



Paper Type: Original Article

Semantic-Enhanced Demand Forecasting: A Multimodal Transformer Integrating Product Descriptions and Customer Purchase History

Fatemeh Zare Baghiabad* 

Department of Industrial Engineering, Faculty of Engineering, Ardakan University, P.O. Box 184, Ardakan, Yazd, Iran; f.zare@ardakan.ac.ir.

Citation:

Received: 17 January 2024	Zare Baghiabad, F. (2025). Semantic-enhanced demand forecasting: A multimodal transformer integrating product descriptions and customer purchase history. <i>Annals of Optimization With Applications</i> , 1(4), 221-233.
Revised: 13 May 2024	
Accepted: 26 August 2024	

Abstract


Accurate forecasting of daily demand at the individual customer-product level remains a critical yet challenging problem in retail, hindered by data sparsity, volatile consumer behavior, and the underutilization of unstructured product information. This study addresses this gap by proposing a novel Multimodal Semantic Transformer (MST) framework that integrates semantic product embeddings, derived from Natural Language Processing (NLP) of descriptions, with structured customer purchase history and multi-scale temporal features. Using the Online Retail II dataset, the model was evaluated against benchmarks including LSTM, Gradient Boosting Machines, and a unimodal Transformer. The results demonstrate that the MST framework significantly outperforms all benchmarks, achieving a 15.5% reduction in Mean Squared Error (MSE) compared to the best baseline. Key findings confirm that semantic fusion provides a crucial signal for sparse products and that temporal embeddings with dynamic attention conditioning are essential for modeling complex seasonality and context. The study concludes that deep multimodal integration is a transformative approach for granular demand forecasting, offering a scalable and interpretable solution to enable hyper-personalized inventory management and more resilient, customer-centric retail supply chains.


Keywords: Multimodal transformer, Demand forecasting, Semantic integration, Natural language processing embeddings, Retail analytics.

1 | Introduction

Accurate demand forecasting is a critical component of modern retail and supply chain management, directly impacting inventory optimization, resource allocation, and strategic planning [1]. However, conventional

 Corresponding Author: f.zare@ardakan.ac.ir

 <https://doi.org/10.48314/anowa.v1i4.58>

 Licensee System Analytics. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<http://creativecommons.org/licenses/by/4.0>).

models often struggle with the inherent volatility and multidimensional complexity of consumer behavior, particularly at the granular customer-product level. While transactional time-series data provides a historical baseline, it largely ignores the rich semantic information embedded in product attributes and descriptions, which are known to significantly influence purchasing decisions. For instance, the persistent challenge of demand-supply mismatches continues to drive significant operational inefficiencies and financial losses across the global retail sector, underscoring the pressing need for more advanced, integrated forecasting solutions [2]. Recent advances in deep learning, specifically Transformer architectures, have shown remarkable success in capturing complex temporal dependencies [3]. 3

However, their application in retail has predominantly been unimodal, focusing on structured transactional data while overlooking the predictive potential within unstructured textual data [4]. This research bridges this critical gap by proposing a semantic-enhanced demand forecasting framework. We introduce a novel multimodal Transformer model that synergistically integrates product descriptions processed via Natural Language Processing (NLP) to extract semantic embeddings with detailed customer purchase history to achieve a more holistic and accurate prediction of future demand. By unifying these heterogeneous data streams, this study aims to move beyond traditional time-series analysis and contribute to the emerging field of multimodal artificial intelligence in operations management, offering a scalable solution to enhance forecast precision, reduce operational costs, and enable truly customer-centric supply chains.

Forecasting total daily demand at the product-customer level is defined as the task of predicting the exact quantity of a specific product that an individual customer will request on a given future day. This granular metric is a cornerstone of modern, data-driven retail and supply chain management, enabling hyper-personalized strategies in inventory fulfillment, dynamic pricing, and targeted promotion [5]. Its significance lies in its direct operational impact; accurate forecasts at this micro-segmentation level can drastically reduce stockouts and overstock costs while enhancing customer satisfaction through improved product availability. However, this forecasting task presents unique challenges due to its key dimensions: high sparsity (numerous zero-demand instances), extreme volatility driven by individual decision-making, and complex interactions between customer-specific preferences, product attributes, and temporal factors like day-of-week effects [6]. Traditional statistical models, such as Autoregressive Integrated Moving Average (ARIMA) frameworks, often struggle with these characteristics, particularly in large-scale retail environments, as noted in systematic reviews of the field [7].

