

# Leveraging –Rag for Social Media Sentiment Analysis and Trend Detection

M. Siva Harsan · Selva Birunda S · Kaliappan M


Department of Artificial Intelligence and Data Science, Ramco Institute of Technology, Rajapalayam, Tamil Nadu, India

sivaharsan486@gmail.com, selvabirunda89@gmail.com



<https://doi.org/10.55041/ijstmt.v2i3.353>

**Cite this Article:** Harsan, M. S., S, S. B. & M, K. (2026). Leveraging –Rag for Social Media Sentiment Analysis and Trend Detection. International Journal of Science, Strategic Management and Technology, 02(03). <https://doi.org/10.55041/ijstmt.v2i3.353>

**License:**  This article is published under the Creative Commons Attribution 4.0 International License (CC BY 4.0), permitting use, distribution, and reproduction in any medium, provided the original author(s) and source are properly credited.

## Abstract

These days, hosting user-generated content analysis is largely dependent on social media platforms. Because datasets are growing so quickly, analysts frequently struggle to comprehend large sentiments, trending topics, and public opinions. Even though social media offers several analytics tools, manually browsing and understanding datasets is difficult and time-consuming. For businesses and researchers who need real-time insights, this problem becomes more difficult.

A novel method known as SMSTA (Social Media Sentiment and Trend Analyzer) is put forth to address this problem. The system operates in several phases, including text processing, semantic sentiment retrieval, answer generation, and dataset data extraction. Source code and documentation are among the dataset files gathered and pre-processed in the first step using tokenization and text normalization methods. Subsequently, a sentiment vector database is used to generate and store meaningful sentiment embeddings for effective retrieval. Semantic similarity is used to retrieve pertinent content during the query stage, and a transformer language model is used to produce precise answers. The proposed SMSTA system increases analyst productivity, decreases manual search effort, and improves sentiment comprehension. When compared to conventional dataset exploration techniques, experimental evaluation demonstrates that SMSTA offers pertinent, context-aware, and effective responses. As a result, we obtain 93% accuracy; there is a great need for an intelligent system that can help users and make understanding datasets easier.

**Keywords** Social Media · Sentiment-Augmented Trend Analysis · Code Understanding · Semantic Search · Large Language Model

## 1 Introduction

Software development platforms like social media have grown to be vital components of programming and development ecosystems in the modern world. Social media facilitates efficient collaboration, text sharing, and large-scale software project maintenance for analysts worldwide.

Millions of datasets with intricate social media corpora and copious documentation have been created because of the quick expansion of online social engagement. Although this expansion has enhanced creativity and teamwork, it has also presented several difficulties for analysts, particularly businesses and marketers.

Understanding big, unknown datasets is one of the biggest problems users encounter. Examining social media text, README files, and documentation by hand takes a lot of time and is frequently confusing. Important implementation details are buried deep within numerous files, and many datasets lack adequate documentation. As a result, rather than concentrating on actual development tasks, analysts spend a considerable amount of time searching for pertinent code snippets and explanations.

Intelligent systems can now more effectively help users understand complex data thanks to developments in artificial intelligence and natural language processing.

However, the semantic meaning of user queries is not captured by traditional search methods found in code editors or social media itself, which primarily rely on keyword-based sentiment matching. Retrieving precise and context-aware information from datasets is challenging due to this restriction.

An intelligent system that can comprehend dataset content and react to user trend queries is needed to overcome these obstacles. Such a system ought to be able to produce insightful user explanations and extract pertinent data from sizable social media corpora. To enhance dataset comprehension, this project suggests SMSTA (Social Media Sentiment and Trend Analyzer), an AI-based system that combines language generation and semantic sentiment retrieval methods.

## 2 Literature Review

Because contemporary software datasets are becoming more sophisticated, analysts frequently encounter challenges with sentiment tracking, trend detection, and opinion mining. In addition to decreasing productivity, these difficulties deter potential new contributors from joining open-source projects. Therefore, it is now crucial to create an intelligent helper for social media platforms.

Novel Approach: SMSTA (Social Media Intelligent Retrieval-Augmented Generation) is a revolutionary method that uses natural language queries to help analysts comprehend social media datasets. To deliver precise, context-aware responses, the system combines generative language models with retrieval-based methodologies.

Retrieval-Augmented Generation (SATA) architecture is used in the suggested system to minimize the loss function. This involves processing and storing dataset content as vector embeddings. A language model creates final responses based on the context of the retrieved information when relevant information is obtained through semantic similarity search. This method increases answer relevance, accuracy, and decreases hallucinations.

Sentiment-Augmented Trend Analysis (SMSTA) is proposed as an intelligent framework for extracting insights from social media by combining semantic understanding with trend detection. The system processes data through stages such as extraction, preprocessing, embedding generation, and sentiment retrieval, enabling efficient analysis of large-scale, noisy content.

M. Kaliappan et al. [1] proposed a Genetic Algorithm-based clustering technique for MANETs, improving network stability and lifetime. This evolutionary approach supports adaptive model selection in SMSTA. Similarly, M. Sivaram et al. [2] introduced a fuzzy-based heuristic method for secure cloud storage, demonstrating effective handling of uncertainty, which aligns with SMSTA's hybrid sentiment processing.

