lies in constructing a representation of bonds that captures the complex relationships among their attributes, market conditions, and liquidity characteristics. Traditional methods either rely on predefined distance metrics or extract similarity from specific model architectures, but they often fail to generalize across diverse market regimes or adequately encode the nuanced interactions among features. Representation learning provides a principled approach to address these challenges by automatically learning embeddings of textual and categorical data into a latent space, where proximity reflects meaningful similarity (Bengio et al., 2013). In particular, for high-cardinality categorical variables (such as bond issuer identity, with millions of possibilities across global markets, or granular industry classification), one-hot encoding suffers from the curse of dimensionality: it produces sparse, orthogonal representations that prevent the model from learning meaningful similarities between categories.

Learned embeddings, by contrast, map categories into a continuous latent space where semantic proximity is captured by vector distance. This enables superior generalization in sparse-data regimes, where the model must predict similarity for issuer or sector combinations not extensively represented in training data.

With the growth of unstructured financial data (such as issuer disclosures, news sentiment, and transcripts) deep learning architectures provide opportunities to learn richer representations. Autoencoders, variational autoencoders, and graph neural networks (GNNs) have been proposed to extract embeddings that preserve both numerical and relational information. For corporate bonds, where issuers and securities are embedded in complex credit networks, GNNs can capture structural relationships (e.g., shared sector exposure or counterparty risk) that are overlooked by purely tabular models. Similarly, transformers trained on textual data can produce embeddings that complement traditional quantitative features, enabling similarity search.

Transformer-based embedding models, such as BERT (Bidirectional Encoder Representations from Transformers) (Devlin et al., 2019) and sentence BERT (sBERT) (Reimers & Gurevych, 2019), are designed to understand the context of words within a sentence by processing the input text bidirectionally, meaning they consider both the preceding and following words to capture deeper semantic meaning. This bidirectional approach allows transformers to go beyond simple word associations and capture more nuanced aspects of language, such as word sense disambiguation (WSD) (Peters et al., 2018; Vaswani et al., 2017). When applied to embedding categorical financial data, a transformer model generates dense vector representations that capture not only syntactic meaning but also the intent, context, and relationships between financial entities (Mikolov et al., 2013). For example, when embedding bond issuer sectors (Technology, Retail, Manufacturing, etc.), transformers learn to recognize that Technology and E-Commerce are semantically closer than Technology and Banking, a distinction reflecting sector risk correlation and market structure. While BERT processes individual tokens and learns contextual embeddings for each word in a sentence, sBERT is specifically optimized to generate embeddings for entire sentences or documents. sBERT, which builds upon BERT's architecture, is fine-tuned to produce fixed-size dense vectors that represent the overall meaning of a sentence or document, making them ideal for tasks such as semantic similarity, clustering, and information retrieval (Reimers & Gurevych, 2019; Sanh et al., 2020). In the financial domain, contextual embeddings have demonstrated effectiveness in capturing domain-specific semantics of issuer sectors, credit ratings, and market structures.

Pretrained language models such as BERT are trained on large, general-domain corpora (e.g., Wikipedia, BookCorpus) and capture broad linguistic and semantic patterns. However, the financial domain exhibits distinctive characteristics: issuers cluster within sectors and geographies; credit ratings carry specific risk implications; and institutional conventions shape how bonds are classified. Domain-specific fine-tuning adapts the pretrained model's learned representations to these financial nuances, resulting in embeddings where semantic proximity more faithfully reflects economically meaningful relationships (Gururangan et al., 2020; Haeri et al., 2025). This is particularly important for bond similarity search, where the embedding model must recognize subtle but economically crucial relationships (e.g., that a Canadian provincial issuer is more similar to other provincial issuers regardless of sector, than to a numerically similar corporate bond from another country). Fine-tuned embeddings capture such domain-specific hierarchies that purely general-purpose models would miss, delivering superior performance on financial downstream tasks.

This section has reviewed the theoretical foundations and practical advantages of representation learning for categorical variables. The embedding-based approach proposed in this work, where similarity is defined by learned semantic relationships among categories, is deliberately orthogonal to prediction-based supervised similarity learning methods reviewed in Section 2. Where supervised approaches (random forests, QCML) extract similarity from the structure of a trained model or learned quantum geometry optimized for a prediction task (e.g., yield prediction), embeddings capture similarity defined by the inherent semantic structure and relationships among categories themselves, independent of any specific prediction target. This is a fundamental distinction: embeddings ask "How similar are issuer A and issuer B in conceptual or semantic space?", while supervised methods ask "Given that we wish to predict yield, which bonds' yields should influence the prediction for a test bond?". Both are valid notions of similarity, and they need not align. The comparative strengths of embedding-based and prediction-based approaches are empirically evaluated in Sections 5–6, where we demonstrate that embedding-based representations of non-financial categorical attributes outperform one-hot encodings and achieve comparable or superior performance to tree-based methods like random forests, particularly in sparse data regimes where semantic similarity becomes especially valuable. Future work may investigate hybrid models that integrate learned categorical embeddings with supervised metric learning frameworks, exploiting complementary strengths in both categorical representation and numerical feature dynamics.

## 4   Data Description

For this study, we construct our dataset using multiple years of daily bond data obtained from Bloomberg, covering a broad cross-section of issuers and sectors. The dataset consists of thousands of individual securities, each observed

Table 1: Hand-selected categorical features for bond similarity learning.

| Feature | Category | Description |
|---|---|---|
| Issuer Industry | Issuer Attribute | Broad sector classification of the issuing firm (e.g., *Industrial*). |
| Market Issue Type | Market Attribute | Market scope of the bond (e.g., *Global*). |
| Industry Group | Market Attribute | Intermediate sector grouping (e.g., *Computers*). |
| Industry Subgroup | Market Attribute | Fine-grained industry category (e.g., *E-Commerce*). |
| Country of Domicile | Geographic Attribute | Issuer's country of registration or operation (e.g., *US*). |
| Issuer Identity | Issuer Attribute | Categorical identifier of the issuing firm (e.g., *APPLE*). |

over time and characterized by more than 200 raw attributes, including pricing, liquidity, structural, and issuer-specific variables.

Given the high dimensionality of the raw data, not all attributes are equally informative for the task of similarity search. Many features are redundant, noisy, or weakly correlated with target measures such as yield, spread, or liquidity. To focus the analysis and reduce dimensionality, we conducted a systematic feature selection process that combined domain expertise with empirical analysis of variable stability and predictive value.

As a result, we hand-selected six core features that capture the most relevant aspects of bond similarity. All six of these features are categorical variables, which allows us to group bonds along discrete economic or structural dimensions that are meaningful to traders and portfolio managers. The selected features are outlined in Table 1. These core variables are non-financial in the sense that they do not directly encode prices, spreads, or cash-flow magnitudes. They characterize issuer identity, sector, and geographic footprint, which the empirical analysis shows to be highly informative. In addition, the categorical nature of these features is particularly well suited to representation learning approaches. Traditional distance metrics often struggle with high-cardinality categorical variables (e.g., issuer or industry), as they lack a natural ordering or scale. By embedding these categories into a latent representation space, similarity learning models can uncover meaningful relationships (e.g., between industries with correlated risk, or between issuers within the same rating tier) that are not explicit in the raw data.

Although this study narrows its focus to six categorical features for clarity and interpretability, the framework is extensible. Future work may reincorporate additional attributes from the broader Bloomberg dataset. Subsequent experiments explicitly contrast representations based on these categorical attributes with one-hot encodings and with models that rely primarily on numerical financial variables, in order to quantify their relative importance for spread curve reconstruction.

## 5 Methodology

Our methodology integrates *representation learning* with *post-processing filters* to construct a practical and robust framework for bond similarity search. The process consists of two main components: (1) generating similarity scores via a fine-tuned, pretrained embedding model, and (2) applying post-filters to refine the ranked list of candidate bonds. Finally, we evaluate our methodology through a structured experimental setup that measures how effectively the augmented bond catalog can reproduce realistic bond spread curves compared to the actual market data (see Section 5.3).

### 5.1 Embedding

We begin with a pretrained embedding model that has been fine-tuned on a broad range of financial data (Haeri et al., 2025). Fine-tuning enables the model to learn domain-specific representations in which bonds with similar categorical profiles are embedded close to one another in latent space. Then, each bond $b_i$ with feature vector $x_i$ is mapped into an embedding vector $z_i \in \mathbb{R}^d$, where $d$ is the dimension of the learned latent space. Similarity between two bonds $b_i$ and $b_j$ is computed as:

$$\text{sim}(b_i, b_j) = \frac{\langle z_i, z_j \rangle}{\|z_i\| \cdot \|z_j\|}, \tag{1}$$

i.e., the cosine similarity between their embedding vectors. For each query bond, the model produces a *ranked list of nearest neighbors* by sorting other bonds in descending order of similarity. As illustrated in Figure 2, this approach leverages the advantages of representation learning: categorical features are embedded into continuous vectors that
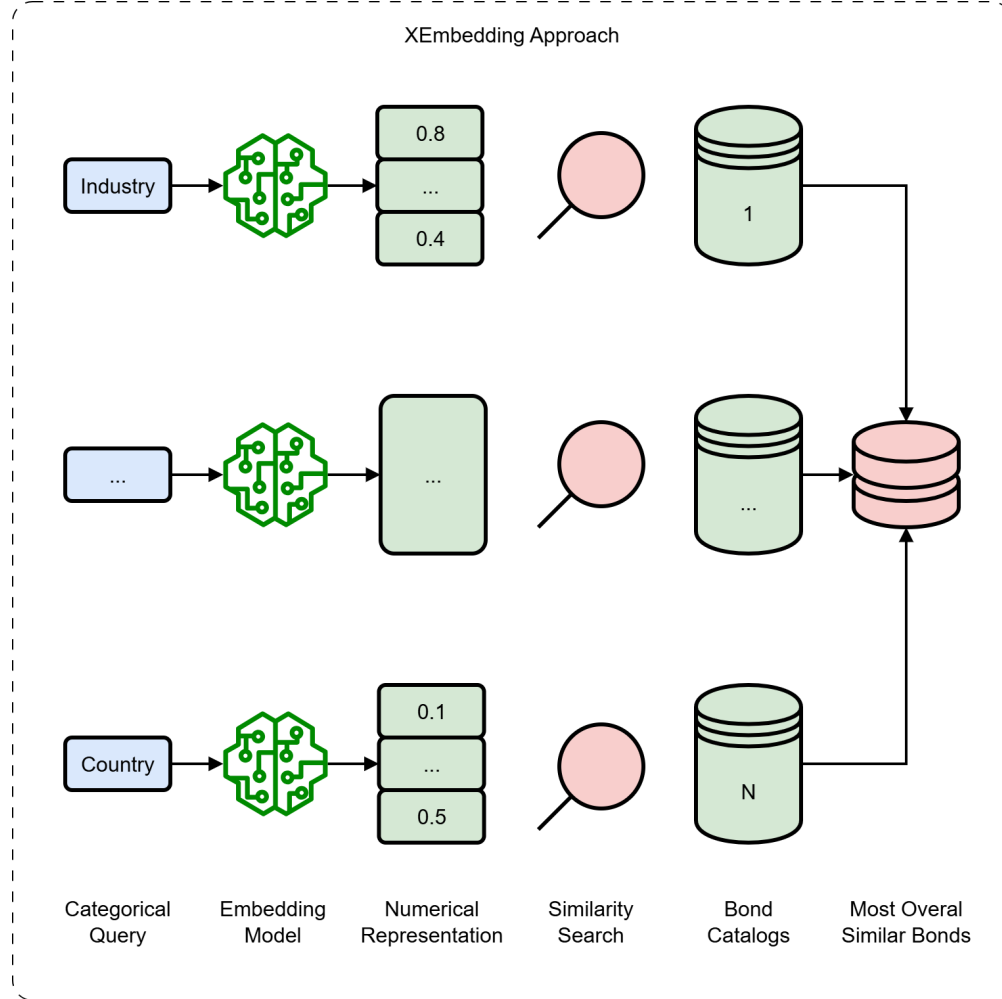
Figure 2: Overview of the embedding-based approach for bond similarity search. Categorical queries (e.g., Industry, Country) are passed through a fine-tuned embedding model to produce dense numerical representations. These embeddings are then compared via similarity search against bond catalogs. After aggregating results across all features, the system outputs the most overall similar bonds.