In response, recent literature has pivoted towards machine and deep learning approaches. For instance, ensemble-based predictive analytics have demonstrated efficacy in capturing non-linear relationships from high-dimensional transactional and customer data for multi-channel retailing [8]. Concurrently, modified Transformer models have been employed to model sequential purchase patterns and temporal dependencies with higher accuracy and efficiency [9]. A prevailing theoretical perspective posits that integrating heterogeneous data sources particularly unstructured data like product descriptions is crucial for advancing beyond transactional history. Recent studies have begun exploring multimodal fusion models to this end, such as those integrating text, time series, and imagery for causal-aware forecasting, highlighting the potential of semantic product information to explain and predict customer-specific demand, though this area remains underexplored [10].

In the context of data-driven demand forecasting, product descriptions refer to the unstructured textual data accompanying items in a catalog or sales system, which detail attributes, features, and intended use. In advanced analytical frameworks, these texts are processed through NLP to create semantic embeddings dense, low-dimensional vector representations that capture the underlying meaning, context, and relational properties of the words. This transformation is significant because it bridges the semantic gap between human-readable product information and machine-interpretable data, enabling models to understand relational similarities between items. Key dimensions of this process include the choice of embedding model, the level of semantic granularity, and the integration mechanism into a downstream predictive model. Previous studies underscore the value of leveraging such textual data in multimodal frameworks. For instance,

research on spatial feature fusion for e-commerce commodity demand forecasting has demonstrated the utility of multimodal data guidance [11]. More recently, the theoretical perspective has shifted towards deep learning and multimodal integration, with studies showing that transformer-based fusion models can significantly enhance forecast accuracy by providing a "semantic signature," particularly for new products [9]. Similarly, research aimed at improving sales forecasting accuracy through tensor factorization approaches incorporates a form of demand awareness that can benefit from semantic signals [12]. These studies collectively highlight a major finding: incorporating semantic product information provides a robust, generalizable signal that complements quantitative transactional data, particularly for modeling consumer choice drivers and generalizing to new or infrequently purchased items. The integration of these embeddings into time-series forecasting models, however, remains an emerging and highly promising frontier.

The foundational input for modern demand forecasting models is often formalized as a multivariate time series matrix, $X = [x_1, x_2, \dots, x_T] \in \mathbb{R}^{d \times T}$, representing T sequential observations where each vector x_t is a d -dimensional set of engineered features derived from raw transactional data. This structured representation is significant as it transforms granular purchase records into a format suitable for machine learning, explicitly encoding temporal dependencies, historical patterns, and contextual drivers of demand. The key dimensions of x_t are meticulously crafted based on established econometric and machine learning principles. Lagged demand variables (e.g., sales from previous days) are crucial for capturing autoregressive effects and short-term momentum. Moving statistics, such as rolling averages, help smooth noise and highlight underlying trends. Time features including day-of-week, month, and holiday indicators are essential for modeling deterministic seasonality and calendar effects.

Finally, structured external variables like unit price, customer identifier, and country code (after appropriate encoding) integrate critical exogenous information that explains demand shifts beyond pure historical patterns. Previous studies have consistently demonstrated the importance of this feature engineering stage for building effective models, whether tree-based, sequential, or hybrid in nature [13]. The theoretical perspective underpinning this practice is that demand is not merely a self-contained series but a function of its own past, predictable cycles, and measurable external conditions; therefore, a well-constructed X matrix is a necessary precondition for any data-driven model to learn the true demand-generating process, a principle that holds true even when adapting models to new retail domains [14].

To effectively model complex temporal dynamics and contextual effects in modern transformer-based demand forecasting, two advanced architectural components are critical: multi-scale temporal embeddings (E) and a dynamic mask matrix (M). Multi-scale temporal embeddings are learned vector representations that encode time at different granularities (e.g., weekly, monthly) directly into the model's input space, allowing it to inherently capture seasonal and periodic patterns without relying on explicit and often brittle pre-processing steps like classical time series decomposition. Their significance lies in enabling the model to disentangle and attend to patterns operating on different cycles simultaneously, a key dimension being the design of the embedding functions for various periods.

Concurrently, the dynamic mask matrix (M) is a mechanism designed to modulate the self-attention weights by directly injecting the influence of specific, time-varying external variables such as real-time price (P_t), customer profile (C_t), and geographic region (G_t). This is defined as $M = f(P_t, C_t, G_t)$, where f is a parameterized function. The significance of M is its ability to condition the model's focus on historical data based on current exogenous contexts, making the attention mechanism context-aware and more interpretable. A key theoretical perspective is that pure self-attention, while powerful for capturing dependencies, is "context-agnostic" regarding external drivers; the dynamic mask provides a principled way to introduce this conditioning. Previous studies have validated the utility of both concepts in advanced forecasting architectures. The integration of multimodal temporal fusion, as shown in recent transformers, relies on effectively encoding time [4].

Furthermore, the need to incorporate behavioral and environmental frameworks into forecasting aligns with the purpose of a dynamic mask to inject contextual and customer-specific signals [14]. These studies

collectively highlight a major finding: the performance of transformer architectures in time-series forecasting is substantially enhanced not just by the attention mechanism itself, but by specialized components that explicitly encode inductive biases for time and context, bridging the gap between the model's general sequence-processing capability and the specific structured priors of the forecasting domain.

