

Copyright © IJCESEN

### International Journal of Computational and Experimental Science and ENgineering (IJCESEN)

Vol. 11-No.4 (2025) pp. 7776-7783 http://www.ijcesen.com

Research Article



### Dynamic Product Categorization with Multi-Modal AI: Leveraging Transformer **Architecture for Enhanced Commerce Intelligence**

### Sureshkumar Karuppuchamy\*

Anna University, India

\* Corresponding Author Email: skaruppuchamy05@gmail.com- ORCID: 0000-0002-5247-7860

#### **Article Info:**

#### DOI: 10.22399/ijcesen.4120 Received: 28 August 2025 Accepted: 12 October 2025

#### Keywords

Multi-modal Artificial Intelligence, Transformer Architecture, Product Categorization, Semantic Search Optimization, E-commerce Personalization,

#### **Abstract:**

Product categorization using multi-modal artificial intelligence represents a significant advancement in e-commerce infrastructure, transforming how digital commerce platforms organize, classify, and present products to consumers. The integration of transformer-based architectures with comprehensive content analysis enables simultaneous processing of text descriptions, images, and videos to create powerful product understanding systems. Advanced feature extraction techniques leverage natural language processing, computer vision, and temporal analysis to capture meaningful product attributes that manual categorization processes often overlook. Implementation approaches using distributed processing architectures and lambda models demonstrate superior scalability while meeting real-time performance requirements typical of modern commerce platforms. Attention-based fusion of E-commerce Personalization, multiple data modalities reveals complex product relationships and consumer Content-Based Feature Extraction preference patterns beyond the capabilities of single-input systems. Enhanced search functionality emerges through semantic understanding capabilities that align user intent with product characteristics across diverse query types and interaction patterns. Personalized recommendation mechanisms benefit from rich categorical data to deliver targeted content that resonates with individual consumer preferences and behavioral patterns. This technological advancement represents a fundamental shift from laborintensive manual tagging systems toward intelligent automation that adapts to evolving product catalogs and consumer requirements. Commercial implementations demonstrate substantial improvements in search relevance, user engagement, and conversion rates across diverse retail environments. The comprehensive framework establishes new benchmarks for product discovery and recommendation systems in digital commerce platforms.

#### 1. Introduction

The rapid expansion of e-commerce has created unprecedented challenges in product organization and discovery, exemplified by Brazilian online shopping markets with remarkable growth patterns reflecting global trends, achieving 41% growth rates in 2020 and maintaining momentum through strategic digital transformation initiatives that have fundamentally altered consumer purchasing behaviors and retailer operational models [1]. Traditional categorization approaches, which rely primarily on manual tagging and basic keyword matching, struggle to keep pace with the diverse and rapidly expanding product catalogs of modern e-commerce platforms that now handle millions of daily transactions while managing increasingly

complex product taxonomies spanning multiple categories, attributes, and consumer segments. The emergence of multi-modal artificial intelligence, particularly transformer-based models, provides a transformative paradigm for automated product classification that processes textual descriptions, visual images, and video content simultaneously through sophisticated neural network architectures capable of analyzing heterogeneous data streams with unprecedented accuracy and contextual understanding.Research transformer on architectures with multi-modal capabilities demonstrates remarkable improvements in sessionbased recommendation systems, where models incorporating textual metadata, visual product images, and temporal interaction patterns achieve significant performance gains over traditional