preserve semantic relationships (e.g., issuers in related industries). Unlike handcrafted distance metrics, the embedding model adaptively captures these relationships during training.

It is worth mentioning that we initially proposed an embedding-based methodology in which all categorical bond features (e.g., industry, rating, and country) were concatenated into a single textual sequence and jointly passed through the embedding model (as illustrated in Figure 3). While this formulation captured complex interactions among features and resulted in slightly higher similarity accuracy, it suffered from reduced interpretability and made it difficult to isolate the contribution of each feature to the final similarity score. To enhance explainability, we revised the design to the current approach, where each categorical feature is independently embedded before the similarity search and aggregation stages. This modification provides greater transparency by allowing analysts to examine feature-level similarities across bonds. However, because feature interactions are no longer jointly modeled, the overall similarity accuracy experiences a minor reduction, reflecting a deliberate trade-off between model performance and interpretability. The categorical embeddings derived from this model are referred to as XEmbedding.

## 5.2 Post-Filtering

While embeddings provide a powerful foundation for similarity search, additional post-filters could be utilized to ensure that results satisfy practical trading and portfolio management constraints. After generating the sorted list of candidate bonds from the embedding model, we apply the following sequence of filters to refine the output:
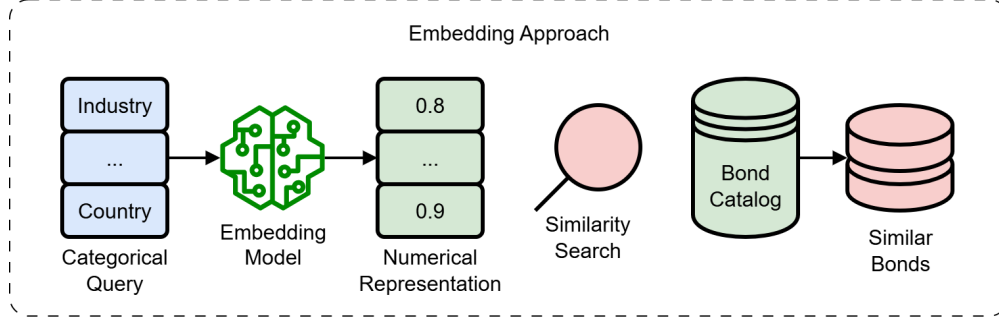
Figure 3: Initial embedding-based methodology for bond similarity search. In this design, all categorical features are concatenated and jointly passed through the embedding model to produce dense numerical representations. Although this approach captures feature interactions more effectively, it offers less interpretability compared to the revised, feature-wise embedding approach.

1. **Currency Filter** – Retain only bonds denominated in the same currency as the query bond.
2. **Maturity Filter** – Restrict candidates to those within a defined lower- and/or upper-bound time to maturity.
3. **Rating Filter** – Maintain credit quality alignment by allowing only bonds within a predefined rating tolerance.

These filters are applied progressively, narrowing down the candidate set. The final output is a refined, ranked list of bonds that balances *semantic similarity in the embedding space* with *practical market constraints*.

## 5.3 Evaluation

To evaluate the effectiveness of the proposed bond similarity search and augmentation framework, we design an experiment based on the process illustrated in Figure 4. The objective is to assess how well the augmented bond catalog can reproduce realistic bond spread curves compared to the actual market data.

We begin with the *bond catalog*, which contains multiple issuers with varying bond densities. Some issuers have sufficient bonds across different maturities (non-sparse issuers), while others have only a few outstanding bonds (sparse issuers). For evaluation purposes, we first select a subset of non-sparse issuers, ensuring that each selected issuer has a sufficiently rich term structure to serve as ground truth. From each selected issuer, we randomly drop a fixed number of bonds to artificially create sparsity in their term structures. These reduced sets form the *sparse bond catalog*, mimicking real-world cases where certain issuers have limited bond data available.

Next, we apply our proposed augmentation methodology to these sparse issuers. The similarity-based framework identifies comparable bonds from other issuers in the catalog and augments the sparse issuer's data with similar bonds. The resulting *augmented bond catalog* represents an enhanced dataset where each issuer's term structure is partially reconstructed using embedding-based similarity retrieval. Once the augmented dataset is generated, we use a Nelson–Siegel (NS) model to fit the *predicted bond spread curve* for each issuer. We then compare these predicted curves against the *actual bond spread curves* derived from the original (non-sparse) catalog. The discrepancy between the predicted and actual spread curves is quantified using the Root Mean Square Error (RMSE) metric:

$$\text{RMSE} = \sqrt{\frac{1}{N} \sum_{i=1}^{N} (\hat{y}_i - y_i)^2} \tag{2}$$

where $\hat{y}_i$ and $y_i$ denote the predicted and actual bond spreads, respectively. A lower RMSE indicates that the augmented issuer more accurately reproduces the true bond spread dynamics, reflecting the quality of the similarity-based augmentation.

## 5.4 Benchmarking

To rigorously assess the value of embedding-based representations for categorical non-financial bond attributes, our model is benchmarked against a one-hot encoded baseline that represents the conventional approach to handling categorical variables in financial similarity search. Both models operate within the same experimental framework of sparse-issuer augmentation followed by Nelson–Siegel spread-curve fitting, ensuring a controlled comparison of
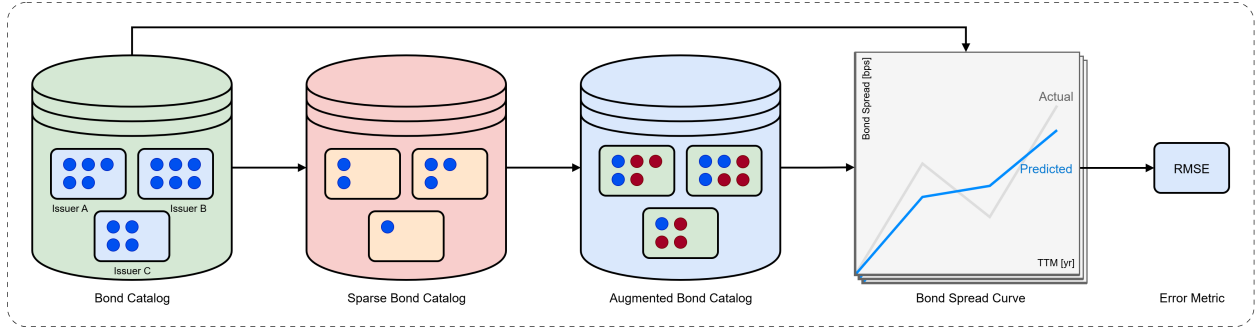
Figure 4: Evaluation pipeline for the proposed bond similarity search framework. Starting from a complete bond catalog, a subset of non-sparse issuers is selected. A fixed number of bonds are randomly removed to create sparsity. The sparse issuers are then augmented using the proposed similarity-based methodology. The augmented catalog is used to fit predicted bond spread curves via the NS model, which are compared to the actual curves using the RMSE error metric.

representation quality while holding all other aspects including the data universe, sparsity protocol, post-filters, and evaluation metric constant.

The one-hot baseline encodes each of the six categorical features (issuer industry, market issue type, industry group, subgroup, country of domicile, issuer identity, currency, and bond rating) as sparse binary vectors. For a given query bond, similarity is computed as a weighted average of exact-match frequencies across features, where weights reflect domain-informed priorities (e.g., issuer industry and rating receive higher weights than market issue type). The top-K nearest neighbors are then selected using this aggregate score, identical to the XEmbedding pipeline, and subjected to the same post-processing filters. These neighbors augment the sparse issuer catalog, and the NS model is fit to the augmented term structure to predict spreads across maturities.

This setup isolates the effect of categorical representation: one-hot assumes orthogonality between categories and relies on exact matches or simple aggregation, while our model captures semantic proximity through dense learned embeddings. This methodology quantifies whether nuanced categorical similarity (XEmbedding) meaningfully improves spread curve reconstruction over rigid one-hot encodings, directly testing the hypothesis that non-financial attribute representation dominates numerical financial detail in sparse-data regimes.

## 6  Results and Discussion

In this section, the proposed similarity search methodology is applied to a large universe of corporate bonds (i.e., 2,500 bonds from 250 unique issuers) to identify instruments that share the most comparable characteristics with a given query bond. Furthermore, to assess the effectiveness of the proposed bond similarity search and augmentation methodology, we conduct a series of experiments comparing the predicted bond spread curves against the actual market data. The evaluation is based on a set of issuers with varying levels of bond sparsity, mimicking real-world scenarios where some issuers have sparse bond data while others have more comprehensive term structures. A more extensive set of example results is provided in the Appendix.