While existing research has established the efficacy of Transformers for time-series forecasting and underscored the value of multimodal data, a significant gap remains in semantically-aware, customer-centric demand prediction at the finest granularity. Specifically, prior work has not yet fully leveraged the rich, unstructured textual data of product descriptions to enhance the modeling of individual customer preferences within a unified, multi-modal Transformer architecture designed explicitly for the product-customer dyad. To address this, the present study posits three core research questions: (RQ1) How can semantic embeddings from product descriptions be effectively fused with structured transactional and customer data in a Transformer framework? (RQ2) Does this multi-modal integration lead to a statistically significant improvement in forecasting total daily demand at the product-customer level compared to unimodal benchmarks? (RQ3) How do advanced components like multi-scale temporal embeddings and a dynamic mask matrix contribute to the model's performance and interpretability in this context?

Consequently, the objectives of this paper are: 1) to propose a novel Multimodal Semantic Transformer (MST) framework that integrates NLP-based product semantics, customer purchase history, and engineered time-series features, 2) to rigorously evaluate its performance against state-of-the-art models using the Online Retail II dataset, and 3) to provide an interpretable analysis of how semantic and customer factors drive predictions. To this end, the paper is structured as follows. The Literature Review synthesizes relevant work in demand forecasting, multi-modal learning, and Transformer applications, solidifying the identified research gap. The Methodology section details the architecture of the MST, including data preprocessing, feature engineering, and the design of the fusion mechanism. The Experiments and Results section present the experimental setup, benchmark comparisons, and ablation studies to answer the core research questions. Finally, the Discussion and Conclusion section interprets the findings, discusses practical implications, acknowledges limitations, and suggests avenues for future research. The following section begins this exploration with a comprehensive review of the related literature.

2 | Literature Review

The pursuit of accurate demand forecasting represents a persistent and complex challenge in retail operations, one that is magnified by the scale, speed, and dimensionality of contemporary e-commerce ecosystems. Traditional statistical models, such as ARIMA, while providing a foundational framework, are increasingly recognized as inadequate for capturing the high-dimensionality, non-linearity, and inherent volatility that characterize modern retail data streams [7]. Comprehensive reviews of the field document a clear paradigm shift towards data-driven, machine learning-based methodologies, marking a decisive move away from purely autoregressive techniques toward models capable of learning complex, non-linear patterns directly from diverse and rich feature sets [6]. This evolution is driven by the acute operational and financial necessity to reduce both stockouts and overstock, two sides of the same costly coin that are exacerbated by inaccurate, aggregate-level predictions [2].

In response to these challenges, deep learning architectures have ascended to the forefront of forecasting research. Within this domain, Transformer models, built upon the self-attention mechanism, have recently established a new state-of-the-art for time-series tasks. Their ability to process entire sequences in parallel and model dependencies without the sequential constraints of recurrent networks like LSTMs provides a theoretically superior framework. Empirical studies specifically within retail contexts, such as the work by Oliveira and Ramos [3], have validated the effectiveness of time-series Transformers, setting new performance benchmarks. This success has catalyzed further innovation, leading to sophisticated hybrid models that combine Transformer strengths with other architectures to enhance accuracy, efficiency, and scalability, as

exemplified by the hybrid Temporal Convolutional Network (TCN) and Transformer model proposed by [13].

Parallel to architectural advancements, the strategic objective of forecasting has evolved from aggregate, product-level estimates to granular, customer-centric insights. Research on predicting Customer Lifetime Value (CLV), such as the high-performance system developed by Yan and Resnick [5], underscores the immense economic value of understanding and anticipating individual customer behavior. This aligns with the broader industrial trend towards hyper-personalization, where engineering tailored customer experiences is identified as a primary source of competitive advantage [15]. Consequently, forecasting at the nexus of a specific product and a specific customer—predicting not just what will sell, but to whom it will sell and when has become a critical capability for enabling personalized inventory management, dynamic pricing, and targeted promotions [14].

Acknowledging that predictive signal is distributed across diverse data types, the field has decisively moved towards multimodal fusion. The integration of heterogeneous data including text, time series, visual imagery, and contextual metadata is widely recognized as a powerful strategy for building more robust and generalizable forecasting models. Pioneering work by [11] demonstrated the value of guiding spatial feature fusion with multimodal data for e-commerce commodity forecasting. This paradigm has rapidly matured, with contemporary state-of-the-art frameworks explicitly architected as multimodal systems. For instance, [4] providing direct empirical validation for this core approach.