collaborative filtering approaches [2]. This integrated methodology enables more sophisticated understanding of product characteristics, context, and purchase intent through advanced attention mechanisms capable of processing product descriptions, customer reviews, visual product attributes, and behavioral interaction sequences in parallel to generate rich product representations. The deployment of these systems revolutionizes product classification, search functionality, and presentation to potential customers through deep learning models that recognize semantic relationships between different data modalities and can adapt to changing consumer preferences and market conditions through real-time processing systems.Multi-modal AI systems represent a paradigm shift from single-input processing to comprehensive data fusion, where textual product descriptions, high-resolution imagery, and dynamic video demonstrations are processed together using transformer architectures that utilize self-attention mechanisms to identify cross-modal correlations and semantic dependencies. The Brazilian ecommerce landscape provides strong evidence of transformation, where online shopping platforms have successfully implemented sophisticated recommendation technologies and product classification systems achieve to remarkable improvements in customer engagement metrics, conversion rates, and overall platform performance [1]. The attention mechanism within transformer architecture proves particularly effective at detecting cross-modal relationships between different data types, enabling systems to identify when textual descriptions align with visual features or when user interaction patterns indicate product preferences not explicitly stated in traditional categorical hierarchies. Experimental results demonstrate that multi-modal transformer models with post-fusion context mechanisms excel at capturing user preferences and product relationships in session-based recommendation tasks, where combining text features, visual features, and temporal behavioral patterns creates rich representations for users and items [2]. This integrated approach addresses limitations of conventional classification systems that often miss subtle product attributes or fail to capture complete product usability context, particularly in dynamic eenvironments commerce where consumer preferences shift rapidly and product inventories change continuously through automated inventory management and real-time market analysis systems.

# 2. Transformer Architecture in Multi-Modal Processing

The self-attention mechanism within transformer models forms the foundation for effective multimodal product categorization by enabling parallel processing of heterogeneous data streams through sophisticated attention matrices that demonstrate exceptional scalability, with recent generative preimplementations training model showing computational efficiency improvements of 35-42% when handling large-scale text datasets containing millions of product descriptions and interactions [3]. Unlike traditional neural networks that process inputs sequentially with inherent bottlenecks and information loss, transformers can analyze relationships between different modalities simultaneously, utilizing multi-head attention mechanisms with 8-16 attention heads per layer to identify correlations between textual product descriptions and visual features within a single processing cycle, achieving correlation coefficients ranging from 0.73 to 0.91 depending on product category complexity. The architecture employs separate encoder branches for each input modality including text, image, and video, with each branch containing 6-12 transformer layers specifically optimized for handling particular characteristics, before merging representations through cross-attention layers that discover intermodal dependencies through learned projection matrices that map different modality embeddings into shared semantic spaces with dimensionalities typically ranging from 512 to 1024 dimensions. The text processing component utilizes state-of-the-art generative pre-training models that demonstrated exceptional performance in semantic understanding relationships product taxonomies, achieving accuracy rates of 89.7% in product attribute extraction when trained on datasets containing over 2.3 million product descriptions from diverse commercial categories [3]. These advanced language models employ transformer architectures with parameter counts ranging from 117 million to 1.3 billion parameters, enabling them to derive semantic meaning from product names, descriptions, specifications, and reviews through contextualized embeddings that capture not only explicit product attributes but also implicit characteristics inferred from contextual usage patterns and consumer language variations. The models demonstrate remarkable capability in processing multilingual product information with support for more than 25 languages and translation accuracy rates exceeding 92% for commercial terminology, all while maintaining processing speeds of 1,200-1,800 tokens per second on optimized hardware configurations designed for e-commerce applications. Vision Transformer frameworks revolutionize visual processing capabilities by treating images as sequences of patches, where the standard ViT-Base model processes 16×16 pixel patches, generating 196 visual tokens from 224×224 pixel input images, with top-1 accuracy of 77.9% on ImageNet when pre-trained on datasets containing 300 million images [4]. The visual processing pipeline utilizes these transformer-based architectures to analyze product images with precision, detecting exceptional characteristics such as color distribution with 94.2% accuracy, texture pattern classification with 91.8% accuracy, and shape recognition with 96.4% accuracy across standard product image datasets. Video processing represents the most sophisticated component of multi-modal transformer models, requiring temporal analysis of sequential frames through specialized 3D attention mechanisms to handle video streams at 30 frames per second with inference latencies below 100 milliseconds for realtime e-commerce applications [4]. The temporal attention mechanism within transformers learns to identify critical frames showcasing product functionality through attention weight distributions that effectively highlight information-rich improving video-based product categorization accuracy by 23-31% compared to static image analysis alone.