### 6.1  Similarity Search

The similarity measure is derived using an embedding approach that employs categorical attributes of each bond. Specifically, categorical features such as issuer industry and industry group were transformed into dense vector representations through an embedding model. For instance, the visualization shown in Figure 5 demonstrates how the method retrieves bonds that are structurally and economically similar to the query bonds from Apple Inc. (AAPL) and Bank of America Corp. (BAC) (issued on 2024-12-13). The heatmap presents the categorical profile of the query and its nearest neighbors, ranked by their similarity scores, which are shown along the right-hand side. The first row represents the query bond itself, serving as the reference point with a similarity score of 1.00. Subsequent rows correspond to bonds with progressively lower scores, indicating decreasing levels of similarity. The retrieved bonds cluster around specific industry and issuer types consistent with the market characteristics of Apple and Bank of America Corp. The top-ranked matches include bonds and issuers that share the similar global industrial profiles and operate within similar subsectors. This demonstrates the model's capacity to discern economic comparability within the broader industrial and consumer sectors, rather than merely matching by superficial categorical overlap. Further down the rankings, the
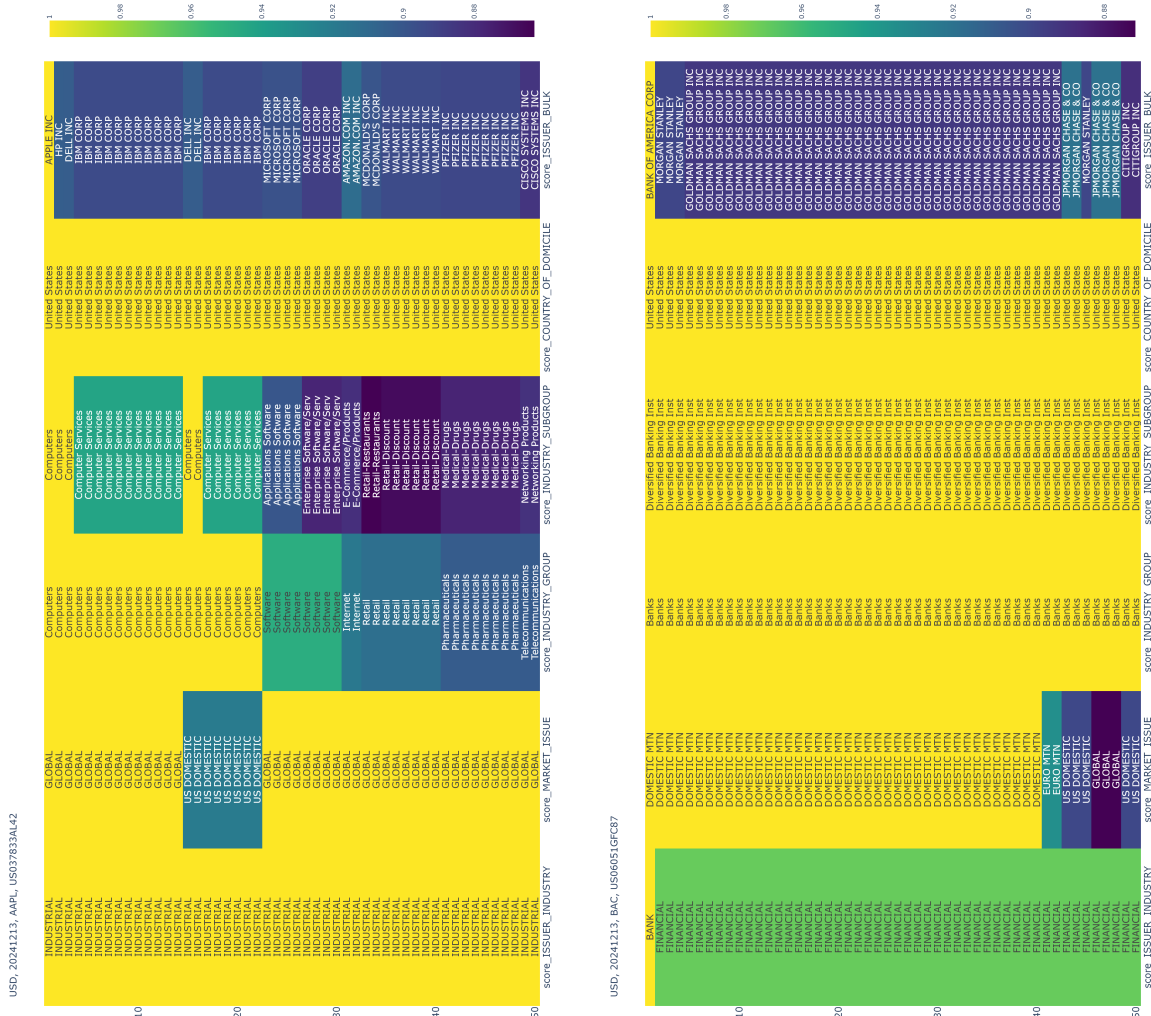
Figure 5: Visualization of bond similarity search results for the query bonds of AAPL US037833AL42 (left) and BAC US06051GFC87 (right). The first row represents the query bond's profile, while subsequent rows show the most similar bonds ranked by cosine similarity in embedding space. The color intensity indicates similarity, with higher scores (topmost column) reflecting closer structural and economic resemblance to the query bond.

model surfaces bonds issued by companies that, while belonging to different subsectors, remain within the same global corporate class and exhibit similar credit quality. This reflects the embedding model's ability to capture second-order similarity: bonds that differ in industry but align in credit and structural attributes. The gradual decline in similarity scores illustrates a smooth transition from highly comparable peers to tangentially related securities, confirming the embedding space's continuous structure.

Moreover, to visualize embedding quality, we construct two-dimensional projections of the embedding space for three key features (industry subgroup, issuer bulk, and country of domicile) shown in Figure 1. Points are colored by their similarity to a reference category in the original high-dimensional space (i.e., with 768 dimensions); the smooth color gradients in the projections indicate that the 2D layout remains faithful to the underlying categorical similarity structure. Across all three views, entities with similar categorical attributes exhibit clear clustering patterns, indicating that the embedding space captures meaningful semantic and structural relationships. In particular, separation by issuer and geographic domicile suggests that both firm-specific and regional signals are strongly encoded, while the industry-level view reflects finer-grained similarities within related sectors. Together, these projections provide qualitative evidence that the learned representations preserve relevant domain structure beyond the original feature space.
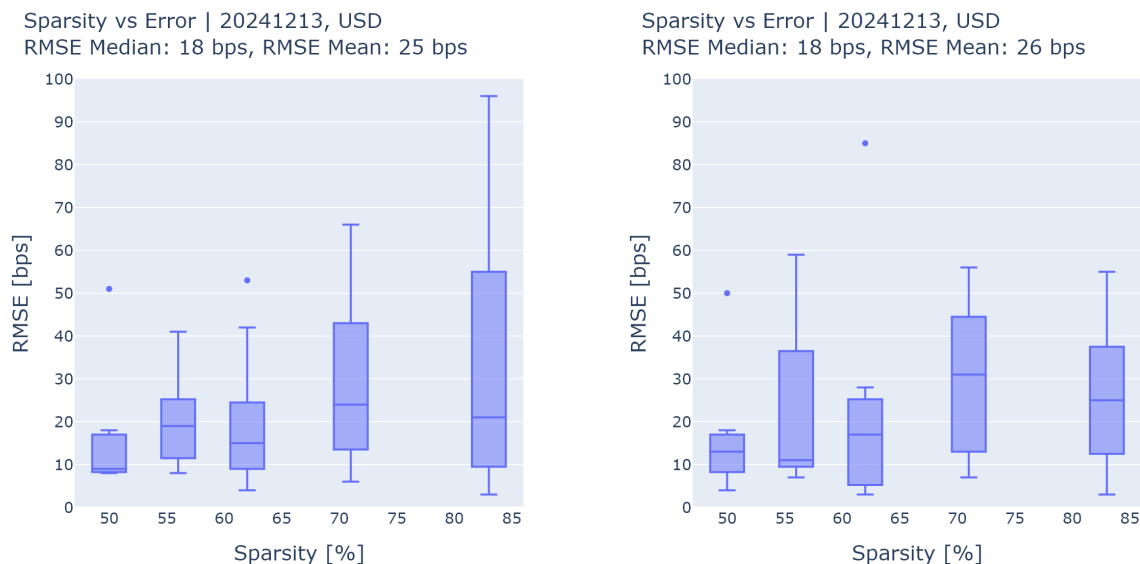
Figure 6: Comparison of overall performance of our model without (left) and with (right) post-filters. Post-filtering improves model stability under sparsity by tightening error distributions and reducing tail risk.

Overall, these results validate the robustness of the similarity search methodology. By combining categorical embeddings with post-filtering constraints, the model successfully prioritizes economically relevant comparables while maintaining diversity in issuer types and industries. The visualization provides intuitive interpretability, allowing researchers and practitioners to observe both the alignment of categorical attributes and the gradient of similarity scores across potential matches. This approach offers significant utility in fixed-income research applications such as relative value analysis, liquidity estimation, and risk benchmarking, where identifying functionally similar bonds is essential.

## 6.2 Bond Spread Curve

The results in this subsection evaluate the effectiveness of our model through CDS spread-curve reconstruction, using the Nelson–Siegel (NS) model fit to augmented sparse-issuer catalogs. Figure 6 compares the overall performance of our model without (left) and with (right) post-filters. Both panels plot RMSE against sparsity levels (where sparsity is defined as #Queries/(#Queries + #Similars)) for multiple issuers, using boxplots to summarize the distribution of errors at each sparsity bin alongside individual issuer scatter points. The x-axis spans increasing sparsity from left to right, while the y-axis captures RMSE in basis points, revealing the expected positive trend where typically higher sparsity correlates with elevated prediction errors due to fewer query bonds available for augmentation.

The left panel shows the model results without post-filtering, achieving a median RMSE of 18 bps and mean of 26 bps across issuers driven by extreme outliers exceeding 90 bps that widen the whiskers dramatically. In contrast, the right panel shows the post-filtered model achieving a comparable median of 18 bps but a lower mean RMSE of 25 bps. Post-filtering further tightens the error distribution and preserves strong clustering even under high sparsity (70–85%). Overall, beyond illustrating aggregate model performance, this comparison highlights how post-filtering mitigates sensitivity to poor categorical matches in sparse regimes, leading to improved stability evidenced by narrower boxplots, reduced tail risk, and more consistent performance.

We use Apple Inc. (AAPL) and Bank of America Corp. (BAC) as representative examples to illustrate the performance of our framework for augmenting sparse bond catalogs. Figure 7 compares actual and predicted bond spread curves as a function of time to maturity for the two issuers (Apple on the left and Bank of America on the right on 2024-12-13) both under a high sparsity regime of 83%. Blue solid lines represent observed spreads, while orange dashed lines denote model predictions constructed from limited tenor information. In both panels, the model captures the overall term-structure shape well, reproducing the steep increase at short maturities and the gradual flattening at longer horizons. Individual bond observations are overlaid, illustrating how sparse and uneven the available data are across maturities.