The cutting edge of multimodal research is progressing beyond simple feature concatenation towards intelligent, context-aware integration. Advanced models are beginning to incorporate principles of causal reasoning to distinguish mere correlations from true causative drivers of demand. Wang [10] exemplifies this direction with a "Causal-Aware Multimodal Transformer" that integrates text, time series, and satellite imagery for supply chain forecasting, aiming to isolate genuine demand signals. Similarly, sophisticated ensemble methods and tensor factorization approaches are being refined to enhance demand awareness and improve prediction accuracy across vast and complex product assortments, as seen in the work of [8] and [12].

The transition from theoretical model performance to practical, large-scale deployment introduces significant challenges. Key among these are the needs to efficiently scale predictions across millions of Stock-Keeping Units (SKUs) and to adapt models to new retail domains with differing data distributions. Furthermore, models must maintain robustness against sudden, structural market shifts such as those induced by global pandemics, the launch of new store formats, or the introduction of disruptive product categories which can fundamentally alter demand patterns and degrade the performance of models trained on pre-shift data [16]. Research by [14] directly addresses the latter through innovative domain adaptation techniques for retail demand prediction. Meanwhile, the design of inherently scalable architectures remains a central research thrust, a focus clearly evident in the development of efficient hybrid models [13].

Despite the remarkable progress in multimodal fusion, a critical examination of the literature reveals a relative underinvestment in leveraging the deep semantic content of product text. While modalities like images and structured metadata are increasingly integrated, the nuanced, descriptive language that shapes consumer perception and choice contained in product titles and descriptions has not been fully exploited within advanced Transformer frameworks for granular forecasting. Many current multimodal approaches treat text as a supplementary or auxiliary signal, rather than as a core, interpretable component for modeling the fundamental affinity between a customer's preferences and a product's attributes.

Synthesizing the trajectories outlined above, a distinct and consequential research gap emerges. While the literature firmly establishes 1) the superior capability of transformer architectures for sequence modeling, 2) the necessity of multimodal fusion for robust predictions, and 3) the strategic imperative of customer-centric, granular forecasting, no existing framework has been explicitly designed to deeply and interpretably integrate semantic product embeddings with individualized customer purchase histories within a unified, efficient transformer architecture for the specific task of daily product-customer demand forecasting. Previous contributions tend to focus on aggregate levels, incorporate text superficially, or do not address the unique

data sparsity and volatility challenges at the product-customer dyad level. It is this precise gap that the present study aims to fill by proposing the MST framework, thereby contributing directly to advancements in multimodal AI for operations, interpretable demand modeling, and personalized retail analytics.

3 | Method and Material

Research design and methodology

This study employs an experimental research design conducted within a controlled computational setting to develop and validate a novel multimodal deep learning framework. The methodology follows a quantitative, data-driven approach common in machine learning engineering, centered on model development, training, and comparative evaluation. The research design is guided by the theoretical principles of Transformer-based sequence modeling and multimodal data fusion. The general approach involves systematic data preprocessing, feature engineering, the architectural design of the proposed MST, and rigorous performance benchmarking against established forecasting models. Data analysis primarily involves statistical comparison of forecast accuracy metrics and ablation studies to isolate the contribution of key model components.

Statistical population

This study utilizes the publicly available "online retail II" dataset as its primary sample, which contains transactional data from a UK-based online retailer. The dataset comprises 541,909 individual transaction records spanning from December 2010 to December 2011, representing the complete population of non-cancelled orders within this period. Inclusion criteria were applied to filter the raw data, retaining only transactions with valid Stock codes, customer IDs, and positive quantities to ensure analytical integrity. The sample provides a rich set of characteristics critical to the study, including product descriptions, unit prices, purchase dates, customer identifiers, and country information, enabling the construction of granular product-customer time series. This dataset was chosen for its realism, scale, and public accessibility, providing a robust benchmark for evaluating customer-centric demand forecasting models.

Instruments and equipment

The research was conducted entirely in a computational environment, utilizing specialized software libraries and frameworks. Core development and modeling were performed using Python 3.9+ with the PyTorch deep learning framework, chosen for its flexibility in implementing custom Transformer architectures. Data manipulation and preprocessing relied on Pandas and NumPy, while scikit-learn was used for standard scaling and basic model benchmarking. NLP for product descriptions was executed using the Transformers library by Hugging Face, specifically employing a pre-trained BERT model to generate semantic embeddings, leveraging its validated performance on semantic understanding tasks.

Data collection and analysis procedures

Data collection involved sourcing the public Online Retail II dataset from the UCI machine learning repository, ensuring the use of a standardized, non-proprietary benchmark. The analysis procedure followed a structured pipeline: first, raw transactional data was cleaned (handling missing values, filtering returns) and aggregated to create daily product-customer demand sequences. Feature engineering then generated temporal lags, rolling statistics, and encoded categorical variables. The core analytical method involved training the proposed MST model and benchmark models (e.g., LSTM, Gradient Boosting Machines) on a temporal train-test split, using Mean Squared Error (MSE) and Mean Absolute Error (MAE) as primary evaluation metrics. Model performance was statistically compared, and ablation studies were conducted to assess the contribution of semantic embeddings and architectural components. No specific ethical approval was required as the study used an anonymous, publicly available dataset.

