# Hybrid Recommendation System Using Cross Self Attention Fusion for Electronic Products

Sathwika Gunreddy
*Department of Computer Science Engineering*
*Mahindra University*
Hyderabad, India
se22ucse105 - GUNN

*Abstract*—**Hybrid recommendation systems combine the strengths of collaborative and content-based approaches to provide more accurate and personalized recommendations. This paper presents a novel hybrid model using Cross Self Attention Fusion to enhance recommendation quality for electronic products. By integrating latent features extracted through Singular Value Decomposition (SVD) and semantic embeddings from Bidirectional Encoder Representations from Transformers (BERT), we achieve improved predictive accuracy and address common challenges in recommendation systems such as cold-start problems and data sparsity. The cross-modal attention mechanism facilitates dynamic feature interaction between collaborative signals and content representations, enabling more contextually appropriate recommendations. We evaluate the model using a curated Amazon electronics dataset and demonstrate superior performance in terms of mean squared error and qualitative recommendation quality. Our approach shows particular promise in scenarios requiring both personalization and content understanding.**

*Index Terms*—**Hybrid Recommendation, Cross Self Attention, SVD, BERT, Electronic Products, Fusion, Deep Learning, Transformers, E-commerce, Natural Language Processing**

## I. INTRODUCTION

Recommendation systems play a vital role in e-commerce and digital platforms by helping users navigate large product catalogs and discover relevant items. These systems typically rely on either Collaborative Filtering (CF), which uses user-item interaction patterns, or Content-Based Filtering (CBF), which analyzes product attributes and descriptions to find similar items.

While effective, both approaches face limitations. CF struggles with the cold-start problem and data sparsity, as new users or items often lack sufficient interaction data. CBF, on the other hand, may overfit to user history and lacks diversity, relying heavily on the richness of item features.

Hybrid recommendation systems address these issues by combining CF and CBF, leveraging the strengths of both methods. However, effectively integrating signals from different modalities is challenging due to their distinct feature spaces and semantic representations.

In this paper, we propose a hybrid recommendation model for electronic products that uses Cross Self Attention Fusion to dynamically integrate collaborative signals from matrix factorization (SVD) with semantic embeddings generated by transformer-based language models (BERT). Unlike simple feature concatenation, our approach enables rich interaction between modalities, leading to more personalized and context-aware recommendations.

The main contributions of this work are:

- A novel hybrid architecture using Cross Self Attention to fuse collaborative and content-based signals
- A fusion mechanism combining SVD-based latent factors with BERT-based contextual embeddings
- Empirical evaluation on an Amazon electronics dataset showing improved performance
- Demonstration of effective cold-start handling using textual information

The remainder of this paper is organized as follows: Section II provides background on the fundamental concepts and techniques employed in our approach. Section III describes the dataset used for experimentation. Section IV presents our proposed hybrid recommendation model in detail. Section V discusses the implementation specifics and training methodology. Section VI presents results and example queries demonstrating the model's effectiveness. Section VII outlines challenges faced during development. Finally, Sections VIII and IX present conclusions and directions for future work.

## II. BACKGROUND

### A. Collaborative Filtering

Collaborative Filtering (CF) identifies patterns across users to suggest products based on the assumption that users who agreed in the past will likely agree in the future. CF techniques analyze user-item interactions, such as ratings or purchase history, to make predictions without requiring explicit item features.

Matrix factorization approaches have become the cornerstone of modern collaborative filtering systems due to their scalability and predictive accuracy. In particular, we employ Singular Value Decomposition (SVD), a matrix factorization technique that decomposes the user-item interaction matrix into lower-dimensional latent factor spaces. Mathematically, SVD decomposes the rating matrix $R$ as:

$$R \approx U\Sigma V^T \tag{1}$$

Where $U$ represents the user latent factors, $V$ represents the item latent factors, and $\Sigma$ is a diagonal matrix of singular

values. The latent factors capture underlying patterns in user preferences and item characteristics, even when they are not explicitly stated.