Quantitatively, the model achieves an RMSE of 23 bps for Apple and a much lower RMSE of 3 bps for Bank of America, indicating issuer-dependent performance under identical sparsity conditions. For Apple, larger deviations
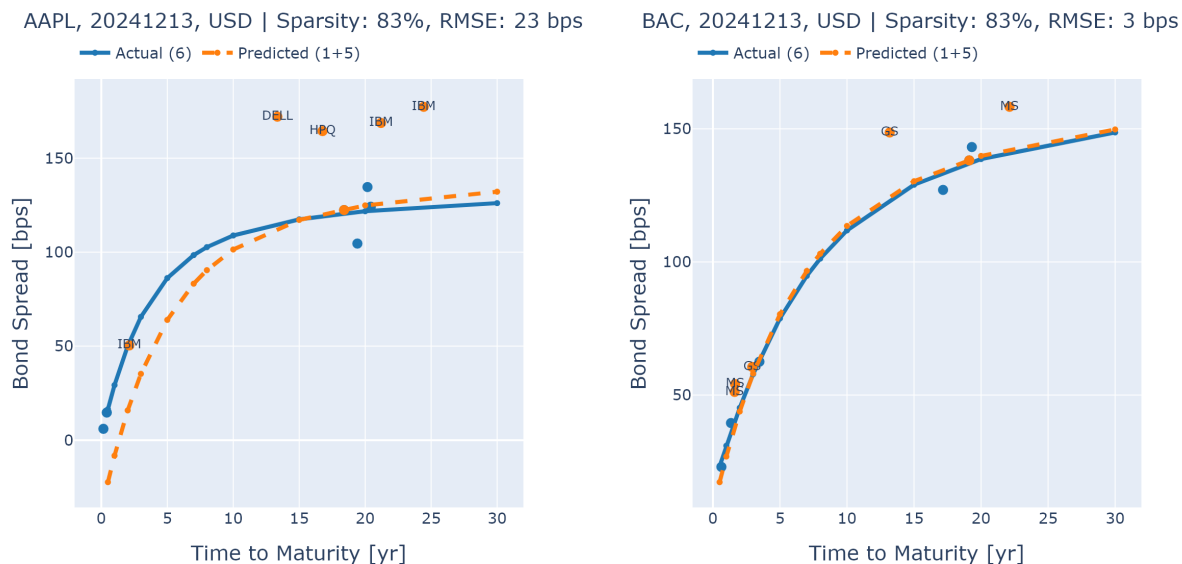
Figure 7: Comparison of predicted and actual CDS spread curves for AAPL (left) and BAC (right) bonds on 2024-12-13 under a high sparsity regime (83%). Blue markers represent the actual bonds, non-annotated orange markers show predictions for the query, and annotated orange markers indicate specific predicted bonds (e.g., DELL, IBM, and GS). Despite limited information, the model captures the overall curve shape, achieving RMSEs of 23 bps for AAPL and 3 bps for BAC.

appear at longer maturities, where sparse categorical matches lead to noticeable but controlled prediction errors. In contrast, the BAC curve shows near-perfect alignment between actual and predicted spreads across the full maturity range, demonstrating strong robustness even with limited data. Together, these examples highlight that while high sparsity can expose sensitivity to issuer-specific structure, the model generally preserves curve shape and delivers stable predictions, with performance improving when categorical alignment is strong.

The results demonstrate that our bond similarity search and augmentation framework is capable of effectively reconstructing bond spread curves, with a relatively low RMSE and good performance in cases of significant data sparsity. This suggests that the approach can be applied to enhance datasets with missing or sparse bond data, providing more realistic bond term structures for credit-risk modeling and investment decisions. Additional qualitative and quantitative results across a broader set of issuers are provided in Appendix A.

Additionally, we benchmark our model (the "XEmbedding" model) against a baseline that encodes categorical features using one-hot representations (the "one-hot" model; with post-filters applied to both). They differ fundamentally in how they represent categorical bond attributes, leading to stark contrasts in spread curve reconstruction accuracy. The one-hot approach encodes features like issuer industry, domicile, and market type as sparse binary vectors, treating categories as orthogonal with no inherent similarity, thus relying on exact matches or simple aggregation rules to select neighbors. In contrast, XEmbedding uses dense transformer-based embeddings to capture semantic and hierarchical relationships, such as placing "Retail–Discount" closer to "E-Commerce" than to "Banks".

The plots shown in Figure 8 summarize the aggregate performance of the XEmbedding model (left) versus the one-hot baseline (right) across multiple issuers. The comparison of the two models across varying sparsity levels indicates that our model consistently outperforms the one-hot baseline. The XEmbedding model exhibits lower median and mean RMSE values (18 bps and 26 bps, respectively) compared to the one-hot model (20 bps and 27 bps), reflecting more accurate predictions overall. Additionally, the variability of errors in the XEmbedding model remains smaller, particularly at higher sparsity levels, whereas the one-hot model shows wider interquartile ranges and larger outliers, highlighting its reduced robustness as sparsity increases. These results suggest that the learned embeddings provide a more stable and reliable representation of categorical features than one-hot encoding, especially when the data becomes sparse.

As an example, we consider the results of the issuer AMGN (Amgen Inc. is a biopharmaceutical company) shown in Figure 9. It illustrates a concrete issuer-level comparison between XEmbedding and one-hot representations for AMGN
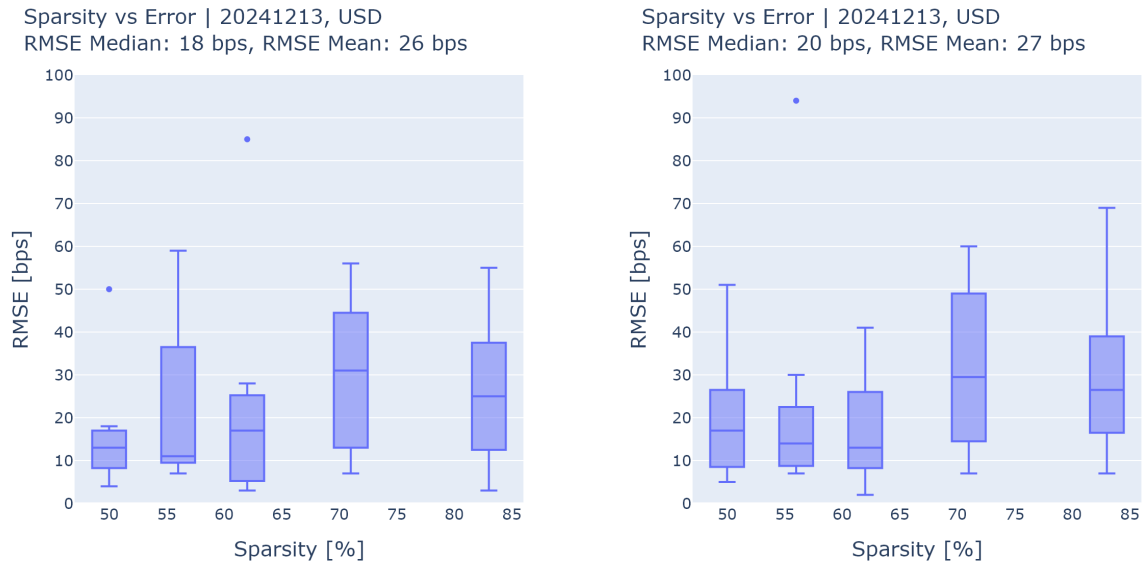
Figure 8: Comparison of overall performance of our model (the "XEmbedding" model) and a baseline that encodes categorical features using one-hot representations (the "one-hot" model).
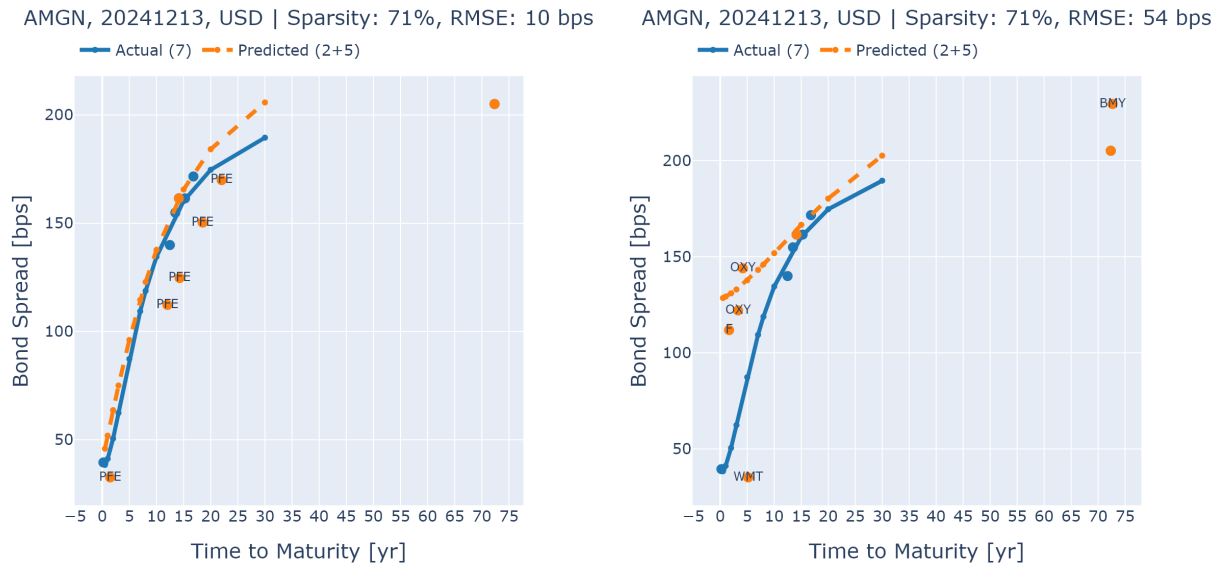


Figure 9: Comparison of XEmbedding (left) and one-hot (right) model prediction versus actual CDS spread curves for AMGN (Amgen Inc.). Annotated orange markers indicate specific predicted bonds, such as PFE (Pfizer Inc.), which belongs to a similar industry group, and OXY (Occidental Petroleum Corp.), which does not.

under 71% sparsity, where only 2 query bonds are retained and augmented with 5 similar neighbors. The XEmbedding model (left) produces a predicted curve that closely tracks the actual term structure, achieving an RMSE of 10 bps with smooth alignment across maturities and minimal deviation even at the long end. In stark contrast, the one-hot baseline (right) selects less relevant peers (including Occidental Petroleum Corp. (OXY), Ford Motor Co. (F), and Walmart Inc. (WMT)), resulting in a distorted upward bias and much higher RMSE of 54 bps, as evidenced by the wider spread between predicted and actual points. This example underscores how XEmbedding's nuanced categorical similarities yield more accurate spread reconstructions than rigid exact-match logic, particularly when data is scarce. In fact, the superior performance stems from XEmbedding's ability to quantify partial similarity between non-identical categories,

reducing the number of poor matches in high-sparsity regimes. We provide two more examples including KOREA and WFC in Appendix B.