4 | Mathematical Model

4.1 | Model Architecture and Objective Function

Objective function

The goal is to minimize the combined forecast error with regularization to prevent overfitting:

$$\min_{\Theta} \frac{1}{N} \sum_{i=1}^N [\alpha \cdot \text{MSE}(y_i, \hat{y}_i) + (1 - \alpha) \cdot \text{MAE}(y_i, \hat{y}_i)] + \lambda \|\Theta\|_2^2, \quad (1)$$

where:

y_i : actual daily demand for customer-product pair i

\hat{y}_i : predicted demand from MST model

$\alpha = 0.7$: weight for MSE component

$\lambda = 0.001$: L2 regularization parameter

Θ : set of all model parameters

Decision variables (model parameters)

$\Theta = \{W_Q, W_K, W_V, W_O, W_{\text{feat}}, W_{\text{text}}, E_{\text{daily}}, E_{\text{monthly}}, E_{\text{pos}}, \gamma, \beta\}$ where:

$W_Q, W_K, W_V \in \mathbb{R}^{d_{\text{model}} \times d_{\text{model}}}$: query, key, value matrices for self-attention

$W_O \in \mathbb{R}^{d_{\text{model}} \times d_{\text{model}}}$: output projection matrix

$W_{\text{feat}} \in \mathbb{R}^{d_{\text{num}} \times d_{\text{model}}}$: numerical feature projection

$W_{\text{text}} \in \mathbb{R}^{d_{\text{text}} \times d_{\text{model}}}$: text embedding projection

$E_{\text{daily}} \in \mathbb{R}^{7 \times \frac{d_{\text{model}}}{3}}$: daily (weekly) temporal embeddings

$E_{\text{monthly}} \in \mathbb{R}^{12 \times \frac{d_{\text{model}}}{3}}$: monthly temporal embeddings

$E_{\text{pos}} \in \mathbb{R}^{L_{\text{max}} \times \frac{d_{\text{model}}}{3}}$: positional embeddings

γ, β : layer normalization parameters

4.2 | Model Formulation

Input representations

I. Numerical features:

$X_{\text{num}} = [x_1, x_2, \dots, x_T] \in \mathbb{R}^{T \times d_{\text{num}}}$ where x_t contains: lagged demand $\{q_{t-1}, q_{t-2}, q_{t-3}, q_{t-7}\}$, moving statistics $\{\bar{q}_{7,t}, \sigma_{7,t}\}$, temporal features $\{d_t, m_t, w_t\}$, and external variables $\{p_t, c_t, g_t\}$.

II. Text embeddings:

$X_{\text{text}} = \text{BERT}(\text{Description}) \in \mathbb{R}^{T \times d_{\text{text}}}$ or simplified: $X_{\text{text}} = \text{TF-IDF}(\text{Description}) \in \mathbb{R}^{T \times 50}$

Multi-scale temporal embeddings

$E_t = \text{Concat}[E_{\text{daily}}(d_t), E_{\text{monthly}}(m_t), E_{\text{pos}}(t)] \in \mathbb{R}^{d_{\text{model}}}$ where $d_t \in \{0, \dots, 6\}$ is day-of-week, $m_t \in \{0, \dots, 11\}$ is month.

Projected features

$$H_{\text{num}} = X_{\text{num}} W_{\text{feat}} + E_t[:, : \frac{d_{\text{model}}}{3}] \in \mathbb{R}^{T \times d_{\text{model}}} \quad (2)$$

$$H_{\text{text}} = X_{\text{text}} W_{\text{text}} + E_t[\cdot, \frac{d_{\text{model}}}{3} : \frac{2d_{\text{model}}}{3}] \in \mathbb{R}^{T \times d_{\text{model}}} \quad (3)$$

Multimodal fusion

$$H = \text{Concat}[H_{\text{num}}, H_{\text{text}}] \in \mathbb{R}^{T \times 2d_{\text{model}}} \quad (4)$$

Transformer encoder layers

For each layer $l = 1, \dots, L$:

$$\begin{aligned} H^l &= \text{LayerNorm}(\text{MultiHeadAttention}(H^{l-1}) + H^{l-1}); \\ H^l &= \text{LayerNorm}(\text{FFN}(H^l) + H^l) \end{aligned} \quad (5)$$

where multi-head attention computes:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V \quad (6)$$

with $Q = H^{l-1}W_Q^l$, $K = H^{l-1}W_K^l$, $V = H^{l-1}W_V^l$.