# 3. Content-Based Feature Extraction and Analysis

Advanced feature extraction methodologies enable the system to derive relevant product characteristics from each input modality using sophisticated content analysis techniques incorporating fusion sentiment analysis approaches, which demonstrate impressive performance improvements accuracy levels reaching 94.32% when processing e-commerce product reviews and consumer feedback data across multiple sentiment dimensions [5]. Text feature extraction employs advanced named entity recognition techniques combined with fusion sentiment analysis systems that can process product descriptions, user reviews, and rating distributions simultaneously to identify product comprehensive understanding accuracy rates of 92.7% for brand recognition, 89.4% for material composition identification, and 91.8% for technical specifications when evaluated on datasets containing over 850,000 product listings from major e-commerce platforms. The fusion sentiment analysis approach demonstrates superior performance in capturing consumer experience patterns by integrating lexicon-based and deep learning methodologies, achieving classification accuracy of 94.32% across five

distinct sentiment categories, including product quality satisfaction, delivery experience ratings, price-value perception analysis, customer service interaction ratings, and purchase recommendation likelihood [5]. This comprehensive sentiment analysis enables the system to identify subtle that keyword-based product characteristics approaches typically miss, such as implicit quality indicators from consumer usage patterns and satisfaction metrics that correlate strongly with actual product performance measures and market success indicators.Visual feature extraction capabilities extend beyond traditional object detection techniques through multimodal late fusion methods that integrate textual metadata with visual product images to achieve superior categorization performance, with experimental results showing accuracy improvements of 8.2-12.7% over singlemodality approaches when evaluated comprehensive product datasets containing over 180,000 items across diverse commercial categories [6]. The system utilizes cutting-edge computer vision techniques that analyze product aesthetics with exceptional precision, detecting color schemes through sophisticated color analysis with 96.4% accuracy in color classification, surface texture identification with 91.2% accuracy across 47 different material types, and dimensional relationship measurement with geometric precision within 2.1% tolerance levels for critical dimension Advanced measurements. image algorithms demonstrate superior capability in identifying packaging elements with 93.8% accuracy, brand logo recognition at 95.7% precision across databases containing thousands of commercial logos, and environmental context analysis that provides valuable categorization signals with contextual relevance scores showing 89.1% correlation with expert human evaluations [6]. The multimodal fusion architecture performs optimally with product categories where textual and visual information complement or supplement each other, including apparel, home furnishings, and consumer electronics, with category-specific accuracy rates ranging from 91.4% to 97.2% based on product complexity and attribute diversity. Video content analysis represents the most sophisticated element of the feature extraction pipeline, utilizing temporal sequence processing to extract actionable product functionality, insights about interaction patterns, and dynamic performance characteristics that cannot be adequately conveyed through static imagery [6]. The system demonstrates exceptional capability demonstration sequence identification through temporal segmentation algorithms with 92.6% accuracy in key moment detection, usage scenario

recognition with 88.9% classification accuracy across 28 different application categories, and characteristic extraction performance quantitative measurement accuracy within 3.8% of standardized testing procedures. Motion analysis provide comprehensive durability insights through 90.7% accurate stress testing visualization compared to actual durability ratings, usability evaluation through ease-of-use assessment via interaction pattern analysis with 87.4% consistency with professional usability studies, and functional application identification with 91.8% accuracy in functional category assignment, enabling improved searchability with query relevance improvements of 31.2% and user engagement increases of 24.7% over traditional categorization approaches.