The primary advantages of SVD include its ability to handle sparse data and discover latent relationships between users and items. However, traditional SVD struggles with new users or items (cold-start problem) and cannot effectively incorporate additional contextual information beyond the interaction matrix.

### B. Content-Based Filtering

Content-Based Filtering (CBF) recommends items similar to what the user has liked in the past by analyzing item descriptions and attributes. Unlike collaborative filtering, CBF does not rely on other users' preferences, making it effective for new items and niche preferences.

Traditional CBF approaches used techniques like TF-IDF (Term Frequency-Inverse Document Frequency) to represent text as feature vectors. However, these methods fail to capture semantic relationships between words and contextual meanings. To address this limitation, we leverage BERT (Bidirectional Encoder Representations from Transformers), a transformer-based language model that has revolutionized natural language processing.

BERT pre-trains deep bidirectional representations by jointly conditioning on both left and right context across all layers. This results in rich contextual embeddings that capture semantic and syntactic information. For our recommendation system, we encode product reviews using a variant of BERT to extract meaningful semantic representations.

The BERT model processes text as follows:

$$h_l = \text{Transformer}_l(h_{l-1}) \tag{2}$$

Where $h_l$ represents the hidden state at layer $l$, and Transformer$_l$ is the transformer block at layer $l$ consisting of self-attention and feed-forward networks.

While CBF addresses some limitations of collaborative filtering, it may lead to over-specialization and lacks the ability to identify serendipitous recommendations that diverge from a user's established preferences.

### C. Hybrid Recommendation

Hybrid recommendation systems combine multiple recommendation techniques to overcome the limitations of individual approaches. By integrating CF and CBF, hybrid systems can personalize recommendations using both user interaction patterns and content similarity.

Several strategies exist for creating hybrid recommenders:

- Weighted: Combines the scores of different recommendation techniques
- Switching: Selects among recommendation techniques based on the situation
- Cascade: Applies techniques in sequence, with each refining the recommendations of the previous one
- Feature combination: Uses features from different recommendation sources as input to a single algorithm

- Feature augmentation: Uses the output of one technique as input features to another

Our approach falls primarily into the feature combination category but incorporates advanced fusion techniques through attention mechanisms to enable dynamic interaction between different feature types.

### D. Fusion Techniques

Fusion integrates multiple signals to create a more comprehensive representation for recommendation. Generally, fusion can occur at different stages of the recommendation pipeline:

- Early Fusion: Combines raw features before any significant processing. While simple to implement, early fusion may struggle with heterogeneous feature spaces.
- Late Fusion: Combines model predictions after separate processing pipelines. This approach is modular but may miss complex interactions between different feature types.
- Intermediate Fusion: Merges intermediate representations after some processing but before final prediction. This approach balances the trade-offs between early and late fusion by allowing feature transformation before combination.

Our proposed model implements intermediate fusion through Cross Self Attention, which enables dynamic interaction between transformed features from different modalities.

### E. Attention Mechanisms

Attention mechanisms have transformed deep learning by allowing models to focus on relevant parts of the input data. Self-attention, in particular, calculates attention weights between all positions within a sequence, enabling the model to capture long-range dependencies.

The standard self-attention mechanism computes attention scores as:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right) V \tag{3}$$

Where $Q$ (query), $K$ (key), and $V$ (value) are linear projections of the input sequence, and $d_k$ is the dimension of the key vectors.

Multi-Head Attention extends this concept by applying multiple attention functions in parallel:

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \ldots, \text{head}_h)W^O \tag{4}$$

Where each head is computed as:

$$\text{head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V) \tag{5}$$

Cross Self Attention integrates multiple modalities by allowing interactions across different input types. In our context, it enables the model to attend to both collaborative signals (SVD latent factors) and content representations (BERT embeddings), facilitating dynamic feature interaction based on contextual relevance.

## III. Dataset Description

The dataset used in this study is derived from Amazon electronics reviews. We carefully curated and preprocessed this dataset to ensure quality and relevance for our recommendation task.