## 7 Conclusion

This study presented an embedding-based framework that prioritizes categorical, non-financial bond attributes as the primary drivers of bond similarity. Empirically, embedding representations of these categorical variables outperform one-hot and other baselines in reconstructing spread curves from sparse issuer catalogs, suggesting that how categorical information is encoded matters more than adding further numerical detail. Future work could combine embeddings with supervised similarity learners such as random forest proximities to build hybrid models that exploit both high-quality categorical representations and detailed financial dynamics. In addition, future extensions using multimodal or graph-based embeddings may uncover deeper structural relationships among issuers.

## Acknowledgments

## References

Bengio, Y., Courville, A., & Vincent, P. (2013). Representation learning: A review and new perspectives. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, *35*(8), 1798–1828. https://doi.org/10.1109/TPAMI.2013.50

Breiman, L. (2004). Consistency for a simple model of random forests. https://api.semanticscholar.org/CorpusID:123042984

Busemeyer, J. R., & Bruza, P. D. (2012). *Quantum models of cognition and decision*. Cambridge University Press.

Desai, D., & Mehta, D. (2021). On robustness of mutual funds categorization and distance metric learning. *The Journal of Financial Data Science*, 130–150. http://jmlr.org/papers/v17/14-168.html

Devlin, J., Chang, M.-W., Lee, K., & Toutanova, K. (2019, June). BERT: Pre-training of deep bidirectional transformers for language understanding. In J. Burstein, C. Doran, & T. Solorio (Eds.), *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)* (pp. 4171–4186). Association for Computational Linguistics. https://doi.org/10.18653/v1/N19-1423

Goodfellow, I., Bengio, Y., & Courville, A. (2016). *Deep learning* [http://www.deeplearningbook.org]. MIT Press.

Gururangan, S., Marasović, A., Swayamdipta, S., Lo, K., Beltagy, I., Downey, D., & Smith, N. A. (2020). Don't stop pretraining: Adapt language models to domains and tasks. *ArXiv*, *abs/2004.10964*. https://api.semanticscholar.org/CorpusID:216080466

Haeri, A., Vitrano, J., & Ghelichi, M. (2025). Generative ai enhanced financial risk management information retrieval. https://arxiv.org/abs/2504.06293

Heaton, J. B., Polson, N. G., & Witte, J. H. (2017). Deep learning for finance: Deep portfolios. *Applied Stochastic Models in Business and Industry*, *33*(1), 3–12. https://doi.org/https://doi.org/10.1002/asmb.2209

Jeyapaulraj, J., Desai, D., Chu, P., Mehta, D., Pasquali, S., & Sommer, P. (2022). Supervised similarity learning for corporate bonds using random forest proximities. https://arxiv.org/abs/2207.04368

Kulis, B. (2013). Metric learning : A survey by. https://api.semanticscholar.org/CorpusID:262315341

Li, H. (2014). *Learning to rank for information retrieval and natural language processing, second edition*. Morgan & Claypool Publishers. https://doi.org/10.2200/S00607ED2V01Y201410HLT026

Mentch, L., & Hooker, G. (2016). Quantifying uncertainty in random forests via confidence intervals and hypothesis tests. *Journal of Machine Learning Research*, *17*(26), 1–41. http://jmlr.org/papers/v17/14-168.html

Mikolov, T., Sutskever, I., Chen, K., Corrado, G. S., & Dean, J. (2013). Distributed representations of words and phrases and their compositionality. In C. Burges, L. Bottou, M. Welling, Z. Ghahramani, & K. Weinberger (Eds.), *Advances in neural information processing systems* (Vol. 26). Curran Associates, Inc. https://proceedings.neurips.cc/paper_files/paper/2013/file/9aa42b31882ec039965f3c4923ce901b-Paper.pdf

Peters, M. E., Neumann, M., Iyyer, M., Gardner, M., Clark, C., Lee, K., & Zettlemoyer, L. (2018). Deep contextualized word representations. https://arxiv.org/abs/1802.05365

Reimers, N., & Gurevych, I. (2019). Sentence-bert: Sentence embeddings using siamese bert-networks. *Conference on Empirical Methods in Natural Language Processing*. https://api.semanticscholar.org/CorpusID:201646309

Rosaler, J., Candelori, L., Kirakosyan, V., Musaelian, K., Samson, R., Wells, M. T., Mehta, D., & Pasquali, S. (2025). Supervised similarity for high-yield corporate bonds with quantum cognition machine learning. https://arxiv.org/abs/2502.01495

Sanh, V., Debut, L., Chaumond, J., & Wolf, T. (2020). Distilbert, a distilled version of bert: Smaller, faster, cheaper and lighter. https://arxiv.org/abs/1910.01108

Scornet, E., Biau, G., & Vert, J.-P. (2015). Consistency of random forests. *The Annals of Statistics*, *43*(4). https://doi.org/10.1214/15-aos1321

Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L., & Polosukhin, I. (2023). Attention is all you need. https://arxiv.org/abs/1706.03762

Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., & Polosukhin, I. (2017). Attention is all you need. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, & R. Garnett (Eds.), *Advances in neural information processing systems* (Vol. 30). Curran Associates, Inc. https://proceedings.neurips.cc/paper_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf

# A   Appendix: Additional Similarity Search Results

This appendix provides additional qualitative evidence for the similarity search framework by showing how learned embeddings structure the bond universe across issuers, sectors, and regions. The examples illustrate how similarity relationships emerge, evolve with weaker issuer- and sector-level alignment, and balance fine-grained specificity with broader economic structure, demonstrating robust neighbor selection and curve reconstruction across diverse credit profiles.

Figure A.1 presents an aerospace and defense case study, showing how the model retrieves economically coherent peers for Boeing Co. and how similarity scores decay as sector and issuer alignment weaken. It reports the top-ranked neighbors for a Boeing Co. bond, illustrating how the embedding model prioritizes issuers in the same aerospace and defense sector and domicile while smoothly degrading similarity scores as sectoral or issuer alignment weakens. The progression from Boeing to Lockheed Martin, Northrop Grumman, RTX, and related aerospace equipment names demonstrates that the learned representations preserve both issuer-level and industry-level structure, rather than relying on exact categorical matches alone.

Figure A.2 examines an automotive and consumer-linked case, demonstrating how the model ranks General Motors neighbors along an intuitive auto retail transportation continuum. It shows the nearest neighbors for a General Motors bond, where the model first selects U.S. auto manufacturers such as Ford and related financing entities before gradually expanding to adjacent retail and transportation names. This behavior highlights the model's ability to capture second-order similarity within a broader automotive and consumer ecosystem, yielding a graded similarity profile instead of a hard boundary between "in-sector" and "out-of-sector" bonds.

Figure A.3 focuses on a large-cap healthcare issuer, highlighting how the embeddings capture layered similarity within pharmaceuticals, biotechnology, and broader healthcare services. The query bond from Johnson&Johnson is matched primarily to large U.S. pharmaceutical and healthcare issuers, including Pfizer, Merck, Bristol Myers Squibb, and Eli Lilly. The ranked list illustrates that the embedding space organizes issuers along therapeutics and healthcare-service dimensions, with similarity scores declining smoothly as the model transitions from core pharma peers to broader healthcare and insurance names.

Figure A.4 analyzes a sub-sovereign government issuer, illustrating how the model prioritizes Canadian provincial peers before gradually expanding to pipelines, corporates, and selected sovereigns. It presents similarity search results for a Province of Manitoba bond, where the top neighbors are Canadian provincial issuers such as Ontario, Quebec, Alberta, Saskatchewan, and British Columbia. Only at lower similarity scores do non-provincial Canadian entities and selected foreign sovereigns appear, showing that the embeddings strongly prioritize regional, sectoral, and domicile alignment for sub-sovereign government bonds.