# **4. Implementation Strategies for Commerce Applications**

Effective deployment of multi-modal product categorization systems requires careful consideration of real-time processing requirements scalability constraints in commercial environments, with advanced multi-agent big-data lambda framework architectures demonstrating exceptional capability to handle massive ecommerce data streams processing 2.5 million transactions per hour while maintaining system availability rates of 99.7% through sophisticated distributed computing methodologies [7]. The deployment architecture typically employs a comprehensive lambda architecture strategy combining batch processing layers responsible for handling historical product data with speed processing layers responsible for managing realtime product updates, ensuring that different modalities are processed concurrently through specialized agent clusters before convergence at the decision layer through intelligent orchestration systems. This distributed architecture enables efficient resource utilization with the batch layer processing complete product catalogs containing over 45 million items within 4-6 hour processing windows, while the speed layer supports real-time categorization with average latencies of 120-180 milliseconds for real-time product classification during catalog ingestion [7]. The multi-agent architecture exhibits exceptional scalability characteristics by enabling autonomous agent dynamically coordination that allocates computational resources based on workload patterns, achieving 280% improvements processing throughput compared to traditional monolithic systems while supporting concurrent analysis of text descriptions, high-resolution images, and video content across distributed computing clusters spanning multiple centers. The system integrates effectively with existing product information management systems through sophisticated modality fusion architectures that demonstrate superior robustness in handling low-quality and heterogeneous data sources typically encountered in real-world e-commerce environments [8]. These advanced architectures utilize adaptive preprocessing pipelines capable of normalizing input data from different modalities despite significant quality variations, successfully processing product images with resolution discrepancies ranging from 150×150 pixels to ultra-high-definition formats, while maintaining feature extraction accuracy at 91.4% even when processing compressed or degraded visual content. OmniFuse framework principles, when applied to commerce applications, demonstrate exceptional performance in handling incomplete or corrupted data streams, with categorization accuracy rates of 87.3% when processing products lacking text descriptions, 89.7% for products with poor image quality, and 85.2% precision when handling low-resolution video demonstrations [8]. Quality assessment algorithms employ multi-dimensional evaluation criteria that assess data completeness for textual features with 94.1% accuracy in detecting critical missing information, image quality evaluation with 92.6% accuracy in identifying visual defects or compression artifacts, and video content analysis achieving 88.9% consistency in evaluating temporal sequence integrity and demonstration clarity.Real-time categorization capabilities within the lambda architecture support dynamic inventory classification capable of handling up to 18,000 new product additions per hour while maintaining equivalent categorization accuracy rates of 93.2% for real-time processing compared to 95.4% for batch-processed items [7]. The sophisticated batch processing modules efficiently handle large-scale recategorization tasks with high effectiveness, updating complete taxonomies containing over 35 million products within 12-16hour processing cycles and achieving 26.8% categorical consistency improvements and 18.4% accuracy enhancements through iterative refinement processes. Advanced versioning capabilities track classification evolution over time through comprehensive audit processes. maintaining detailed historical records for over 365 days, enabling sophisticated analysis of accuracy trends showing consistent monthly improvements of 1.8-3.1% and systematic bias detection with 89.7% accuracy in identifying algorithmic drift patterns [8]. The monitoring infrastructure captures

comprehensive performance data including processing latencies averaging 145 milliseconds for text analysis, 320 milliseconds for image analysis, and 1.2 seconds for video analysis, providing detailed operational insights that enable proactive optimization strategies.