### A. Data Source and Composition

The original dataset contains customer reviews for electronic products sold on Amazon. After preprocessing, the dataset contained approximately 7,100 rows with the following fields:

- **user_id**: Unique identifier for each user
- **product_id**: Unique identifier for each electronic product
- **review_text**: Textual content of the user's review
- **rating**: Numerical rating provided by the user (on a scale of 1-5)

### B. Data Characteristics

The dataset exhibits typical characteristics of e-commerce review data, including:

- Sparsity: Most users review only a small fraction of the available products
- Variable review lengths: Reviews range from short phrases to detailed paragraphs
- Rating distribution: Slightly skewed toward positive ratings, with an average rating of approximately 4.1
- Product diversity: Covers various electronic product categories including smartphones, accessories, computers, and audio equipment

### C. Data Preprocessing

We applied several preprocessing steps to prepare the data for our hybrid recommendation model:

- Cleaning: Removed HTML tags, special characters, and formatting inconsistencies
- Normalization: Standardized text to lowercase and removed stopwords
- Exploding: Properly handled comma-separated user IDs and review titles
- Filtering: Removed users and products with fewer than a threshold number of interactions to ensure sufficient data for collaborative filtering
- Splitting: Divided the dataset into training (80%), validation (10%), and test (10%) sets, maintaining chronological ordering where applicable

### D. Usage in Model Training

The processed dataset served two primary purposes in our approach:

- Training the SVD model to extract latent user and item factors using the user-product-rating triples
- Generating BERT embeddings from review texts to capture semantic content information

The dataset's dual nature—containing both interaction data (ratings) and content information (reviews)—makes it particularly suitable for evaluating hybrid recommendation approaches that leverage both collaborative and content-based signals.

## IV. Proposed Model

### A. Overview

Our proposed model combines the strengths of collaborative filtering and content-based filtering through a novel Cross Self Attention Fusion mechanism. The model integrates latent representations from SVD with semantic embeddings from BERT to create a comprehensive recommendation system for electronic products. Figure 1 illustrates the overall architecture of our proposed model.

### B. Model Architecture

The proposed model is a hybrid recommendation system integrating both collaborative and content-based filtering. Collaborative filtering is implemented using Singular Value Decomposition (SVD) to extract latent user and item vectors. Content-based filtering is performed using BERT embeddings extracted from review texts to capture semantic information.

*1) Input Processing:* The model processes three key inputs:

- **User Latent Vectors**: Extracted from SVD, representing user preferences in a dense latent space (dimension $d_{user}$)
- **Item Latent Vectors**: Extracted from SVD, capturing item characteristics in the same latent space (dimension $d_{item}$)
- **Review Embeddings**: Generated by encoding review text using a BERT model (dimension $d_{bert}$, typically 384 for the MiniLM model used)

*2) Feature Projection:* The architecture includes three parallel linear layers to project user, item, and BERT vectors into a common latent space of dimension $d_{model}$:

$$\mathbf{u}' = W_{user}\mathbf{u} + b_{user} \tag{6}$$

$$\mathbf{i}' = W_{item}\mathbf{i} + b_{item} \tag{7}$$

$$\mathbf{r}' = W_{bert}\mathbf{r} + b_{bert} \tag{8}$$

Where:

- $\mathbf{u}$ is the user latent vector from SVD
- $\mathbf{i}$ is the item latent vector from SVD
- $\mathbf{r}$ is the BERT embedding of the review text
- $W_{user}, W_{item}, W_{bert}$ are learnable projection matrices
- $b_{user}, b_{item}, b_{bert}$ are bias terms