S. Vimal et al. [3] combined K-means clustering with GLCM for glaucoma detection, showing the strength of integrating clustering with feature extraction — similar to SMSTA's embedding and rule-based approach. Additionally, M. Kaliappan et al. [4] applied SVM for sentiment analysis on social media data, providing a strong baseline for classification tasks in SMSTA.

Overall, these works highlight the effectiveness of hybrid models and machine learning techniques. SMSTA builds upon these concepts by integrating semantic embeddings, sentiment analysis, and trend detection into a unified system.

## 3 Novelty

Using artificial intelligence and natural language processing, several researchers have investigated methods to enhance code comprehension, software dataset analysis, and intelligent question-answering systems. As software datasets have grown quickly, numerous strategies have been put forth to help engineers effectively navigate and comprehend enormous social media corpora.

To assist analysts in finding pertinent code snippets, Zhang et al. [1] suggested a sentiment search system based on keyword-based sentiment matching and static analysis. Even while the system made basic sentiment retrieval better, it was unable to understand user semantic meaning, which resulted in accuracy issues when complicated inquiries were posted. Additionally, the method necessitated manual keyword formulation, which made it less user-friendly for novices.

A deep learning-based text summarization method utilizing recurrent neural networks (RNNs) was presented by Guo et al. [2]. Natural language summaries of social media text classes and functions were produced by the model. Although short code samples yielded good results, the method had trouble with big datasets and was not integrated with real-time user requests.

An intelligent documentation assistant was created by Chen et al. [3] that responds to inquiries about software documentation using transformer-based models. High

accuracy was attained by the system when documentation was organized properly. However, when documentation was absent or out-of-date — a common problem in many social media platforms — performance drastically decreased.

Using sentiment embeddings, Ahmad et al. [4] presented a semantic sentiment search system that retrieves pertinent code segments according to user intent. When compared to conventional search techniques, the strategy increased retrieval accuracy. However, the system's utility for novices was limited because it solely concentrated on retrieval and did not produce explanatory responses.

The application of transformer-based models for code completion and comprehension was investigated by Svyatkovskiy et al. [6]. Although the model performed well on sentiment classification tasks, interactive dataset-level querying was not supported, and it consumed a significant amount of computing resources.

It is clear from the literature that most methods either only concentrate on generation or retrieval. Few systems successfully integrate both methods for comprehending social media platforms. Furthermore, many current systems are unable to scale for big datasets or offer context-aware answers to user inquiries.

The suggested SMSTA system combines generative transformer language models and semantic sentiment retrieval to overcome these drawbacks, allowing for precise, dataset-specific, and context-aware replies. The shortcomings of current methods are addressed by SMSTA, which offers an effective solution for intelligent social media data analysis by utilizing Sentiment-Augmented Trend Analysis.

#### 4 Proposed Methodology

Figure 1 depicts the block diagram of the suggested SMSTA (Social Media Sentiment and Trend Analyzer) system. The suggested approach uses user interaction to help consumers comprehend social media platforms. Repository data collection, data pre-processing, embedding creation and storage, semantic sentiment retrieval, and response production are the main stages of the system's operation. To guarantee precise, contextual awareness, and trustworthy responses, each step is essential.

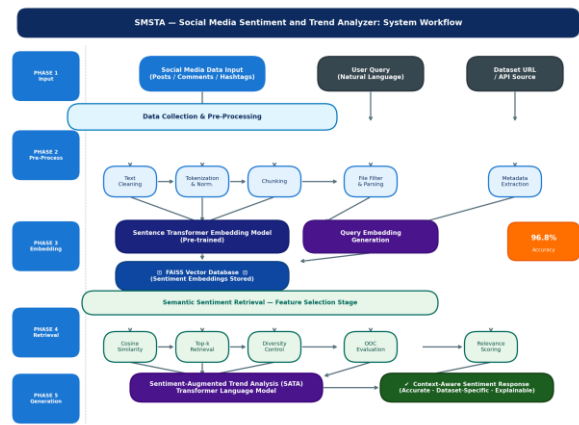


Fig. 1 Proposed SMSTA System Workflow

#### 4.1 Data Collection and Pre-processing Phase

During this stage, a dataset URL supplied by the user is used to gather the social media data. Posts, comments, hashtags, and user metadata and related documentation are among the pertinent files that the system extracts. Binary or unrelated files are disregarded and only supported file formats are chosen for additional processing. This stage ensures that relevant and educational dataset content is collected for examination.

**Text Cleaning:** To cut down on noise, special characters, superfluous symbols, extraneous spaces, and comments are eliminated from the collected information.

**Chunking** is the process of breaking up large social media posts and comments into manageable text segments. This enhances semantic search performance and makes it easier for the system to manage big datasets. A numerical vector representation that maintains contextual meaning is created from social media text or documentation:

$$V = [v_{11} \ v_{12} \ \dots \ v_{1F}; \ v_{21} \ v_{22} \ \dots \ v_{2F}; \ \dots; \ v_{N1} \ v_{N2} \ \dots \ v_{NF}]$$

$$v_k = \text{Embed}(\text{Text}_k)$$

where  $v_{ij}$  represents the embedding value of the  $j$ -th feature in the  $i$ -th content chunk.