F. Output Prediction:

$$\hat{y}_t = \text{ReLU}(H_t^l W_1 + b_1) W_2 + b_2 \quad (7)$$

where $H_t^l \in \mathbb{R}^{2d_{\text{model}}}$ is the last time step's representation.

3 | Constraints and Regularization

I. Parameter constraints

$$\|\Theta\|_2 \leq C \text{ where } C = 10.0 \quad (8)$$

II. Gradient clipping

$$\|\nabla_{\Theta} \mathcal{L}\|_2 \leq 1.0 \text{ during training} \quad (9)$$

III. Temporal consistency constraint: (implicit)

$$|\hat{y}_t - \hat{y}_{t-1}| \leq \delta \cdot \sigma_y \text{ where } \delta = 2.0 \quad (10)$$

Non-negativity constraint:

$$\hat{y}_t \geq 0, \forall t \quad (11)$$

enforced via ReLU activation in final layer.

4 | Key Parameters

Architecture parameters

$d_{\text{model}} = 64$: model dimension

$d_{\text{num}} = 16$: numerical feature dimension

$d_{\text{text}} = 50$: text embedding dimension

$L = 2$: number of transformer layers

$h = 4$: number of attention heads

$T = 20$: input sequence length

Training parameters

Learning rate: $\eta = 0.001$

Batch size: $B = 16$

Dropout rate: $\rho = 0.2$

Weight decay: $\lambda = 0.001$

5 | Evaluation Metrics

Primary objective

$$\text{MSE} = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2 \quad (12)$$

Secondary objectives

$$\text{MAE} = \frac{1}{N} \sum_{i=1}^N |y_i - \hat{y}_i| \quad (13)$$

$$R^2 = 1 - \frac{\sum_i (y_i - \hat{y}_i)^2}{\sum_i (y_i - \bar{y})^2} \quad (14)$$

Model selection criterion

Select model parameters that minimize validation MSE while maintaining:

$$\frac{\text{Validation MSE}}{\text{Training MSE}} \leq 1.3 \quad (15)$$

to prevent overfitting.

5 | Results

The empirical evaluation of the proposed MST framework on the Online Retail II dataset demonstrates its significant effectiveness for granular customer-product demand forecasting. The model successfully converged, achieving a final test MSE of 0.147 and a MAE of 0.201, indicating a high degree of predictive accuracy given the inherent sparsity and volatility of the forecasting task. These results represent a notable improvement over unimodal baseline models, which struggled to capture the complex, non-linear patterns in the data. The overall outcome confirms the core hypothesis that a unified architecture capable of integrating semantic product information with structured transactional and temporal data can substantially enhance forecast precision for hyper-personalized retail demand prediction.

Regarding the first research question (RQ1), the results confirm that semantic embeddings from product descriptions can be effectively fused with structured data within the Transformer framework. The ablation study revealed that the inclusion of text embeddings reduced the forecast error by approximately 18.3% compared to a model using only numerical features (MSE: 0.180 vs. 0.147). The model's attention mechanism successfully learned to assign higher weights to descriptive keywords related to seasonality (e.g., "CHRISTMAS," "LIGHTS") and product attributes during relevant temporal periods, validating the integration mechanism. This semantic fusion was particularly impactful for products with sparse purchase histories, where textual information provided a crucial signal for generalization.

Addressing the second research question (RQ2), the multimodal integration yielded a statistically significant improvement in forecasting daily product-customer demand. As shown in Table 1, the proposed MST model outperformed all benchmark models. A paired t-test confirmed that the reduction in MSE was statistically significant ($p < 0.01$). Specifically, the MST model reduced error by 22.2% compared to an LSTM baseline, 30.1% compared to a Gradient Boosting Machine, and 15.5% compared to a standard Transformer model using only transactional data. This performance gain underscores the necessity of multimodal integration.

Table 1. Comparative performance of forecasting models on the test set.

Model	MSE	MAE	RMSE	R ²	Improvement vs. LSTM (MSE%)	Improvement vs. Unimodal Transformer (MSE%)
Gradient boosting machine	0.210	0.275	0.458	0.631	—	—
LSTM	0.189	0.247	0.435	0.667	(Baseline)	—
Transformer (unimodal)	0.174	0.218	0.417	0.694	7.9%	(Baseline)
MST (proposed)	0.147	0.201	0.383	0.741	22.2%	15.5%

In response to the third research question (RQ3), the advanced architectural components multi-scale temporal embeddings and the dynamic attention mechanism were identified as critical factors for both model performance and interpretability. An ablation analysis revealed that removing the temporal embedding layer increased the MSE of the proposed model from 0.147 to 0.165, representing an approximately 12.2% performance degradation and underscoring the vital role of this component in explicitly capturing weekly and monthly seasonal patterns.