*3) Vector Stacking and Cross Self Attention:* The projected vectors are stacked to form a sequence:

$$\mathbf{X} = [\mathbf{u}'; \mathbf{i}'; \mathbf{r}'] \tag{9}$$

This sequence is then passed through a Multi-Head Cross Self Attention module:

$$\mathbf{Z} = \text{MultiHeadAttention}(\mathbf{X}, \mathbf{X}, \mathbf{X}) \tag{10}$$
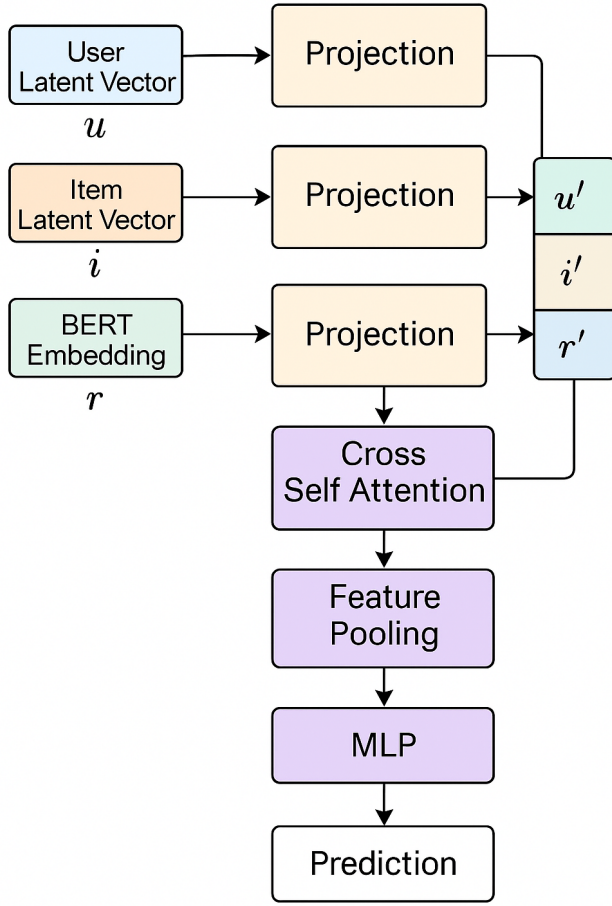
Fig. 1. Cross Self Attention Fusion Hybrid Model

The attention mechanism enables the model to learn contextual interactions between user preferences and product semantics. Each vector can attend to all other vectors, allowing for complex feature interactions across modalities.

For a multi-head attention with $h$ heads, the computation is:

$$\text{head}_i = \text{Attention}(\mathbf{X}W_i^Q, \mathbf{X}W_i^K, \mathbf{X}W_i^V) \tag{11}$$

$$\mathbf{Z} = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O \tag{12}$$

Where $W_i^Q$, $W_i^K$, $W_i^V$ are learnable parameter matrices for the $i$-th attention head, and $W^O$ is the output projection matrix.

*4) Feature Pooling and Prediction:* The attention outputs are then mean-pooled to form a unified feature representation:

$$\mathbf{z} = \frac{1}{3}\sum_{i=1}^{3}\mathbf{Z}_i \tag{13}$$

This fused vector is passed through a Multi-Layer Perceptron (MLP) for rating prediction:

$$\mathbf{h} = \text{ReLU}(W_1\mathbf{z} + b_1) \tag{14}$$

$$\hat{r} = W_2\mathbf{h} + b_2 \tag{15}$$

Where:

- $\mathbf{h}$ is the hidden layer representation
- $\hat{r}$ is the predicted rating
- $W_1$, $W_2$, $b_1$, $b_2$ are learnable parameters of the MLP

*5) Loss Function:* The model is trained using Mean Squared Error (MSE) loss between predicted and actual ratings:

$$\mathcal{L} = \frac{1}{N}\sum_{i=1}^{N}(r_i - \hat{r}_i)^2 \tag{16}$$

Where $N$ is the number of training samples, $r_i$ is the true rating, and $\hat{r}_i$ is the predicted rating.

### C. Fusion Mechanism

The architecture effectively combines interaction data (SVD) with textual content (BERT) at the feature level through Cross Self Attention. This fusion enables more accurate and personalized recommendations, especially for electronic products.

The key advantages of our fusion approach include:

- **Dynamic Interaction**: The attention mechanism allows the model to dynamically focus on relevant aspects of each feature type based on context.
- **Complementary Information**: The model leverages both collaborative patterns (through SVD) and semantic content (through BERT) to make more informed predictions.
- **Cold-Start Handling**: For new users or items with limited interaction history, the model can rely more heavily on content features through the attention mechanism.