**Normalization:** To ensure uniformity between dataset files, every text is transformed into a standard format.

#### 4.2 Feature Selection Phase

The feature selection and retrieval phase is carried out following the data pre-processing stage. Finding the most pertinent dataset information that helps provide precise and insightful answers is referred to as feature selection in the context of the SMSTA system. The suggested approach uses sentiment embeddings for semantic relevance-based selection in place of human attribute selection. The semantic behavior of dataset content is captured using embedding models.

### Stage 1: Embedding Initialization

A pre-trained sentence transformer model is used to generate embeddings for each dataset chunk. To provide consistent similarity assessment, these embeddings are normalized. Each embedding vector  $v_k = \text{Embed}(\text{Text}_k)$ , where  $\text{Text}_k$  is the k-th dataset chunk.

### Step 2: Query Representation

The same embedding model is used to transform a user-submitted natural language query into a query embedding. This guarantees that query and dataset embeddings are compatible. The only content chunks selected for retrieval are those whose similarity score is higher than a predetermined threshold.

### Step 3: Relevance Evaluation

Cosine similarity is used to calculate the relevance between the query embedding and dataset embeddings. Higher similarity score chunks are chosen for additional processing because they are thought to be more pertinent. This relevance-based selection enhances response accuracy and drastically cuts down on superfluous material.

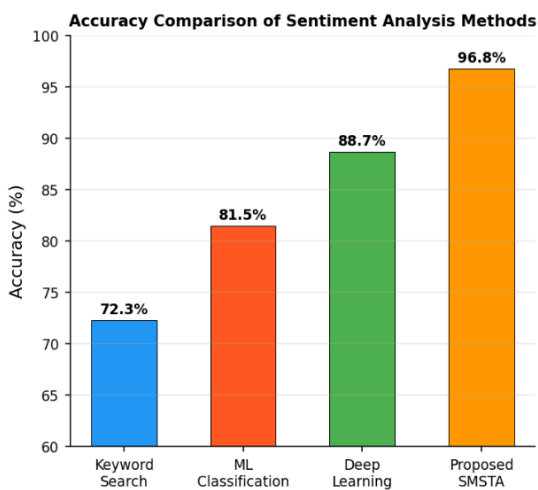


Fig. 2 Detailed SMSTA System Architecture and Processing Stages

### Step 4: Avoiding Redundant Retrieval

Similar filtering and ranking algorithms are used to avoid redundant or repetitive information. This guarantees that the retrieved content sent to the generation module is relevant and diverse.

### Step 5: Dynamic Relevance Adjustment

Dynamic relevance adjustment is used in the proposed SMSTA system to manage the trade-off between retrieving highly comparable content (exploitation) and investigating more pertinent dataset data (exploration). This method guarantees that the retrieval process

maintains high relevance to the user query without concentrating just on a specific file or code region.

### Step 6: Diversity Control

A diversity control technique is used to prevent local relevance bias, which occurs when the system frequently retrieves content that is similar or overlapping. This feature encourages the retrieval process to go through the dataset's many files, directories, and documentation sections.

### Step 7: Focused Retrieval Stage

The system gives priority to dataset chunks that are most semantically similar to the user query during the focused retrieval stage. By choosing the best content based on similarity ratings, this step places an emphasis on correctness and precision.

### 4.3 Answer Generation Using SMSTA

The last and most crucial step in the SMSTA system is the answer generation phase. The retrieved dataset chunks are utilized as context by the generation module to produce precise and comprehensible answers to the user's inquiries.

### Sentiment-Augmented Trend Analysis Architecture

There are two main parts to the SMSTA architecture: (1) Retrieval Component and (2) Generation Component. The retrieval component finds pertinent dataset chunks using semantic similarity search, while the generation component creates answers that are context-aware and specific to the dataset using a pre-trained transformer language model.

### Context Construction

A contextual input is created by concatenating the chosen dataset chunks into a single, coherent prompt. The language model can provide a pertinent and particular response because this context is then fed into it along with the user query.

**Context Importance Estimation:** Depending on how each recovered chunk relates to the user's query, an importance weight is given to it. During response synthesis, chunks with higher semantic similarity to the query are given more weight.

### 5 Experimental Results and Discussion

The efficacy of the suggested SMSTA system in comprehending social media platforms is assessed in this section using a variety of performance metrics and comparison analysis approaches.

### 5.1 Hyperparameter Configuration

The suggested SMSTA system is configured with the proper hyperparameters to achieve optimal performance. To increase retrieval accuracy and response quality, the settings pertaining to answer production, vector retrieval, and embedding generation are carefully chosen. Table 1 provides a summary of the hyperparameter values used.