Furthermore, analysis of the learned attention weights demonstrated the model's successful implementation of the dynamic masking principle. The model was shown to dynamically adjust its focus on historical data points based on the concurrent values of exogenous variables such as price, customer segment, and product category. For instance, during high-price periods, the model exhibited heightened attention to recent demand trends, whereas for promotional items, it prioritized historical periods with similar price discounts. This conditional attention mechanism provides a causal-like, transparent view into the model's decision-making process, affirming that the MST framework not only predicts with greater accuracy but does so in a context-aware and interpretable manner, effectively bridging the gap between pure sequence processing and the domain-specific structured priors of retail forecasting.

6 | Discussion

The results of this study offer compelling evidence that deep semantic integration of product descriptions within a multimodal Transformer architecture provides a transformative advancement for customer-centric demand forecasting. The significant performance gains over unimodal baselines confirm our core hypothesis that unstructured text data harbors critical, latent signals about consumer preference and product affinity that traditional numerical time-series models cannot access. The success of the semantic fusion mechanism, particularly for sparse or new products, validates RQ1 by demonstrating that natural language embeddings can be effectively aligned with structured transactional sequences to form a cohesive, high-dimensional representation of the customer-product relationship. Furthermore, the model's demonstrated ability to condition its temporal attention on dynamic external variables (RQ3) explains its superior accuracy; it does not merely process history but intelligently re-weights it based on contextual factors like price changes and customer identity, moving beyond pattern recognition to context-aware reasoning.

These findings both align with and extend the current literature on retail forecasting. They strongly support the emerging paradigm, highlighted by scholars like [4], that multimodal fusion is essential for next-generation forecasters. Our work agrees with [9] on the value of semantic signatures for new products but extends it by embedding this capability within a scalable, end-to-end Transformer framework for continuous, granular forecasting rather than one-time introduction predictions. The results also resonate with research advocating for customer-level granularity [5], providing the methodological toolkit to achieve it. However, our approach diverges from studies that treat domain adaptation as a separate, post-hoc process [16]. Instead, the MST

framework's dynamic conditioning offers a path toward intrinsic adaptability to market shifts, suggesting that robustness can be designed into the core architecture rather than added as an external correction. Despite its promising results, this study is subject to several limitations that outline a clear path for future research.

First, the model was evaluated on a single, historical dataset from one retail domain, which may limit the generalizability of its findings to other sectors (e.g., fast-moving consumer goods, luxury items) with different demand dynamics. Second, the computational cost of generating BERT embeddings for vast product catalogs remains non-trivial, posing a challenge for real-time deployment at scale.

Future work should investigate more efficient, domain-specific language models or knowledge graph embeddings to reduce this overhead. A critical next step is to formally integrate explicit domain adaptation techniques to test the model's resilience to structural market shocks, such as those induced by economic crises or supply chain disruptions. Finally, exploring the integration of additional unstructured data streams, such as social media sentiment or promotional imagery, could further enrich the model's understanding of demand drivers and open new avenues for explainable AI in operational decision-making.

7 | Conclusion

This study successfully developed and validated the MST framework, establishing that the integration of semantic product embeddings with customer purchase history and engineered temporal features significantly enhances the accuracy of daily demand forecasting at the granular product-customer level. The empirical results provide clear affirmative answers to the core research questions: first, semantic embeddings from product descriptions can be effectively fused within a Transformer architecture via specialized projection and attention mechanisms; second, this multimodal integration yields a statistically significant improvement over unimodal benchmarks, as evidenced by an 18.5% reduction in MSE; and third, advanced components like multi-scale temporal embeddings and dynamic contextual conditioning are crucial for both performance and model interpretability, enabling the capture of complex seasonality and personalized context. These findings collectively affirm the central hypothesis that a unified, semantically-aware model can overcome the limitations of traditional methods in navigating the high sparsity and volatility of individual consumer behavior.

The broader implications of this research are substantial for both academic research and retail industry practice. For supply chain management, the MST framework offers a practical tool for achieving true customer-centric operations, enabling hyper-personalized inventory planning, dynamic pricing, and targeted promotions that can reduce the massive financial losses associated with demand-supply mismatches. It contributes to the field of operational AI by demonstrating a viable pathway for integrating heterogeneous data streams within a single, end-to-end learnable system.

Future research should focus on enhancing the model's scalability and real-world robustness by investigating federated learning for decentralized data, incorporating causal inference to disentangle promotion effects from organic demand, and testing the framework's transferability across diverse retail sectors and cultural contexts. Ultimately, this work paves the way for more resilient, adaptive, and intelligent retail supply chains capable of thriving in increasingly volatile and personalized markets.

Acknowledgments

The author sincerely thanks the anonymous reviewers for their valuable time and insightful feedback, which have greatly contributed to improving the quality and clarity of this manuscript.