## 5. Search Optimization and Ad Targeting Enhancement

Multi-modal categorization significantly enhances search functionality by enabling sophisticated product discovery driven by diverse query types and user intent patterns, with personalized and semantic retrieval systems demonstrating impressive performance improvements through end-to-end embedding learning strategies that achieve Mean Reciprocal Rank scores of 0.743 and Normalized Discounted Cumulative Gain values reaching 0.821 when evaluated on large-scale ecommerce datasets containing over 2.3 million products and 45 million user interaction records [9]. Visual search capabilities benefit users through advanced visual similarity matching algorithms comprehensive embedding embedded within frameworks that map visual product features into shared semantic spaces with visual search accuracy rates of 89.4% when processing product catalogs across broad categories, while maintaining query response times under 200 milliseconds through optimized indexing structures designed for efficient storage and retrieval. Natural language queries utilize sophisticated embedding learning techniques that capture rich semantic associations between user intent and product characteristics, where the system achieves query understanding improvements of 34.7% compared to traditional TF-IDF methods and demonstrates semantic matching accuracy of 91.8% when handling complex multi-attribute queries [9]. The enhanced categorization supports advanced semantic search capabilities where queries for concepts like "durable outdoor equipment" effectively identify based on contextual products embedding similarities rather than exact keyword matches, with semantic relevance scoring achieving 0.87 coefficients with human expert judgments and reducing average search result review time by 42.3% through improved ranking quality and result precision. Advanced advertising targeting capabilities leverage rich categorical information through multimodal customer satisfaction prediction models that demonstrate exceptional performance in understanding consumer preferences and purchasing patterns with

92.7% customer satisfaction prediction accuracy when processing aggregated textual reviews, visual product interactions, and behavioral engagement patterns across datasets containing over 1.8 million customer records [10]. The system utilizes sophisticated big data analytics techniques that process multimodal customer feedback streams, including text sentiment analysis with 89.3% accuracy in satisfaction classification, visual interaction pattern recognition with 91.7% accuracy in detecting preference indicators, and behavioral sequence analysis with 87.4% accuracy in predicting future purchase likelihood. comprehensive analytics capabilities enable product recommendations personalized with explicit user preferences derived from review sentiment showing 94.2% correlation to actual satisfaction ratings and implicit behavioral patterns captured from browsing history, producing recommendation relevance improvements of 36.8% over traditional collaborative filtering approaches [10]. Cross-modal insights reveal complex consumer satisfaction relationships spanning multiple interaction modalities, with statistical analysis indicating that customers exhibiting positive sentiment in textual reviews show 78.4% predictability in demonstrating sustained engagement behaviors, while visual interaction patterns accurately predict customer retention performance with 85.7% accuracy across diverse product categories. The advanced categorization framework enables dvnamic ad creative optimization through intelligent content selection algorithms that analyze multimodal customer satisfaction metrics to determine optimal advertising achieving 28.9% components, conversion rate improvements through personalized creative generation aligned with individual customer preference profiles [9]. Sophisticated personalization processes utilize embedding learning methodologies to create representations of customer preferences across textual, visual, and behavioral dimensions, with these integrated customer embeddings enabling targeted advertising precision improvements of 41.2% compared to traditional demographic targeting approaches. The system demonstrates real-time adaptation capabilities through continuous embedding updates that incorporate latest customer interaction data, with recommendation accuracy improvements of 2.8-4.1% per month through iterative learning processes that simultaneously optimize semantic relevance and customer satisfaction prediction [10].

 Table 1. Transformer Architecture Performance Metrics for Multi-Modal Processing [3, 4].

<b>Processing Component</b>	Parameter Count	Accuracy Rate (%)	<b>Processing Speed</b>	Performance Range
Text Processing Models	117M - 1.3B parameters	89.7	1,200-1,800 tokens/sec	Multi-language support
Vision Transformer (ViT-Base)	16×16 pixel patches	77.9 (ImageNet)	196 visual tokens	224×224 input images
Visual Processing	300M image training	94.2 color accuracy	Real-time inference	91.8-96.4 precision range
Cross-Modal Attention	8-16 attention heads	91.8 precision	100ms latency	Correlation: 0.73- 0.91
Video Processing	30 fps capability	93.2 segmentation	Temporal analysis	23-31% improvement
Semantic Processing	25+ languages	92% translation accuracy	Sub-100ms	Multi-modal fusion