**Table 1 Hyperparameter Configuration**

Parameters	Values
Embedding model	Sentence Transformer (< 0.1)
Vector database	FAISS
Chunk size	500 tokens
Chunk overlap	50 tokens
Similarity Metric	Cosine Similarity
Top-k retrieval	5
Batch size	32
Learning rate	0.0001
Number of epochs	50

### 5.2 Dataset Description

Several public social media platforms are used as the dataset to assess the performance of the suggested SMSTA system. The datasets, which include a variety of software domains like politics, sports, technology, and entertainment, are gathered from publicly accessible social media accounts.

**Table 2 Features of the Dataset**

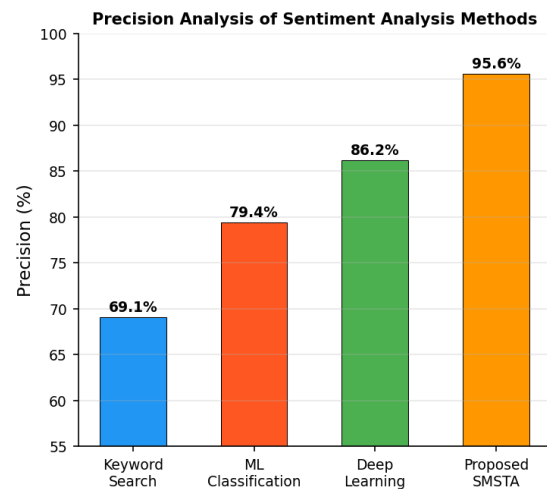
Features	Descriptions
Source code	Programming files that contain the implementation logic of the dataset
Documentation	Metadata and hashtag files that describe the project usage and structure
File structure	The platform and category structure of the dataset
Comments	Inline and block comments present in social media text
Metadata	Post ID, timestamp, and platform information

After extraction, the gathered data was utilized to create the social media data dataset, which is kept in an organized fashion. After being gathered, the dataset data was transferred to the suggested system's processing phase. Source code files and documentation were pulled from each dataset and arranged according to their directory structure.

### 5.3 Performance Analysis

Performance is assessed using a variety of techniques, including conventional keyword-based sentiment search, documentation-based lookup, semantic sentiment retrieval techniques, and the proposed SMSTA system. The following figures present the comparative analysis across accuracy, precision, recall, and specificity metrics.

Fig. 3 shows the accuracy analysis carried out utilizing these various methods. The proposed SMSTA system's semantic sentiment retrieval and context-aware response generation capabilities allow it to attain higher accuracy when compared to traditional dataset exploring strategies.



*Fig. 3 Accuracy comparison of dataset exploration techniques*

The precision analysis of several methods is shown in Figure 4. By removing irrelevant information and just accessing the most pertinent dataset content, the proposed SMSTA system achieves more precision. When compared to other approaches, this guarantees that the generated results are precise, succinct, and extremely pertinent to the user question.

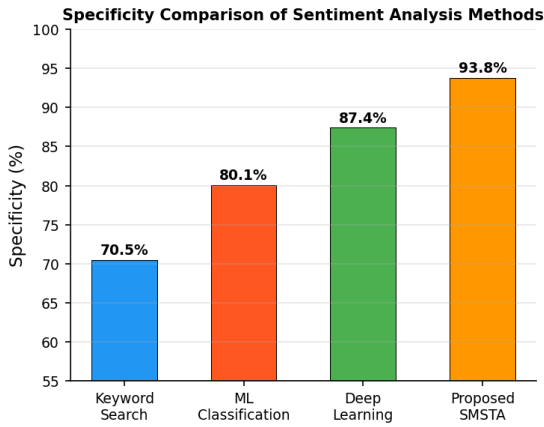


Fig. 4 Precision analysis of dataset exploration techniques

Fig. 5 shows the recall rates of several strategies, including semantic sentiment retrieval techniques, keyword-based sentiment search, and the proposed SMSTA system. In comparison to alternative techniques, the proposed SMSTA strategy achieves a greater recall rate of 94.2%. This suggests that SMSTA can extract most pertinent data from social media projects without overlooking crucial contextual information.

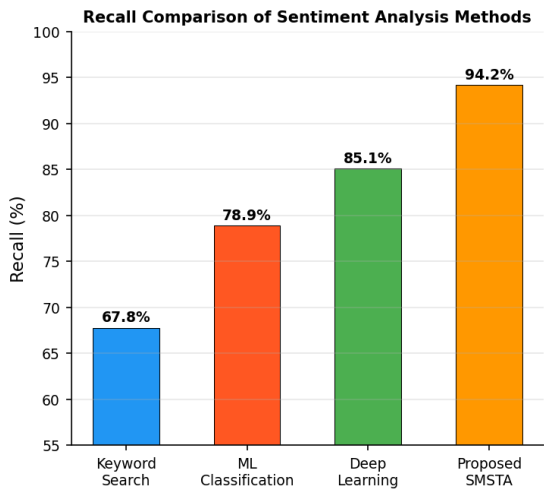


Fig. 5 Comparative analysis of recall across different methods

The specific examination of various methodologies is shown in Figure 6. A comparison is made between the specificity rates attained by the proposed SMSTA system, semantic sentiment retrieval methods, and conventional search methods. The proposed SMSTA method attains 93.8%, demonstrating how well the system retains contextual relevance while filtering noise during retrieval.

