

# Research on Intelligent Public Opinion Monitoring Based on Multi-source Heterogeneous Data Stream Fusion and Deep Semantic Analysis

Huang Mingqi

January 16, 2026

## Abstract

With the diversification of social media, public opinion analysis faces challenges such as multi-source heterogeneity, semantic complexity, and real-time requirements. This paper proposes an intelligent public opinion monitoring method based on multi-source heterogeneous data stream fusion and deep semantic analysis. First, we design a unified data representation model and an adaptive collection framework to achieve real-time alignment and fusion of multi-platform heterogeneous data. Second, we construct a hybrid analysis model that combines TF-IDF-enhanced K-Means topic clustering, rule and dictionary-enhanced sentiment analysis, and introduces the Pangu Embedded Large Language Model (Pangu-Embedded-7B) as the core for semantic understanding and decision-making, enabling accurate parsing of Chinese long texts, complex emotions, and implicit semantics. Furthermore, we propose a vector retrieval-based semantic enhancement method to improve the recall and accuracy of related topic matching. Experiments show that our method achieves a topic clustering ARI of 0.71 and an F1-score of 0.78 for sentiment analysis on multiple public datasets, representing a 9.3% improvement over mainstream baseline models, with good generalization capability and real-time performance. Finally, based on this method, we implement a complete public opinion analysis system that supports multi-channel interaction and automated report generation, providing a reliable decision support tool for governments and enterprises.

**Code Availability:** The source code of this system is publicly available at:

<https://github.com/hmmnxkl/LLM-Based-Intelligent-Public-Opinion-Analytics-Assistant>

**Keywords:** Public opinion analysis, Multi-source data fusion, Topic clustering, Sentiment analysis, Large language model, Pangu model, Hybrid analysis

## 1 Introduction

### 1.1 Research Background and Significance

In the era of social media proliferation, online public opinion has become a critical barometer of social dynamics, public sentiment, and potential risks. The current public opinion ecosystem exhibits four distinct characteristics: **large data volume, multiple source platforms, rapid generation rate, and heterogeneous content forms**. For government agencies, timely and accurate grasp of public opinion trends is essential for social governance and crisis management; for enterprises, public opinion analysis is crucial for brand reputation management and market competition analysis; for research institutions, public opinion data provides valuable resources for observing social mentality and communication patterns.

---

However, traditional public opinion monitoring methods face three major challenges: (1) **Insufficient data fusion capability**: inability to effectively integrate heterogeneous data from multiple platforms; (2) **Limited semantic understanding**: difficulty in capturing complex emotional expressions, implicit semantics, and contextual relationships; (3) **Lack of adaptive intelligence**: heavy reliance on manual rule configuration, unable to achieve automated event discovery and trend prediction.

To address these challenges, this paper focuses on three core research questions:

1. How to achieve unified representation and real-time fusion of multi-source heterogeneous public opinion data?
2. How to combine traditional text mining methods with large language models to enhance semantic understanding and analytical capabilities?
3. How to construct an end-to-end intelligent public opinion analysis framework that balances accuracy, efficiency, and interpretability?

## 1.2 Related Work

### 1.2.1 Multi-source Data Fusion in Public Opinion Analysis

Early research on public opinion data collection primarily focused on single-platform crawlers [1]. With the diversification of platforms, research has shifted toward multi-source data integration. However, most existing works lack systematic methods for handling structural heterogeneity and semantic alignment across different platforms. Recent studies have explored knowledge graph-based fusion methods [3], but their real-time performance is inadequate for dynamic public opinion scenarios.

### 1.2.2 Topic Discovery and Sentiment Analysis Methods

Traditional topic modeling approaches such as LDA have been widely used but perform poorly on short texts. Recent methods based on pre-trained language models (e.g., BERT) combined with clustering algorithms have shown improvements, but they require substantial computational resources and lack interpretability. For sentiment analysis, methods have evolved from lexicon-based approaches [2] to deep learning models. However, these methods struggle with Chinese-specific expressions such as sarcasm and implicit emotions.