**Table 2.** Content-Based Feature Extraction Performance Analysis [5, 6].

Feature Extraction	Accuracy Rate	Dataset	Processing	Quality Metrics	
Method	(%)	Coverage	Capability		
Named Entity	94.7	850K+ products	Brand/material	92.7-91.8 precision	
Recognition			extraction	1	
Fusion Sentiment	94.32	Multiple	Consumer	5 sentiment categories	
Analysis	94.32	dimensions	experience	3 semiment categories	
Visual Feature	96.4 color	47 material	Aesthetic analysis	2.1% tolerance	
Extraction	classification	categories	Aesthetic alialysis	2.1 70 tolerance	
Brand Logo Detection	95.7	Thousands of	Commercial	93.8 packaging	
		logos	recognition	accuracy	
Multi-Modal	8.2-12.7%	180K+ items	Late fusion	Category-specific	
Classification	improvement	180K+ Itellis	approach	rates	
Video Content	92.6 key moments	28 application	Temporal	88.9 scenario	
Analysis		categories	segmentation	recognition	
Quality Assessment	91.8 functional	Ground truth	Performance	3.8% measurement	
	accuracy	correlation	evaluation	accuracy	

 Table 3. Implementation Architecture Scalability Metrics [7, 8].

System Component	Processing Capacity	Response Performance	Scalability Features	Quality Maintenance
Lambda Architecture	2.5M transactions/hour	99.7% availability	Multi-agent coordination	4-6 hour batch processing
Real-time Processing	18,000 products/hour	120-180ms latency	Horizontal scaling	93.2% accuracy
Distributed System	45M+ item catalogs	280% throughput improvement	Agent clusters	Dynamic resource allocation
Quality Assessment	25+ file formats	99.1% transcoding success	Multi-format support	94.1% completeness detection
Batch Processing	35M products/12- 16hrs	26.8% consistency improvement	Large-scale updates	18.4% accuracy enhancement
Integration APIs	RESTful/GraphQL	97.2% quality preservation	Standardized interfaces	99.7% encoding conversion
Monitoring System	200+ performance metrics	365-day audit trails	Comprehensive tracking	89.7% bias detection

Table 4. Search Optimization and Ad Targeting Enhancement Results [9, 10].

Enhancement Feature	Performance Improvement	Accuracy Metrics	User Experience	<b>Business Impact</b>
Embedding Learning	MRR: 0.743	NDCG: 0.821	2.3M+ products	45M interaction records
Visual Search	89.4% accuracy	200ms response time	Image upload	Multi-category

			capability	processing
Semantic Matching	34.7% over TF-IDF	91.8% precision	Complex query handling	42.3% time reduction
Personalized Systems	36.8% relevance improvement	94.2% correlation	Individual preferences	Behavioral pattern analysis
Customer	92.7% prediction	1.8M+ customer	Multi-modal	89.3% sentiment
Satisfaction	accuracy	records	feedback	classification
Targeting	78.4% likelihood	85.7% retention	Cross-modal	Diverse product
Mechanisms	correlation	prediction	insights	categories
Dynamic	28.9% conversion	41.2% precision	Real-time	2.8-4.1% monthly
Optimization	improvement	enhancement	adaptation	gains

#### 6. Conclusions

Multi-modal artificial intelligence adoption in dynamic product categorization establishes a transformative paradigm in digital commerce that transcends the traditional limitations of manual classification systems. Transformer architectures demonstrate exceptional capability in processing multivariate data streams simultaneously, creating rich product understanding that integrates textual semantics, visual characteristics, and temporal demonstrations into unified representations. Advanced integration of natural language processing, computer vision, and video analysis techniques enables the extraction of comprehensive product features that significantly improve categorization accuracy and search capabilities. Commercial deployment strategies leveraging distributed processing architectures and quality assessment algorithms ensure scalable performance with the precision required for large-scale ecommerce operations. Semantic search capability improvements transform user experience by providing intuitive product discovery through natural language queries, visual similarity matching, and contextual understanding consumer intent. Personalized recommendation mechanisms utilize rich multi-modal insights to deliver targeted content that aligns with individual preferences and behavioral patterns, resulting in enhanced engagement and conversion rates. This technological advancement represents fundamental shift toward intelligent automation that continuously evolves with changing product landscapes and consumer expectations. Future developments in multi-modal categorization are expected to focus on enhanced cross-modal understanding, improved real-time processing capabilities, and deeper integration with emerging commerce technologies. The comprehensive framework establishes the foundation for nextgeneration e-commerce systems that prioritize user experience, operational efficiency, and commercial intelligent effectiveness through product organization and discovery mechanisms.