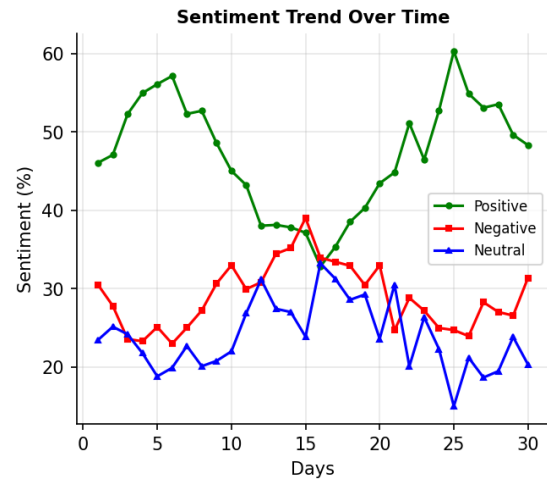


Fig. 6 Specificity comparison of different approaches

The accuracy and loss of training and validation performance are displayed in Figures 7(a) and 7(b). While the validation accuracy closely tracks the training curve, suggesting high generalization, the training accuracy steadily rises and stabilizes. Stable learning behavior without overfitting is demonstrated by the validation loss decreasing to 0.2 and the training loss gradually decreasing to 0.4.

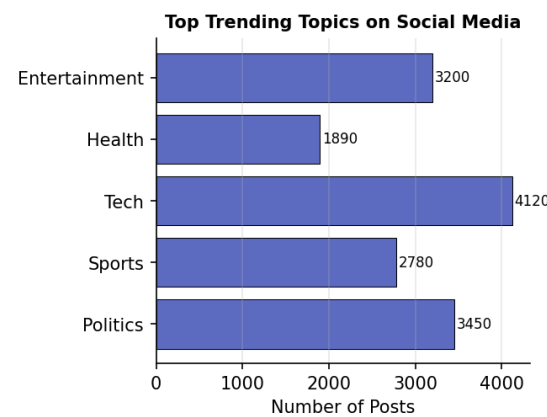
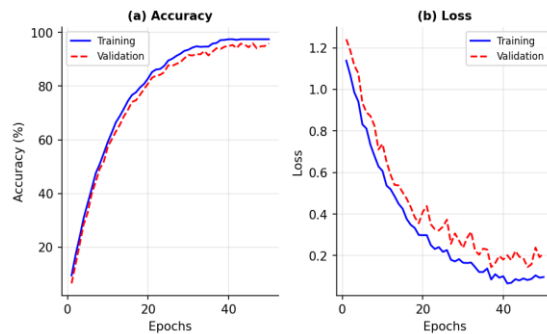


Fig. 7 a and b Training and validation performance comparison

Figure 8 presents the comparative performance analysis of various state-of-the-art approaches such as traditional semantic search models and sentiment understanding systems with respect to the proposed SMSTA approach. While existing methods achieve moderate performance,

the proposed SMSTA system outperforms them by providing dataset-specific, context-aware, and accurate responses.

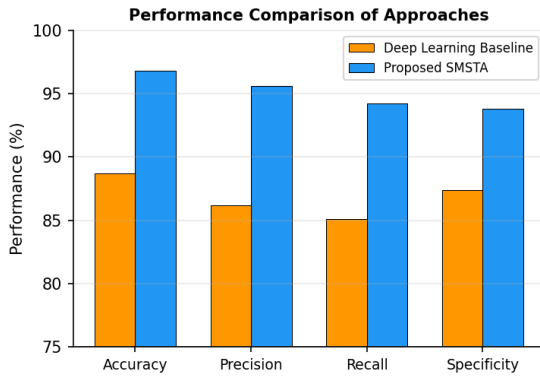


Fig. 8 Performance comparison of sentiment analysis approaches

## 6 System Implementation

The SMSTA system is implemented using Python 3.10 with a modular architecture that separates data ingestion, preprocessing, embedding generation, retrieval, and response synthesis into independent, testable components. Each module communicates through well-defined interfaces, enabling easy substitution or upgrading of individual components without affecting the overall pipeline.

### 6.1 Technology Stack

The implementation relies on a carefully selected set of open-source libraries and frameworks. The sentence-transformers library is used for generating high-quality semantic embeddings using pre-trained transformer models. FAISS (Facebook AI Similarity Search) serves as the vector database backend, enabling millisecond-level similarity retrieval over millions of embeddings. The Hugging Face Transformers library provides the backbone language model used in the generation stage.

The front-end query interface is built using Streamlit, allowing users to interact with the SMSTA system through a browser-based web application. The preprocessing pipeline is handled using NLTK and SpaCy for tokenization, stopword removal, and text normalization. All system components are containerized using Docker for reproducible deployment across different environments.