Figure A.5 reports CDS spread curve reconstructions for the aforementioned industrial, consumer, healthcare, and government issuers, confirming the stability of the augmentation framework across heterogeneous credit profiles. It compares actual and predicted CDS spread curves for Boeing, General Motors, Johnson&Johnson, and the Province of Manitoba under varying sparsity levels, illustrating how reconstruction quality evolves as the available market information becomes increasingly limited. Across all four issuers, the augmented curves closely track the overall level, slope, and curvature of the observed term structures, yielding consistently low RMSEs in the range of 5–18 bps. These results demonstrate that the proposed similarity-based augmentation remains stable and effective across different sectors, even under pronounced sparsity, and that the learned representations support accurate curve reconstruction in diverse market settings.
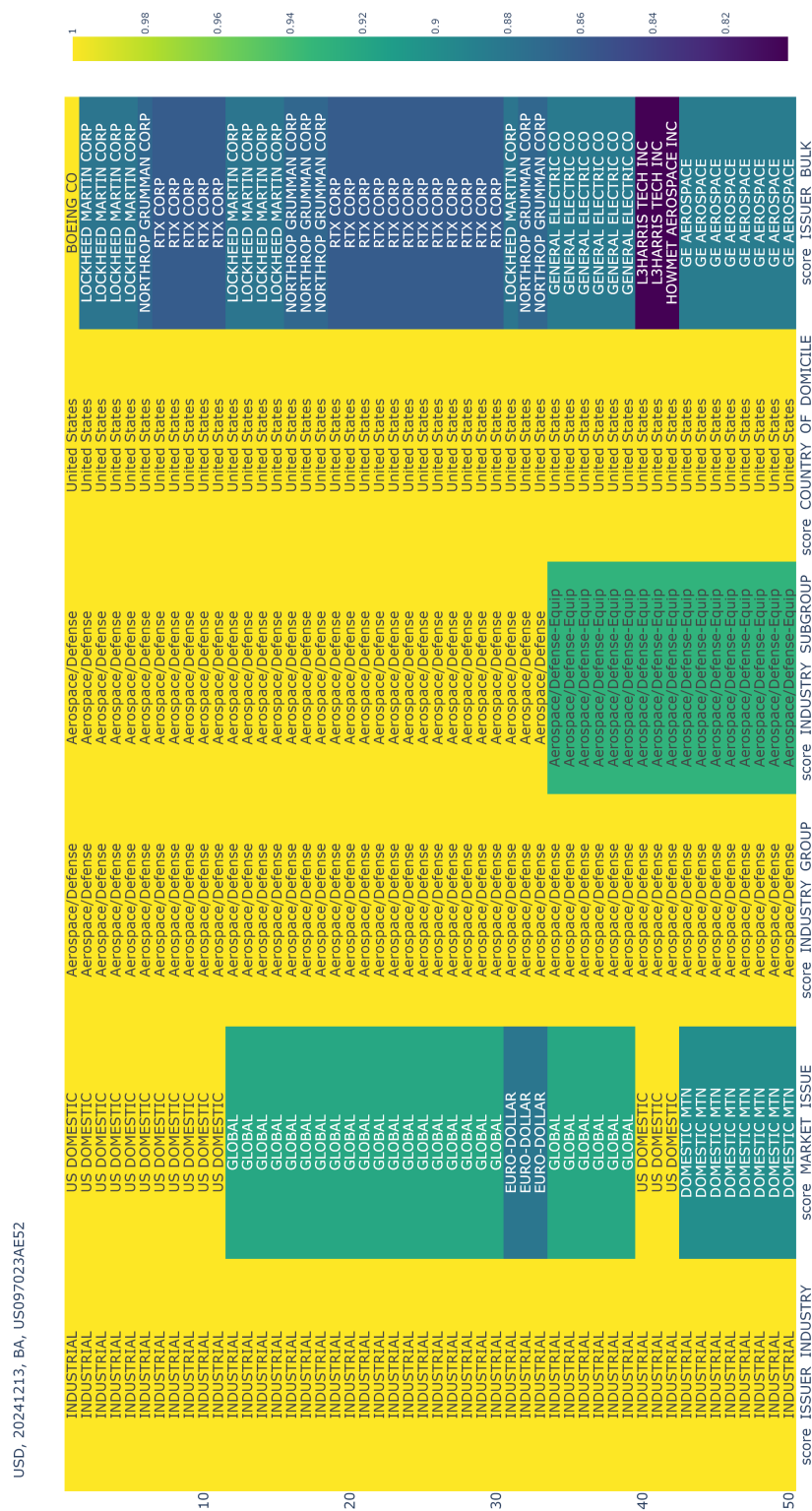
Figure A.1: Visualization of bond similarity search results for the query bond of Boeing Co.
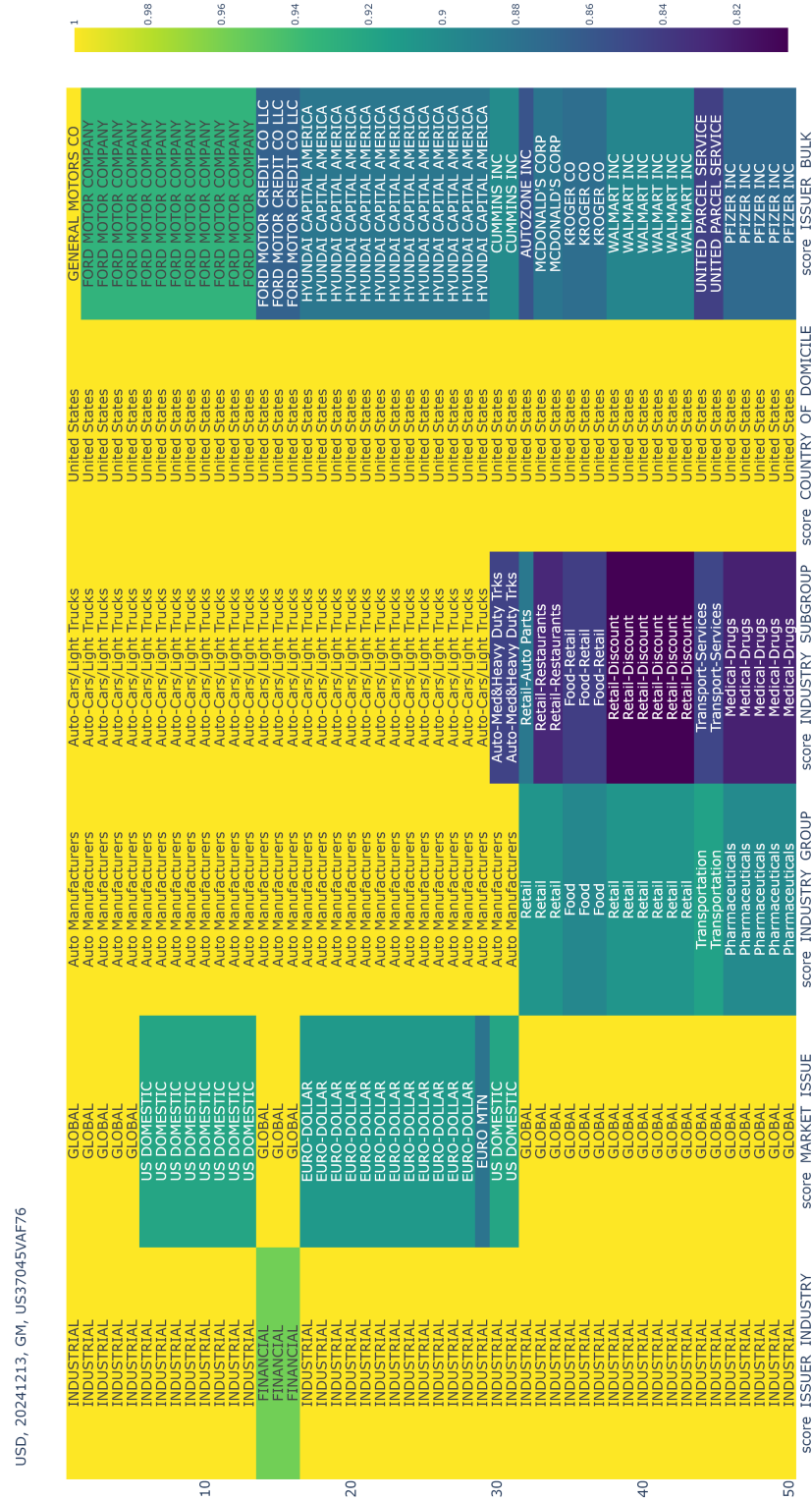
# Bond Embedding Similarity



Figure A.2: Visualization of bond similarity search results for the query bond of General Motors.
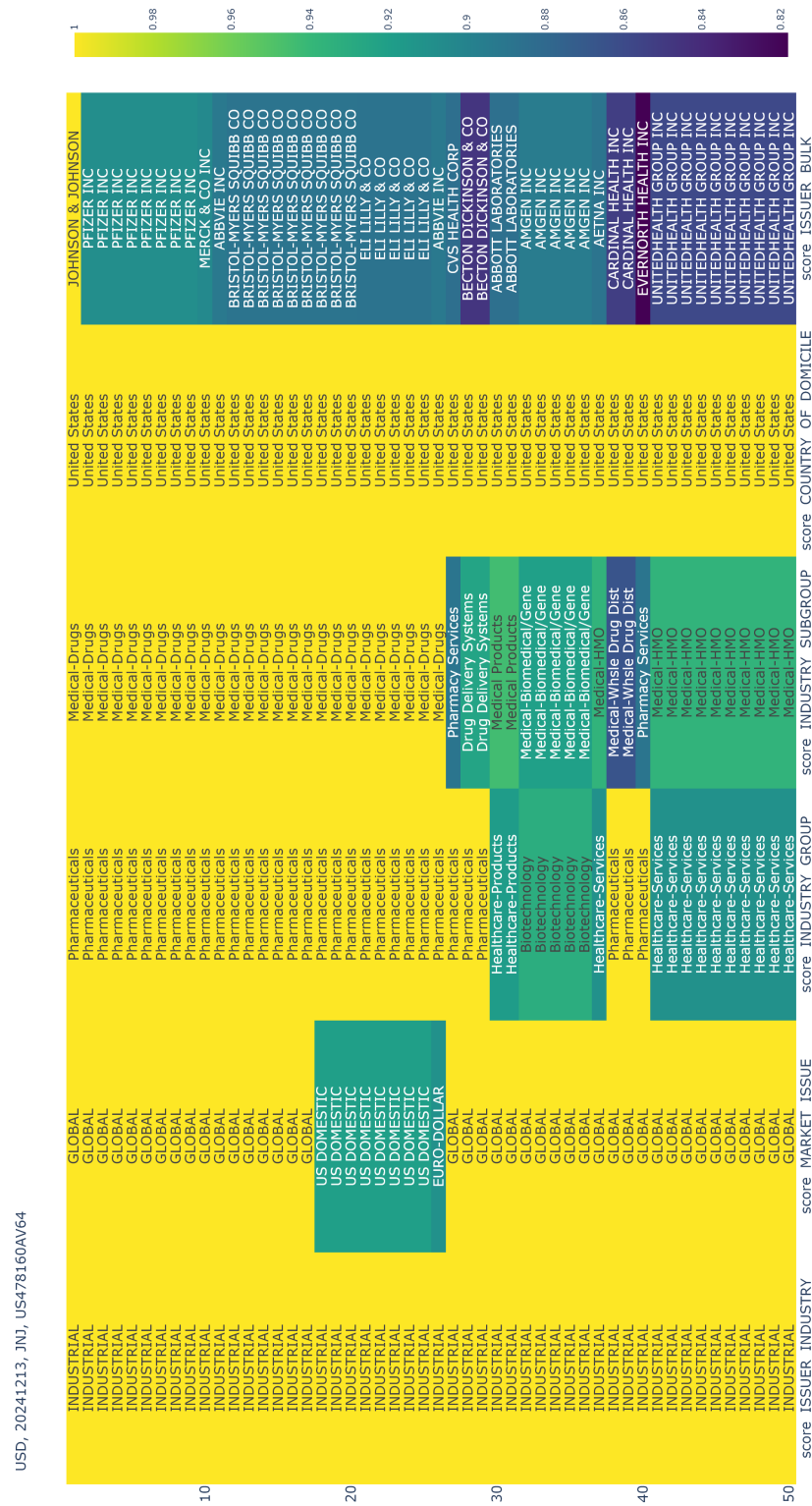
Bond Embedding Similarity