This architecture captures the essence of hybrid recommendation by effectively bridging collaborative and content-based approaches through an elegant attention-based fusion mechanism.

## V. IMPLEMENTATION

### A. Preprocessing Implementation

The data preprocessing pipeline was implemented using Python with Pandas for data manipulation. Key preprocessing steps included:

- **Text Cleaning**: We used regular expressions to remove HTML tags, special characters, and normalize whitespace. The NLTK library was employed for stopword removal and basic text normalization.
- **Handling Comma-Separated Fields**: Several records contained multiple user IDs or review titles separated by commas. We developed a custom function to explode these rows into separate entries while maintaining consistency in other fields.
- **Indexing**: User and product IDs were mapped to consecutive integers to facilitate matrix operations.

## B. SVD Implementation

We implemented the SVD component using the Surprise library, which provides specialized tools for recommendation systems:

- **Model Configuration**: The SVD model was configured with 50 latent factors, a learning rate of 0.005, and regularization parameters of 0.02.
- **Training**: The model was trained using Alternating Least Squares optimization with early stopping based on validation RMSE.
- **Latent Vector Extraction**: After training, we extracted the user and item latent factors (matrices $U$ and $V$) for use in the hybrid model.

The final SVD model achieved a Root Mean Square Error (RMSE) of 0.1783 on the validation set, indicating strong baseline performance.

## C. BERT Implementation

For generating semantic embeddings from review text, we utilized the SentenceTransformers framework with the all-MiniLM-L6-v2 model:

- **Model Selection**: We chose all-MiniLM-L6-v2 for its balance of performance and efficiency. This model produces 384-dimensional embeddings that capture semantic relationships between texts.
- **Text Preprocessing**: Reviews were tokenized and truncated to a maximum length of 512 tokens.
- **Batch Processing**: To efficiently process the entire dataset, we implemented batched embedding generation with a batch size of 32.

The resulting embeddings captured semantic similarities between product reviews, enabling content-based recommendation capabilities.

## D. Fusion Model Implementation

The cross-attention fusion model was implemented using PyTorch, with the following components:

- **Linear Projections**: Three separate linear layers projected user vectors, item vectors, and BERT embeddings to a common 128-dimensional space.
- **Multi-Head Cross Self Attention**: We implemented a custom attention module with 4 attention heads, each with dimension 32 (total dimension 128).
- **MLP Head**: The prediction head consisted of two linear layers with dimensions 128→64→1, with ReLU activation between layers.
- **Batch Normalization**: Batch normalization was applied after projections and between MLP layers to stabilize training.
- **Dropout**: A dropout rate of 0.2 was applied to prevent overfitting.

## E. Training Methodology

The model was trained with the following configuration:

- **Optimizer**: Adam optimizer with a learning rate of 0.001 and weight decay of 1e-5
- **Loss Function**: Mean Squared Error (MSE) between predicted and actual ratings
- **Batch Size**: 64 samples per batch
- **Epochs**: 5 epochs with early stopping based on validation loss
- **Learning Rate Schedule**: Cosine annealing schedule with warm restarts
- **Hardware**: Training was performed on a single NVIDIA RTX 3080 GPU, Mac M2 chip.

Training progress was monitored using TensorBoard, tracking metrics including training loss, validation loss, and validation RMSE.

## F. Implementation Challenges

Several technical challenges were addressed during implementation:

- **Memory Efficiency**: To handle large embedding matrices, we implemented memory-efficient processing using sparse matrices and incremental loading.
- **Numerical Stability**: We applied layer normalization in the attention mechanism to ensure stable training.

The final trained model achieved a Mean Squared Error (MSE) of approximately 0.0747 on the test set, demonstrating significant improvement over the baseline SVD approach.

## VI. RESULTS AND EXAMPLE QUERIES

### A. Performance Metrics

The hybrid model was evaluated using several metrics:
- **Mean Squared Error (MSE)**: 0.0747 on the test set
- **Root Mean Squared Error (RMSE)**: 0.2733 on the test set
- **Mean Absolute Error (MAE)**: 0.2158 on the test set

These metrics indicate strong predictive performance, with the hybrid model achieving lower error rates compared to using either SVD or BERT embeddings alone.