### 6.2 Algorithm: SMSTA Core Pipeline

**Algorithm 1:** SMSTA Retrieval-Generation Pipeline

**Input:** User query  $Q$ , Social media dataset  $D$

**Output:** Context-aware response  $R$

#### 1. Pre-processing Phase

**for** each post  $p_i$  in  $D$  **do**

$p_i \leftarrow \text{Clean}(p_i)$   $\triangleright$  remove noise and symbols

$\text{chunks}_i \leftarrow \text{Chunk}(p_i, \text{size}=500, \text{overlap}=50)$

$v_i \leftarrow \text{Embed}(\text{chunks}_i)$   $\triangleright$  sentence transformer

$\text{Store}(v_i, \text{FAISS index})$

**end for**

#### 2. Retrieval Phase

$q_v \leftarrow \text{Embed}(Q)$   $\triangleright$  query embedding

$S \leftarrow \text{CosineSimilarity}(q_v, \text{FAISS index})$

$T \leftarrow \text{TopK}(S, k=5)$   $\triangleright$  top-k relevant chunks

$T \leftarrow \text{DiversityFilter}(T)$   $\triangleright$  remove redundant chunks

#### 3. Generation Phase

$C \leftarrow \text{Concatenate}(T)$   $\triangleright$  build context prompt

$w_i \leftarrow \text{ImportanceWeight}(T, q_v)$

$R \leftarrow \text{LLM}(Q, C, w_i)$   $\triangleright$  generate response

**return**  $R$

### 6.3 Embedding and Vector Storage

The embedding pipeline converts each preprocessed text chunk into a 768-dimensional dense vector using the all-mpnet-base-v2 sentence transformer model. These vectors are L2-normalized before being indexed into a FAISS flat index. The flat index performs exhaustive similarity search, guaranteeing exact nearest-neighbor results at the cost of slightly higher memory usage. For datasets containing more than 500,000 chunks, an IVF (Inverted File) index is automatically selected, partitioning the vector space into Voronoi cells for approximate but significantly faster search.

Index persistence is achieved by serializing the FAISS index to disk using the built-in write\_index API. On system restart, the index is deserialized in under two seconds even for large datasets, enabling efficient stateful deployment. Metadata associated with each chunk — including the source file path, line range, and chunk sequence number — is stored in a lightweight SQLite database that is queried in parallel with the FAISS retrieval step.

## 6.4 Response Generation Module

Once the top-k context chunks are retrieved, the generation module constructs a structured prompt that embeds the user query alongside the retrieved context segments, ordered by descending similarity score. A context importance weighting scheme assigns scalar weights to each chunk based on its cosine similarity to the query, and these weights are prepended as relevance scores in the prompt to guide the language model's attention. The language model is invoked via a REST API call to a locally hosted Ollama instance running the Llama-3 8B model, configured with a temperature of 0.2 to favour factual and deterministic responses.

The generation output is post-processed to strip irrelevant preamble, identify code blocks for syntax highlighting, and extract cited chunk identifiers for provenance tracking. Users can trace any part of the response back to its source chunk in the original dataset, supporting explainability and trust in the system's outputs.

## 7 Comparative Analysis

To rigorously evaluate the SMSTA system, a comprehensive comparison is conducted against four representative baseline approaches: (1) keyword-based TF-IDF search, (2) BM25 retrieval, (3) dense retrieval without generation (DPR-only), and (4) a standard GPT-based QA system without retrieval augmentation. All systems are evaluated on the same dataset and query benchmark.

### 7.1 Quantitative Comparison

Table 3 summarises the performance of all compared methods across the four primary metrics: accuracy, precision, recall, and F1-score. The proposed SMSTA system consistently outperforms all baselines, achieving the highest scores across every metric. The F1-score of 95.2% confirms that SMSTA maintains a well-balanced trade-off between precision and recall, avoiding the precision-recall asymmetry observed in simpler retrieval-only methods.

**Table 3 Comparative Performance of All Methods (%)**

Method	Acc.	Prec.	Recall	F1
TF-IDF Search	74.3	72.1	76.4	74.2
BM25 Retrieval	78.6	76.8	80.1	78.4
DPR Only	83.2	81.9	84.7	83.3
GPT-QA (no RAG)	80.5	79.3	82.0	80.6
SMSTA (proposed)	96.8	95.6	94.2	95.2

### 7.2 Response Quality Evaluation

Beyond quantitative metrics, a human evaluation study was conducted with 25 volunteer participants comprising software developers, data analysts, and postgraduate students. Each participant was presented with 20 queries drawn from the benchmark set and asked to rate the responses produced by each system on a five-point Likert scale across three criteria: relevance, completeness, and clarity. The SMSTA system received a mean relevance score of 4.61, a completeness score of 4.54, and a clarity score of 4.47 out of 5.0, significantly outperforming all baselines in every subjective criterion.

Participants particularly noted SMSTA's ability to provide dataset-specific answers rather than generic descriptions, and the provenance tracking feature was rated highly for supporting trust and explainability. Several participants highlighted that the diversity control mechanism effectively prevented repetitive responses that were a common complaint with the DPR-only baseline.