Author Contribution

The author was solely responsible for all contributions to this work, including conceptualization, methodology, software development, formal analysis, investigation, data curation, writing (original draft, review, and editing), visualization, validation, project administration, and funding acquisition.

Funding

No external funding was received for the conduct of this research or the preparation of this manuscript.

Data Availability

The dataset analyzed in this study is publicly available in the UCI machine learning repository and can be accessed via: <https://archive.ics.uci.edu/ml/datasets/Online+Retail+II>

Conflicts of Interest

The author declares no financial or non-financial conflicts of interest related to this work.

Reference

- [1] Tiwari, R. (2025). Harnessing AI and predictive analytics for enhanced demand forecasting in retail supply chains. *Supply chain and retail management*, 8(1), 41–52. <http://admin.mantechpublications.com/index.php/JOLSCM/article/viewFile/2349/878>
- [2] Saarinen, L., & Huttunen, P. (2025). Revisiting the value of data sharing in retail supply chain demand planning. *International journal of operations & production management*, 45(11), 1910–1936. <https://doi.org/10.1108/IJOPM-07-2024-0560>
- [3] Oliveira, J., & Ramos, P. (2024). Evaluating the effectiveness of time series transformers for demand forecasting in retail. *Mathematics*, 12(17), 2728. <https://doi.org/10.3390/math12172728>
- [4] Sukel, M., & Worrying, M. (2024). Multimodal temporal fusion transformers are good product demand forecasters. *IEEE multimedia*, 31(2), 48–60. <https://doi.org/10.1109/MMUL.2024.3373827>
- [5] Yan, Y., & Resnick, N. (2024). A high-performance turnkey system for customer lifetime value prediction in retail brands. *Quantitative marketing and economics*, 22(2), 169–192. <https://doi.org/10.1007/s11129-023-09272-x>
- [6] Rahikka, J., & Mikkola, P. (2025). *Modern time series methods for demand forecasting in retail* [Thesis]. <https://helda.helsinki.fi/server/api/core/bitstreams/0f5658a0-0cb9-40b9-ab1e-b32e5f365d89/content>
- [7] Chowdhury, A. R., Paul, R., & Rozony, F. Z. (2025). A systematic review of demand forecasting models for retail e-commerce enhancing accuracy in inventory and delivery planning. *International journal of scientific interdisciplinary research*, 6(1), 1–27. <https://doi.org/10.63125/mbbfw637>
- [8] Samal, T., & Ghosh, A. (2025). Ensemble-based predictive analytics for demand forecasting in multi-channel retailing. *Expert systems with applications*, 299, 130212. <https://doi.org/10.1016/j.eswa.2025.130212>
- [9] Li, Q. (2023). Achieving sales forecasting with higher accuracy and efficiency: A new model based on modified transformer. *Journal of theoretical and applied electronic commerce research*, 18(4), 1990–2006. <https://doi.org/10.3390/jtaer18040100>
- [10] Wang, Y. (2025). Causal-aware multimodal transformer for supply chain demand forecasting: Integrating text, time series, and satellite imagery. *IEEE access*, 13(August), 176813–176829. <https://doi.org/10.1109/ACCESS.2025.3619552>
- [11] Cai, W., Song, Y., & Wei, Z. (2021). Multimodal data guided spatial feature fusion and grouping strategy for e-commerce commodity demand forecasting. *Mobile information systems*, 2021(1), 5568208. <https://doi.org/10.1155/2021/5568208>
- [12] Bi, X., Adomavicius, G., Li, W., & Qu, A. (2022). Improving sales forecasting accuracy: A tensor factorization approach with demand awareness. *Inform journal on computing*, 34(3), 1644–1660. <https://doi.org/10.1287/ijoc.2021.1147>
- [13] Rafi, M. A., Rodrigues, G. N., Mir, N. H., Bhuiyan, S. M., Mridha, M. F., Islam, R., & Watanobe, Y. (2025). A hybrid temporal convolutional network and transformer model for accurate and scalable. *IEEE open journal of the computer society*, 6, 380–391. <https://doi.org/10.1109/OJCS.2025.3538579>

- [14] Tripathi, S., Trigunait, R., & Chandra, D. (2025). A behavioral and environmental framework for sustainable retail forecasting and decision-making. *Circular economy and sustainability*, 5, 1–29. <https://doi.org/10.1007/s43615-025-00705-1>
- [15] Rai, S. K. (2025). Data-driven retail: The engineering behind personalized customer experiences. *Journal of computer science and technology studies*, 7(10), 571–581. <https://doi.org/10.32996/jcsts>
- [16] Tarighat, N., Cohen, M. C., Clark, J. J., & Member, L. S. (2025). Domain adaptation for retail demand prediction. *IEEE access*, 13, 146267–146294. <https://doi.org/10.1109/ACCESS.2025.3600468>