### 1.2.3 Large Language Models in Public Opinion Analysis

The application of large language models (LLMs) in public opinion analysis is still in its early stages. Some studies have explored using GPT-series models for sentiment classification, but they face challenges with Chinese context understanding and data privacy concerns. The Pangu model series has demonstrated advantages in Chinese language processing and local deployment, making it suitable for sensitive public opinion analysis scenarios.

### 1.2.4 Existing Public Opinion Systems

Commercial systems such as Qingbo Big Data and Shangmiao provide basic monitoring functions but lack deep analytical capabilities and customizability. Academic prototype systems often focus on specific aspects (e.g., event detection or sentiment analysis) without providing end-to-end solutions.

---

### 1.3 Main Contributions

This paper makes the following contributions:

1. **Propose a multi-source heterogeneous data fusion framework:** We design a unified data representation model and an adaptive collection mechanism that enables real-time integration of data from over ten mainstream platforms while maintaining data consistency and integrity.
2. **Develop a hybrid public opinion analysis model:** We innovatively combine traditional text mining algorithms (TF-IDF, K-Means) with the Pangu-7B large language model, creating a synergistic analysis pipeline that leverages the strengths of both approaches for topic clustering and sentiment analysis.
3. **Establish a semantic enhancement retrieval method:** We propose a vector retrieval approach with semantic enhancement, improving the accuracy of related topic matching through text enhancement and multi-context fusion.
4. **Conduct comprehensive experimental validation:** We perform extensive experiments on multiple datasets, including comparative analysis with existing systems, ablation studies, and generalization tests, demonstrating the effectiveness and superiority of our approach.
5. **Implement a complete system and provide practical insights:** We develop a fully functional public opinion analysis system and share implementation details and deployment experiences.

### 1.4 Paper Structure

The remainder of this paper is organized as follows: Section 2 reviews related work. Section 3 presents the theoretical foundation and model formulation. Section 4 describes the multi-source data fusion method. Section 5 details the hybrid public opinion analysis model. Section 6 presents experimental design and results. Section 7 discusses system implementation and applications. Section 8 concludes with future research directions.

## 2 Related Work

This section provides a comprehensive review of existing research in multi-source data fusion, public opinion analysis methods, and large language model applications.

### 2.1 Multi-source Data Collection and Fusion

Traditional web crawlers have evolved from general-purpose to platform-specific implementations. The Scrapy framework provides a robust foundation for building distributed crawler systems. However, most existing crawlers focus on single platforms, with limited research on cross-platform data fusion. Recent studies have explored ontology-based integration and graph-based alignment methods, but these approaches often lack real-time performance for dynamic public opinion scenarios.

### 2.2 Topic Modeling and Clustering Algorithms

Latent Dirichlet Allocation (LDA) has been the dominant approach for topic modeling but performs poorly on short texts. Recent advancements include BERT-based clustering and neural topic models. However, these methods often require extensive computational resources and lack interpretability. Our approach combines TF-IDF with K-Means, enhanced by semantic features from the Pangu model, achieving a balance between efficiency and effectiveness.

---

### 2.3 Sentiment Analysis Methods

Sentiment analysis has evolved through three stages: lexicon-based methods, machine learning approaches (SVM, Naive Bayes), and deep learning models (LSTM, Transformer). For Chinese text, specific challenges include handling implicit emotions, sarcasm, and cultural contexts. Recent work has explored fine-tuning pre-trained models for Chinese sentiment analysis, but these methods require substantial labeled data.

### 2.4 Large Language Models in Text Analysis

The success of Transformer-based models has revolutionized natural language processing. The Pangu model series has shown particular strength in Chinese language understanding and generation tasks. Unlike general-purpose LLMs, Pangu models are optimized for Chinese contexts and support efficient local deployment, making them suitable for privacy-sensitive applications like public opinion analysis.

### 2.5 Public Opinion Analysis Systems

Commercial systems (e.g., Qingbo, Shangmiao) focus primarily on data collection and basic visualization. Academic systems often emphasize specific analytical components without providing complete solutions. Our work bridges this gap by integrating multi-source data collection, advanced analysis algorithms, and practical system implementation.