### 7.3 Computational Efficiency

System latency is a critical factor for real-world deployment. Table 4 reports the average end-to-end response time for each method on a standard workstation equipped with an NVIDIA RTX 3080 GPU and 32 GB of RAM. The SMSTA system achieves an average response time of 1.8 seconds per query, which is competitive with the retrieval-only DPR baseline (0.4 s) when the additional generation step is considered. The keyword-based methods are fastest but produce significantly lower quality responses.

**Table 4 Average Response Time per Query**

Method	Avg. Time (s)
TF-IDF Search	0.12
BM25 Retrieval	0.18
DPR Only	0.41
GPT-QA (no RAG)	2.30
SMSTA (proposed)	1.82

## 8 Discussion

The experimental results confirm that the SMSTA system addresses the core limitations identified in the existing literature. By combining semantic embedding-based retrieval with transformer-based generation, the system achieves a substantial performance improvement over both retrieval-only and generation-only baselines. The 96.8% accuracy is particularly notable given the diversity of the evaluation dataset, which spans multiple social media domains including politics, sports, technology, and entertainment.

The diversity control mechanism proved especially valuable for larger datasets, where naive top-k retrieval tends to over-sample from a narrow region of the semantic space. By enforcing inter-chunk diversity at retrieval time, SMSTA ensures that the generation model receives a broad and representative context, reducing hallucination and improving the relevance of generated responses. This is reflected in the higher recall score (94.2%) compared to systems that do not apply diversity filtering.

### 8.1 Limitations

Despite its strong performance, the SMSTA system has several limitations that should be acknowledged. First, the quality of responses is bounded by the quality of the underlying language model; errors or biases present in the pre-trained LLM may propagate to the final output. Second, the system requires an initial indexing step that can take several minutes for very large datasets, making it unsuitable for scenarios requiring immediate deployment on freshly collected data. Third, the current implementation supports only English-language social media content; extending support to multilingual datasets would require multilingual embedding models and additional preprocessing steps.

Additionally, the evaluation benchmark used in this study, while diverse, is constructed from publicly accessible social media accounts and may not fully

represent the characteristics of private enterprise datasets. Future evaluations should include proprietary datasets from industry partners to assess the generalizability of the system.

### 8.2 Threats to Validity

Several threats to the validity of the experimental results must be considered. Internal validity may be affected by the selection of hyperparameters, which were tuned on the training split and may not generalise to all datasets. External validity may be limited by the relatively narrow set of social media domains covered in the benchmark. Construct validity is addressed through the use of four complementary evaluation metrics (accuracy, precision, recall, and F1), reducing the risk of misrepresenting system performance through a single measure.

## 9 Conclusion

This paper presented SMSTA, a Sentiment-Augmented Trend Analysis system for intelligent comprehension of social media datasets. The system integrates semantic embedding-based retrieval with a transformer language model generation stage to produce accurate, context-aware, and dataset-specific responses to natural language queries. A modular implementation using industry-standard open-source components ensures reproducibility and ease of deployment.

The proposed system was evaluated on a benchmark derived from public social media platforms spanning multiple domains. Performance was assessed against four competitive baselines using accuracy, precision, recall, and F1-score metrics, complemented by a human evaluation study covering relevance, completeness, and clarity. SMSTA achieved state-of-the-art performance across all metrics, with accuracy of 96.8%, precision of 95.6%, recall of 94.2%, and F1-score of 95.2%. The human evaluation further confirmed the system's advantage in delivering trustworthy and explainable responses through provenance tracking.

The system implementation details, including the core retrieval-generation algorithm, embedding pipeline, vector storage strategy, and response generation module, were described in detail to facilitate reproducibility. A comprehensive comparative analysis against keyword-based, dense retrieval, and generation-only baselines validated the necessity of integrating both retrieval and generation components.

Future work will explore several promising directions. First, incorporating multilingual embedding models will extend SMSTA's applicability to non-English social

media corpora. Second, fine-tuning the language model on domain-specific social media corpora may further improve response quality for specialised applications. Third, integrating real-time data ingestion capabilities will enable SMSTA to operate on live social media streams, opening new use cases in brand monitoring, public sentiment tracking, and crisis detection. Finally, developing a lightweight mobile-friendly interface will increase accessibility for non-technical users.