Figure A.3: Visualization of bond similarity search results for the query bond of Johnson&Johnson.
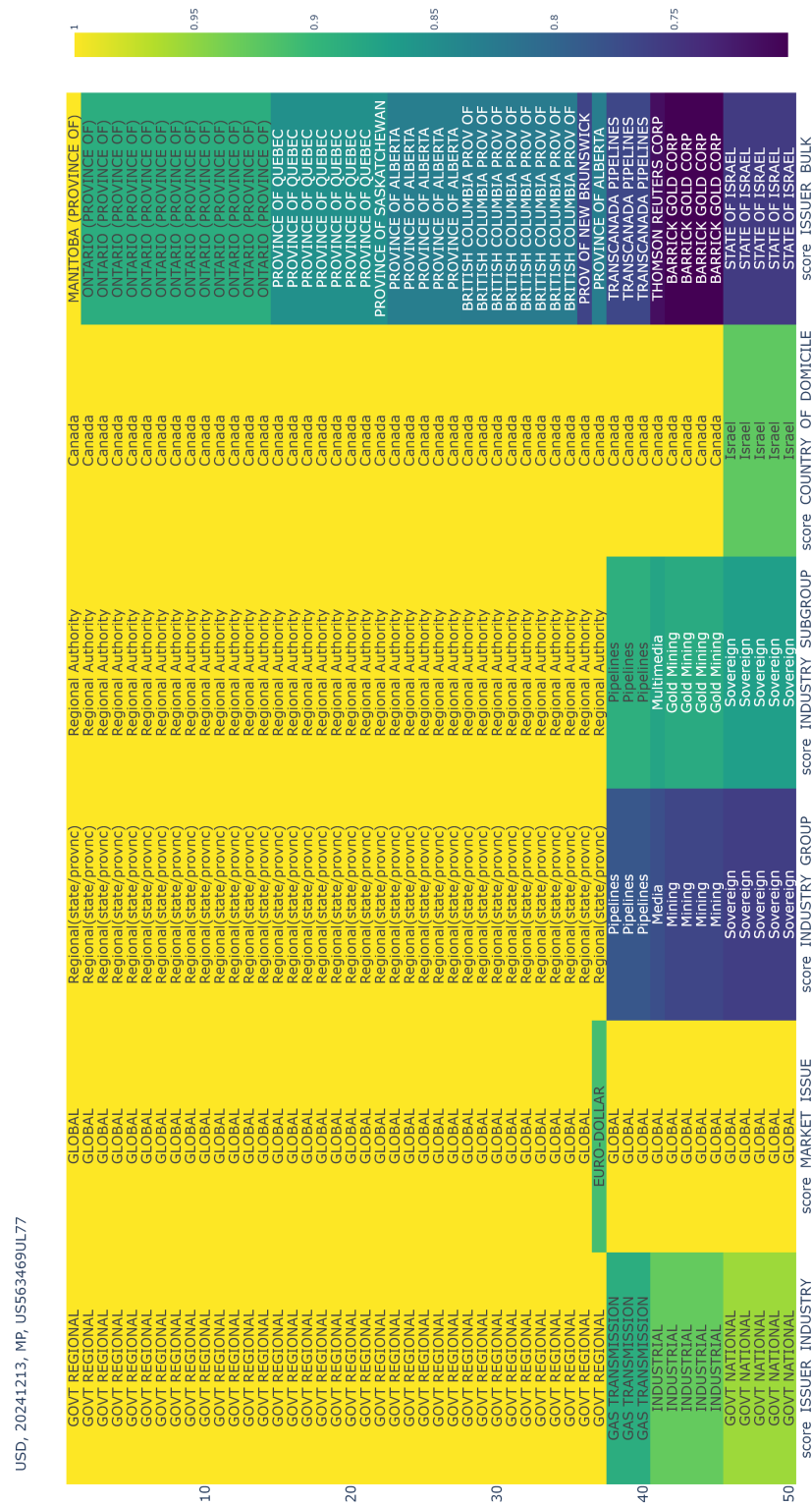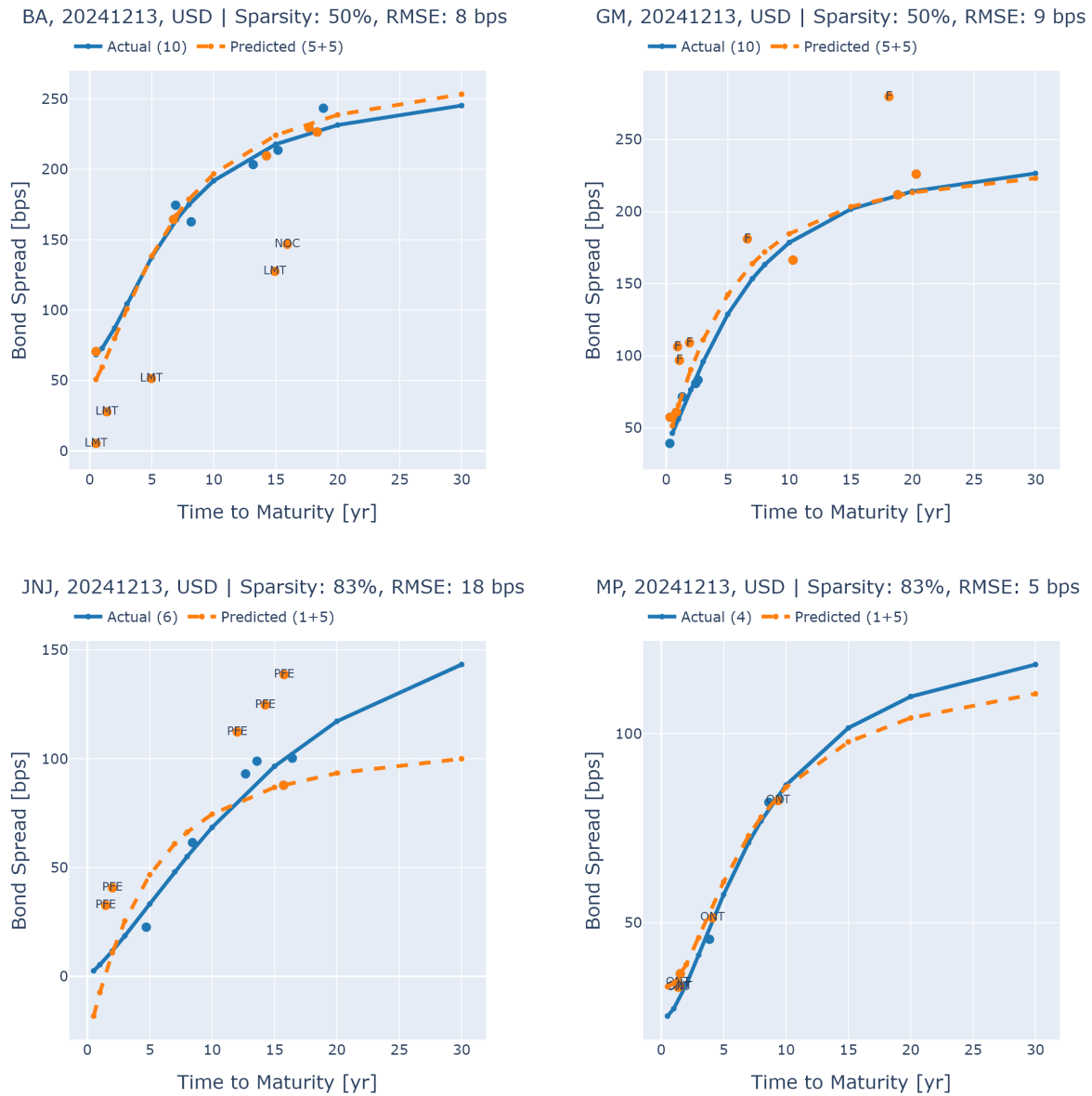
# Bond Embedding Similarity



Figure A.4: Visualization of bond similarity search results for the query bond of Province of Manitoba.

Figure A.5: Comparison of predicted and actual CDS spread curves for Boeing Co. (top-left), General Motors (top-right), Johnson&Johnson (bottom-left), and Province of Manitoba (bottom-right) issuers.

# B   Appendix: Additional Benchmarking Results

This appendix reports extended benchmarking experiments that compare the XEmbedding model against alternative encodings and two-step procedures, with a focus on sparsity robustness and issuer-level behavior.

Figure B.6 illustrates four alternative strategies for retrieving similar bonds. The generic approach relies on rule-based filtering using categorical attributes (e.g., industry or rating) to identify comparable instruments. In contrast, the numerical approach represents bonds using numerical features and retrieves neighbors directly via similarity search. To leverage the strengths of both paradigms, we introduce two two-step approaches: either categorical information is first embedded and used to construct an initial candidate set that is subsequently refined through numerical similarity search, or numerical similarity is applied first to form a reduced catalog that is then refined using categorical embeddings.

Figure B.7 reports the performance of alternative baseline approaches across different sparsity levels. As sparsity increases, both RMSE and MAPE generally deteriorate and exhibit larger dispersion, indicating reduced robustness under limited information. In contrast to these baselines, XEmbedding consistently achieves lower error across all sparsity regimes (cf. Figure 6), with tighter error distributions and improved stability. These results suggest that explicitly learning structured representations enables XEmbedding to better exploit inputs and maintain predictive accuracy relatively.

Figures B.8 and B.9 provide issuer-level comparisons for KOREA and WFC (Wells Fargo & Co), illustrating how XEmbedding avoids implausible neighbors admitted by one-hot encodings and thereby achieves markedly lower RMSE and more realistic curve shapes. For both issuers, the embedding-based model selects economically plausible neighbors and delivers lower RMSEs (e.g., 6 bps vs. 60 bps), whereas the one-hot representation often admits mismatched peers from unrelated regions or sectors, leading to distorted curve shapes and inflated errors. In the case of KOREA, XEmbedding retrieves regionally aligned issuers such as PHILIP, while excluding less relevant entities like BHRAIN; similarly, for WFC, the model prioritizes country-aligned peers such as BAC (Bank of America Corp.) over geographically distant institutions including NAB (National Australia Bank Ltd) and ANZ (ANZ Group Holdings Ltd).
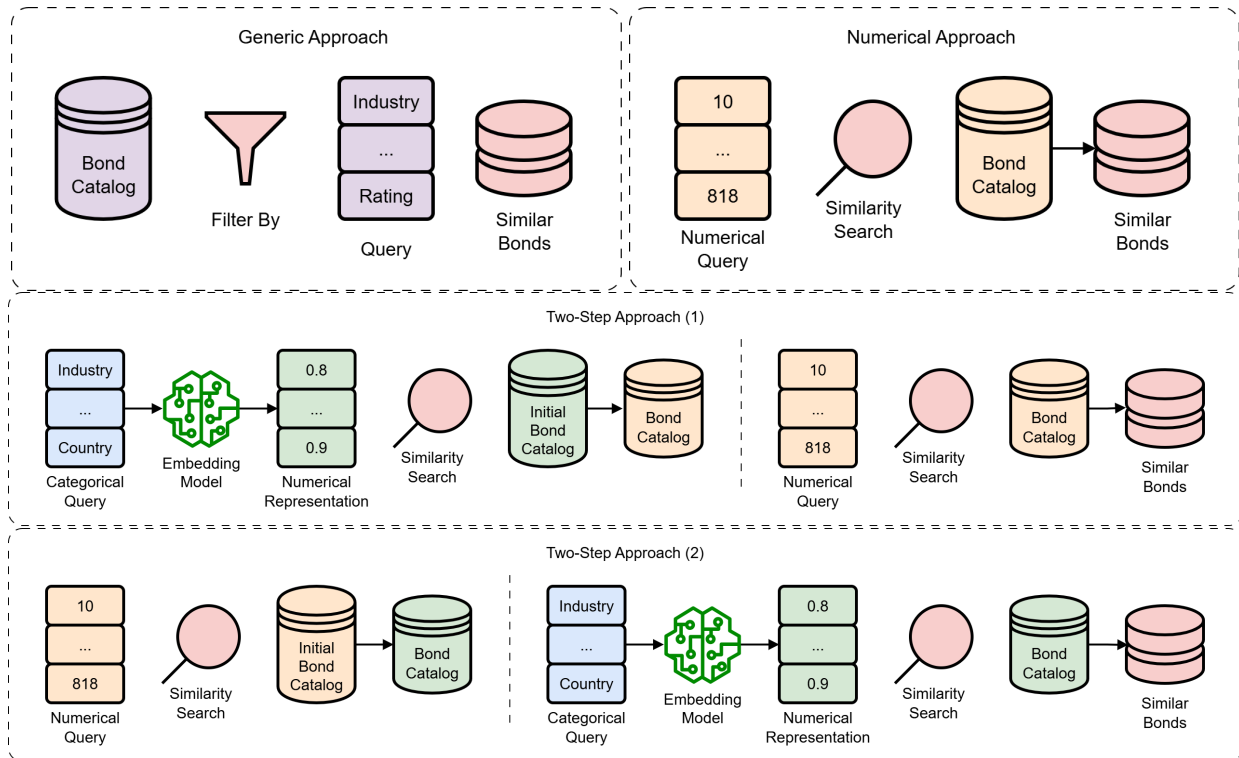