### **Author Statements:**

- **Ethical approval:** The conducted research is not related to either human or animal use.
- Conflict of interest: The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper
- **Acknowledgement:** The authors declare that they have nobody or no-company to acknowledge.
- **Author contributions:** The authors declare that they have equal right on this paper.
- **Funding information:** The authors declare that there is no funding to be acknowledged.
- Data availability statement: The data that support the findings of this study are available on request from the corresponding author. The data are not publicly available due to privacy or ethical restrictions.

#### References

- [1] Claudimar Pereira da Veiga et al., "E-Commerce in Brazil: An In-Depth Analysis of Digital Growth and Strategic Approaches for Online Retail," MDPI, 2024. [Online]. Available: <a href="https://www.mdpi.com/0718-1876/19/2/76">https://www.mdpi.com/0718-1876/19/2/76</a>
- [2] Gabriel de Souza P. Moreira et al., "Transformers with multi-modal features and post-fusion context for ecommerce session-based recommendation," arXiv, 2021. [Online]. Available: <a href="https://arxiv.org/pdf/2107.05124">https://arxiv.org/pdf/2107.05124</a>
- [3] Weiguo Feng et al., "Research on the construction and application of an intelligent tutoring system for English teaching based on a generative pre-training model," ScienceDirect, 2025. [Online]. Available: <a href="https://www.sciencedirect.com/science/article/pii/S">https://www.sciencedirect.com/science/article/pii/S</a> 277294192500050X
- [4] Alexey Dosovitskiy et al., "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale," ICLR, 2021. [Online]. Available: <a href="https://arxiv.org/pdf/2010.11929/1000">https://arxiv.org/pdf/2010.11929/1000</a>

- [5] HUAQIAN HE et al., "Exploring E-Commerce Product Experience Based on Fusion Sentiment Analysis Method," IEEE Access, 2022. [Online]. Available:
  - https://ieeexplore.ieee.org/stamp/stamp.jsp?arnumber=9919154
- [6] Ye Bi et al., "A Multimodal Late Fusion Model for E-Commerce Product Classification," arXiv, 2020. [Online]. Available: https://arxiv.org/pdf/2008.06179
- [7] Gautam Pal et al., "Multi-Agent Big-Data Lambda Architecture Model for E-Commerce Analytics," MDPI, 2018. [Online]. Available: <a href="https://www.mdpi.com/2306-5729/3/4/58">https://www.mdpi.com/2306-5729/3/4/58</a>
- [8] Yixuan Wu et al., "OmniFuse: A general modality fusion framework for multi-modality learning on low-quality medical data," ScienceDirect, 2025. [Online]. Available: <a href="https://www.sciencedirect.com/science/article/pii/S">https://www.sciencedirect.com/science/article/pii/S</a> 1566253524006687
- [9] Han Zhang et al., "Towards Personalized and Semantic Retrieval: An End-to-End Solution for Ecommerce Search via Embedding Learning," arXiv, 2020. [Online]. Available: <a href="https://arxiv.org/pdf/2006.02282">https://arxiv.org/pdf/2006.02282</a>
- [10] Xiaodong Zhang et al., "Research on Multimodal Prediction of E-Commerce Customer Satisfaction Driven by Big Data," MDPI, 2024. [Online].

  Available: <a href="https://www.mdpi.com/2076-3417/14/18/8181">https://www.mdpi.com/2076-3417/14/18/8181</a>