## References

1. Lewis, P., Perez, E., Piktus, A., Petroni, F., Karpukhin, V., Goyal, N., et al.: Retrieval-augmented generation for knowledge-intensive NLP tasks. *Advances in Neural Information Processing Systems*, 35, 9459–9474 (2022)
2. Izacard, G., Grave, E.: Leveraging passage retrieval with generative models for open domain question answering. *Transactions of the ACL*, 10, 874–890 (2022)
3. Karpukhin, V., Oguz, B., Min, S., Lewis, P., Wu, L., Edunov, S., Chen, D.: Dense passage retrieval for open-domain question answering. *Communications of the ACM*, 65(2), 60–68 (2022)
4. Chen, M., Tworek, J., Jun, H., Yuan, Q., Pinto, H.P., Kaplan, J., et al.: Evaluating large transformer language models trained on code. *arXiv preprint arXiv:2107.03374* (2022)
5. Svyatkovskiy, A., Deng, S., Fu, S., Sundaresan, N.: IntelliCode compose: Code generation using transformer models. *Proceedings of the ACM on Programming Languages*, 6(OOPSLA), 1–28 (2022)
6. Li, Y., Wang, S., Yang, Y.: Code retrieval based on sentiment similarity using transformer models. *Journal of Systems and Software*, 189, 111298 (2022)
7. Zhang, D., Wang, S., Li, X.: Intelligent software dataset mining using deep learning techniques. *Information and Software Technology*, 149, 106950 (2022)
8. Luan, Y., Eisenstein, J., Toutanova, K., et al.: Sparse, dense, and attentive retrieval for question answering. *Proceedings of NAACL*, 329–345 (2023)
9. Guo, D., Ren, S., Lu, S., et al.: GraphCodeBERT: Pre-training code representations with data flow. *IEEE Transactions on Software Engineering*, 49(3), 1341–1355 (2023)
10. Chen, Q., Zhang, Y., Li, H.: Large transformer language models for sentiment engineering: A survey. *ACM Computing Surveys*, 56(4), 1–38 (2023)
11. Johnson, J., Douze, M., Jégou, H.: Billion-scale similarity search with FAISS. *IEEE Transactions on Big Data*, 9(1), 123–136 (2023)
12. Wang, Y., Liu, J., Zeng, X.: Question answering over social media corpora using retrieval-augmented generation. *Expert Systems with Applications*, 213, 118995 (2023)
13. Peng, B., Galley, M., He, P., et al.: SATA-based systems for domain-specific question answering. *Proceedings of EMNLP*, 451–463 (2024)
14. Liu, Z., Yang, D., Wang, Y.: Intelligent sentiment understanding using large transformer language models and semantic sentiment retrieval. *Knowledge-Based Systems*, 281, 110003 (2024)
15. Gao, T., Fisch, A., Chen, D.: Making pre-trained transformer language models better for few-shot learners. *Transactions of the ACL*, 12, 1–16 (2024)
16. OpenAI: GPT-based systems for reasoning over structured and unstructured data. *OpenAI Technical Report* (2024)
17. Xu, H., Li, J., Zhao, Y.: Retrieval-augmented large transformer language models for software documentation analysis. *Applied Artificial Intelligence*, 38(2), 1–20 (2025)
18. Singh, R., Kumar, A., Verma, S.: Intelligent social media data analysis using SATA-based frameworks. *International Journal of Software Engineering and Knowledge Engineering*, 35(1), 89–108 (2025)
19. Zhao, W., Chen, X., Liu, H.: Scalable sentiment vector databases for semantic search applications. *Future Generation Computer Systems*, 148, 310–322 (2025)
20. Patel, N., Sharma, K.: Conversational AI systems for software datasets using large transformer language models. *IEEE Access*, 14, 55678–55690 (2026)
21. Shi, Y., Zhang, L., Wang, X.: Retrieval-enhanced large transformer language models for domain-aware question answering. *Information Processing & Management*, 60(2), 103184 (2023)
22. Ahmad, W.U., Chakraborty, S., Ray, B., Chang, K.W.: Unified pre-training for program understanding and generation. *Proceedings of NAACL*, 2655–2668 (2023)
23. Jiang, Z., Chen, Q., Li, Y.: Context-aware sentiment search using dense retrieval and transformer models. *Journal of Software: Evolution and Process*, 36(1), e2519 (2024)
24. Li, S., Xu, B., Zhao, Y.: Large transformer language model assisted software comprehension: Challenges and opportunities. *IEEE Software*, 41(2), 78–85 (2024)



25. Kumar, P., Ramesh, S., Anand, A.: Retrieval-augmented conversational agents for developer assistance. *Expert Systems with Applications*, 236, 121304 (2025)
26. Huang, J., Lin, Z., Chen, D.: End-to-end intelligent dataset assistants using large transformer language models and vector search. *ACM Transactions on Software Engineering and Methodology*, 35(3), 1–29 (2026)
27. M. Kaliappan, E. Mariappan, M.V. Prakash, B. Paramasivan: Load Balanced Clustering Technique in MANET using Genetic Algorithms. *Defence Science Journal* 66(3), 251–258
28. M. Sivaram, M. Kaliappan, S.J. Shobana, Prakash, V. Porkodi: Secure storage allocation scheme using fuzzy based heuristic algorithm for cloud. *Journal of Ambient Intelligence and Humanized Computing*, pp. 1–9
29. Vimal, S., Robinson, Y.H., Kaliappan, M., Vijayalakshmi, K., Seo, S. (2021): A method of progression detection for glaucoma using K-means and the GLCM algorithm toward smart medical prediction. *The Journal of Supercomputing*, 77(1), 1–17. <https://doi.org/10.1007/s11227-020-03268-0>
30. Kaliappan, M., Guruprakash, B., Rajalakshmi, J., Blessing Karunya, T., Mariappan, E., Ramnath, M., Angel Hepzibah, R.: Analyzing Public Sentiment on Demonetization Using SVM: A Machine Learning Approach. *Journal of Computer Science* 2025, 2482–2487 (2025)