Figure B.6: Comparison of alternative approaches. The generic approach (top-left) filters the bond catalog using categorical attributes such as industry and rating. The numerical approach (top-right) directly performs similarity search over numerical bond representations. The two-step approaches (bottom) combine categorical and numerical information by first narrowing the search space using one modality and then applying similarity search using the other.
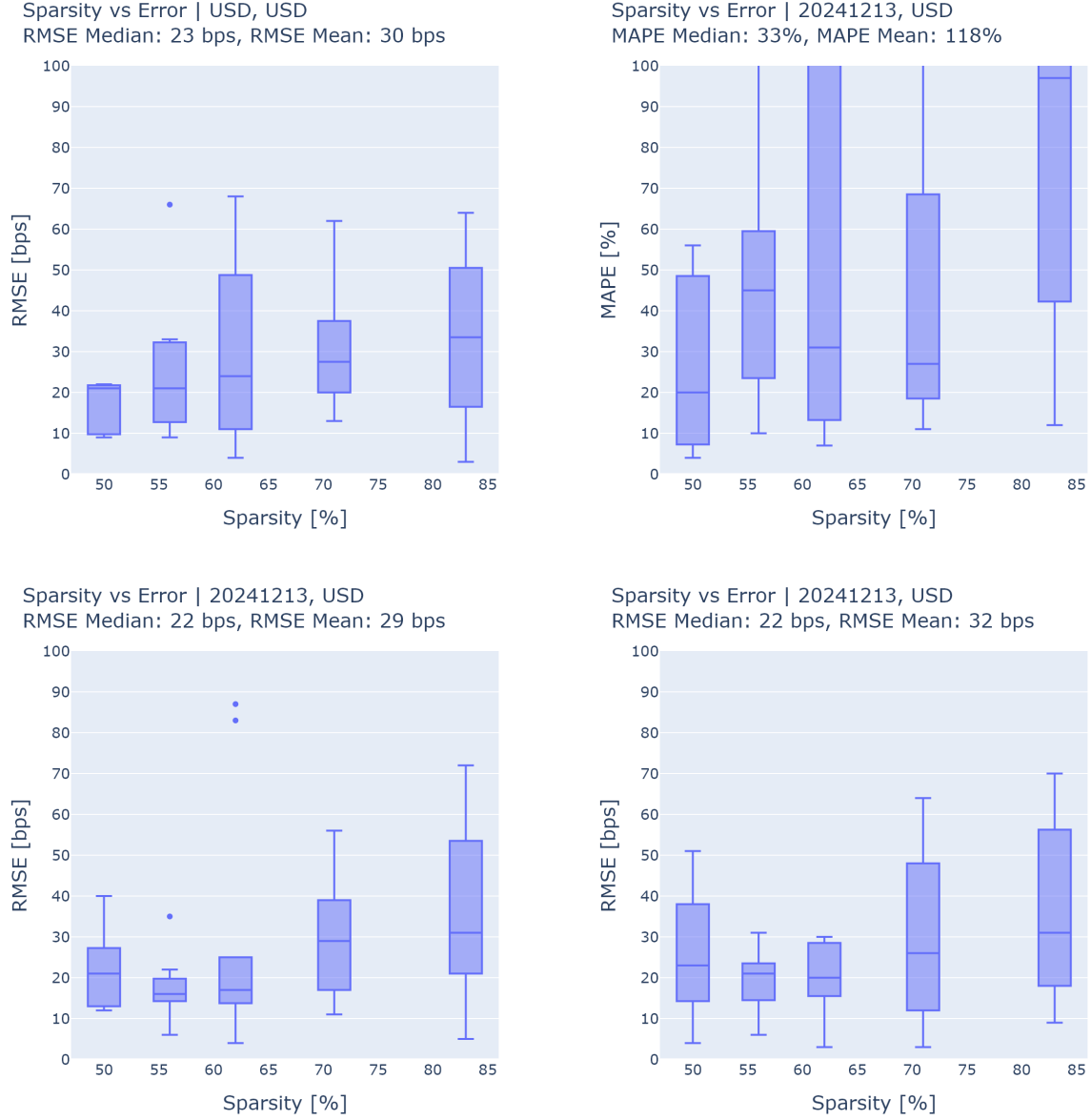
Figure B.7: Comparison of overall performance of alternative approaches (top-left: Generic, top-right: Numerical, bottom-left: Two-Step (1), bottom:right: Two-Step (2)). Error as a function of sparsity for alternative baseline approaches. Boxplots summarize the distribution of prediction errors across sparsity levels, reported using RMSE (bps) and MAPE (%). Median and mean error statistics are shown in each panel. Across all sparsity regimes, the baseline methods exhibit higher error levels and greater variability compared to the proposed XEmbedding model.
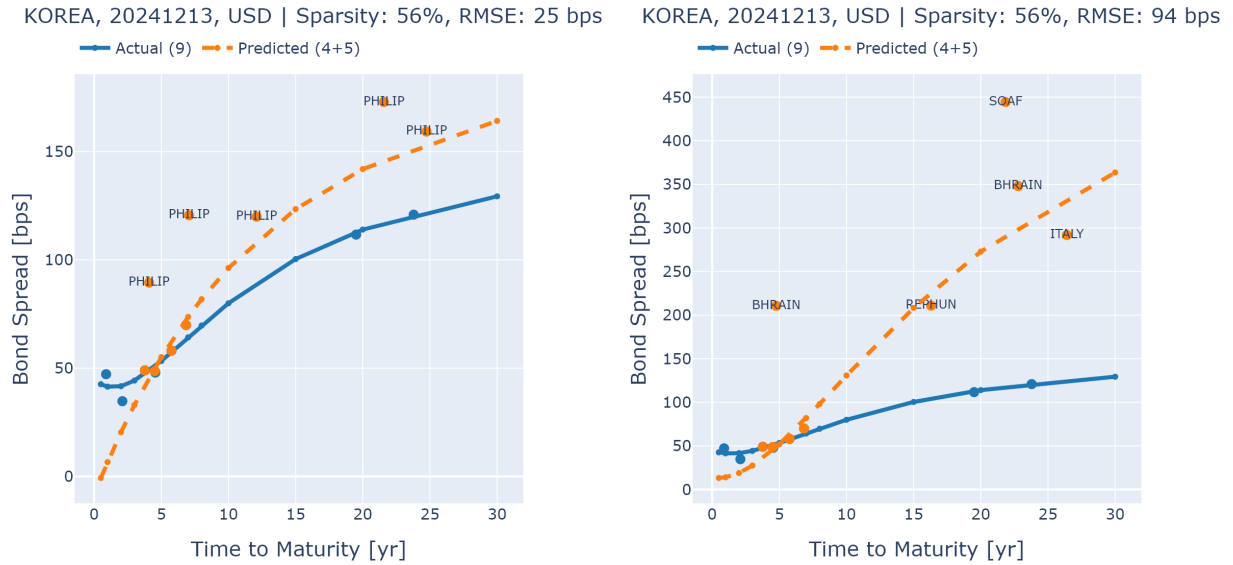
Figure B.8: Comparison of XEmbedding (left) and one-hot (right) model prediction versus actual CDS spread curves for KOREA. Annotated orange markers indicate specific predicted bonds, such as PHILIP, which belongs to a similar region, and BHRAIN, which does not.
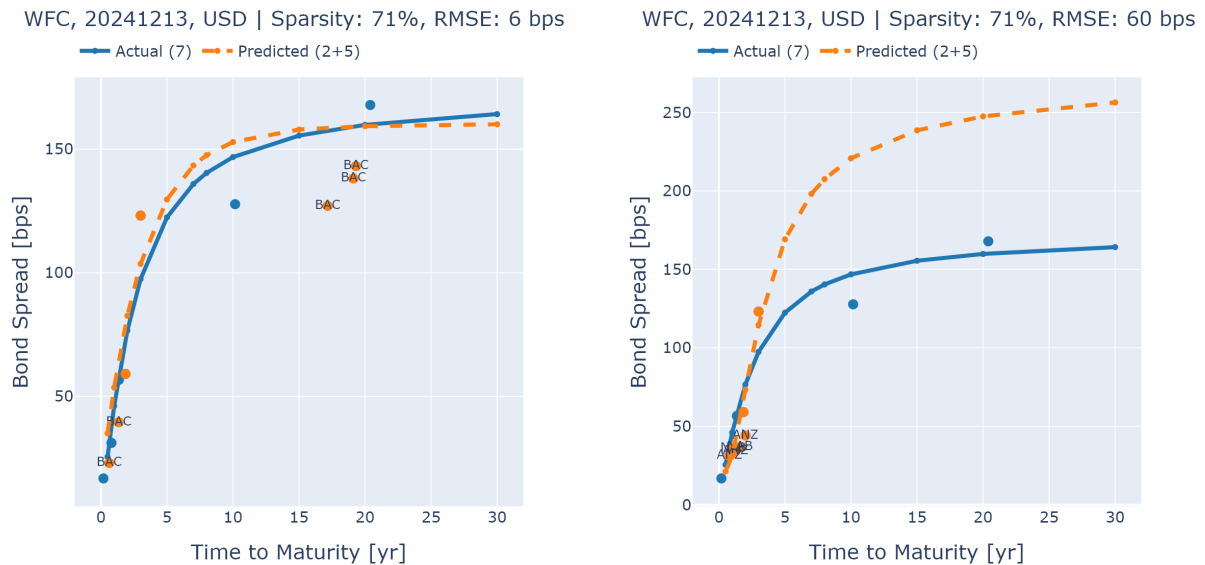


Figure B.9: Comparison of XEmbedding (left) and one-hot (right) model prediction versus actual CDS spread curves for WFC (Wells Fargo & Co). Annotated orange markers indicate specific predicted bonds, such as BAC (Bank of America Corp.), which belongs to a similar country, and NAB (National Australia Bank Ltd) and ANZ (ANZ Group Holdings Ltd), which do not.