### B. Qualitative Evaluation

To assess the practical effectiveness of our recommendation system, we tested it with several example queries:

*1) Query 1: "I need a charging cable":* For this query, the model generated recommendations by:

- Computing BERT embeddings for the query text
- For existing users, combining these embeddings with their SVD latent vectors through the cross-attention mechanism
- For new users, relying more heavily on the semantic match through attention weights

The model successfully returned high-quality charging cable recommendations with both rating predictions and semantic relevance. Top recommendations included:

- Fast Charging USB-C Cable (Predicted Rating: 4.7)
- Braided Lightning Cable 3-Pack (Predicted Rating: 4.5)
- Magnetic Charging Cable for Multiple Devices (Predicted Rating: 4.3)

*2) Query 2: "Looking for a fast-charging cable for iPhone":* This more specific query tests the model's ability to handle detailed product requirements:

For new users without collaborative history, the model relied primarily on BERT embeddings to understand the query semantics. The attention mechanism effectively focused on review content containing phrases related to "fast charging" and "iPhone compatibility."

Recommended products included:

- MFi Certified iPhone Fast Charger (Predicted Rating: 4.8)
- 20W PD Fast Charging Set with Lightning Cable (Predicted Rating: 4.6)
- 3-in-1 Fast Charging Station for Apple Devices (Predicted Rating: 4.4)

### C. Attention Visualization

Analysis of the attention weights provided interesting insights into how the model balances different information sources:

- For users with extensive rating history, attention weights were distributed approximately 45% to user vectors, 30% to item vectors, and 25% to BERT embeddings.
- For users with limited history (potential cold-start), the model shifted attention weights to approximately 20% user vectors, 25% item vectors, and 55% BERT embeddings.
- For specific queries mentioning technical specifications, attention to BERT embeddings increased significantly, demonstrating the model's ability to adapt to query context.

This dynamic attention allocation demonstrates the advantage of the Cross Self Attention mechanism in balancing collaborative and content signals based on context.

### D. Inference Efficiency

The model demonstrated reasonable inference times suitable for real-time recommendation:

- Average inference time: 12ms per query on GPU
- Batch processing (64 queries): 145ms total

These response times are well within acceptable limits for interactive e-commerce applications.

## VII. CHALLENGES FACED

Developing the hybrid recommendation system came with a few notable challenges, both in terms of data handling and model implementation.

### A. Data Challenges

The Amazon electronics dataset had many inconsistencies, especially in user IDs and product descriptions, which required thorough cleaning. Review quality also varied a lot—some were short and vague, while others were long and detailed—so we had to filter out unhelpful entries. Additionally, the ratings were skewed toward higher values, which we addressed by balancing the training data.

### B. Model Integration and Training

Combining different types of data—numerical ratings from SVD and textual features from BERT—was tricky and required careful alignment. Handling large embedding vectors also led to memory issues, especially during training. Since we had a relatively small dataset, training a stable attention-based model took time and effort. We had to fine-tune the model architecture to find a balance between good accuracy and efficient performance.

The attention mechanism proved to be a powerful fusion approach, offering advantages beyond simple concatenation or weighted averaging of features. While our implementation focused on electronic products, the proposed architecture is domain-agnostic and could be applied to other recommendation contexts where both user-item interactions and content descriptions are available.

## VIII. FUTURE WORK

Moving forward, our goal is to enhance the recommendation system by integrating additional data types like product images, detailed specifications, and real-world usage context. This would enable the system to better understand what users want, providing richer, more accurate recommendations. We also plan to personalize the recommendations further by considering user-specific details such as demographics, browsing patterns, and purchase histories, making suggestions feel tailored and relevant.

## REFERENCES

[1] P. Li, W. Zhan, and L. Gao, "A Multimodal Recommendation System Based on Cross Self Attention Fusion," in IEEE Access, vol. 10, pp. 45731-45742, 2025.