## 3 Theoretical Foundation and Model Formulation

### 3.1 Problem Definition

Let  $\mathcal{P} = \{p_1, p_2, \dots, p_N\}$  represent  $N$  data sources (platforms). Each platform  $p_i$  generates a data stream  $D_i = \{d_{i1}, d_{i2}, \dots\}$ , where each data item  $d_{ij}$  contains multiple attributes. The goal is to develop a function  $f : \bigcup_{i=1}^N D_i \rightarrow \mathcal{R}$  that maps heterogeneous data items to a unified representation space  $\mathcal{R}$ , and then perform analysis tasks including topic clustering and sentiment analysis.

### 3.2 Unified Data Representation Model

We define a unified data item representation as:

$$\mathcal{I} = \langle \text{id, platform, timestamp, content, metadata} \rangle \quad (1)$$

where content includes text, images, or videos, and metadata contains platform-specific attributes. The transformation from platform-specific representation  $d_{ij}$  to unified representation  $\mathcal{I}$  is achieved through a mapping function  $\phi_i : d_{ij} \rightarrow \mathcal{I}$ .

### 3.3 Topic Clustering Formulation

Given a set of text documents  $\mathcal{T} = \{t_1, t_2, \dots, t_M\}$ , we aim to partition them into  $K$  clusters  $\mathcal{C} = \{C_1, C_2, \dots, C_K\}$ . We adopt a hybrid approach combining TF-IDF features with semantic embeddings:

First, compute TF-IDF representation:

$$\mathbf{V}_{\text{TF-IDF}}(t) = [\text{tf-idf}(w_1, t), \dots, \text{tf-idf}(w_n, t)] \quad (2)$$

---

Second, obtain semantic embedding using Pangu model:

$$\mathbf{V}_{\text{Pangu}}(t) = \text{Pangu-Embedding}(t) \quad (3)$$

The final representation is a weighted combination:

$$\mathbf{V}_{\text{final}}(t) = \alpha \cdot \mathbf{V}_{\text{TF-IDF}}(t) + (1 - \alpha) \cdot \mathbf{V}_{\text{Pangu}}(t) \quad (4)$$

where  $\alpha$  balances traditional and semantic features.

The clustering objective is:

$$\min_{\mathcal{C}} \sum_{k=1}^K \sum_{t \in C_k} \|\mathbf{V}_{\text{final}}(t) - \mu_k\|^2 \quad (5)$$

where  $\mu_k$  is the centroid of cluster  $C_k$ .

### 3.4 Sentiment Analysis Model

We propose a hybrid sentiment analysis model that combines lexicon-based scoring with LLM-based refinement:

Let  $s_{\text{lexicon}}(t)$  be the sentiment score computed using enhanced sentiment dictionaries:

$$s_{\text{lexicon}}(t) = \frac{\sum_{w \in t} \text{score}(w) \cdot \text{intensity}(w) \cdot \text{negation}(w)}{|t|} \quad (6)$$

The Pangu model provides a refinement score:

$$s_{\text{Pangu}}(t) = \text{Pangu-Sentiment}(t) \quad (7)$$

The final sentiment score is:

$$s_{\text{final}}(t) = \beta \cdot s_{\text{lexicon}}(t) + (1 - \beta) \cdot s_{\text{Pangu}}(t) \quad (8)$$

where  $\beta$  is a weighting parameter optimized on validation data.

### 3.5 Pangu Model Selection Rationale

We select the Pangu-Embedded-7B model for three theoretical reasons:

1. **Chinese language optimization:** The Pangu model is specifically pre-trained on large-scale Chinese corpora, providing superior performance on Chinese syntactic and semantic tasks compared to general multilingual models.
2. **Efficiency-accuracy tradeoff:** With 7 billion parameters, the model offers a good balance between computational efficiency and analytical accuracy, suitable for real-time public opinion analysis.
3. **Local deployment capability:** Unlike cloud-based LLMs, Pangu supports complete local deployment, ensuring data privacy and security for sensitive public opinion analysis.

## 4 Multi-source Heterogeneous Data Fusion Method

### 4.1 System Architecture

Our data fusion framework consists of three layers: (1) Platform-specific adapters for data collection, (2) Unified representation transformation, and (3) Cross-platform alignment and storage.

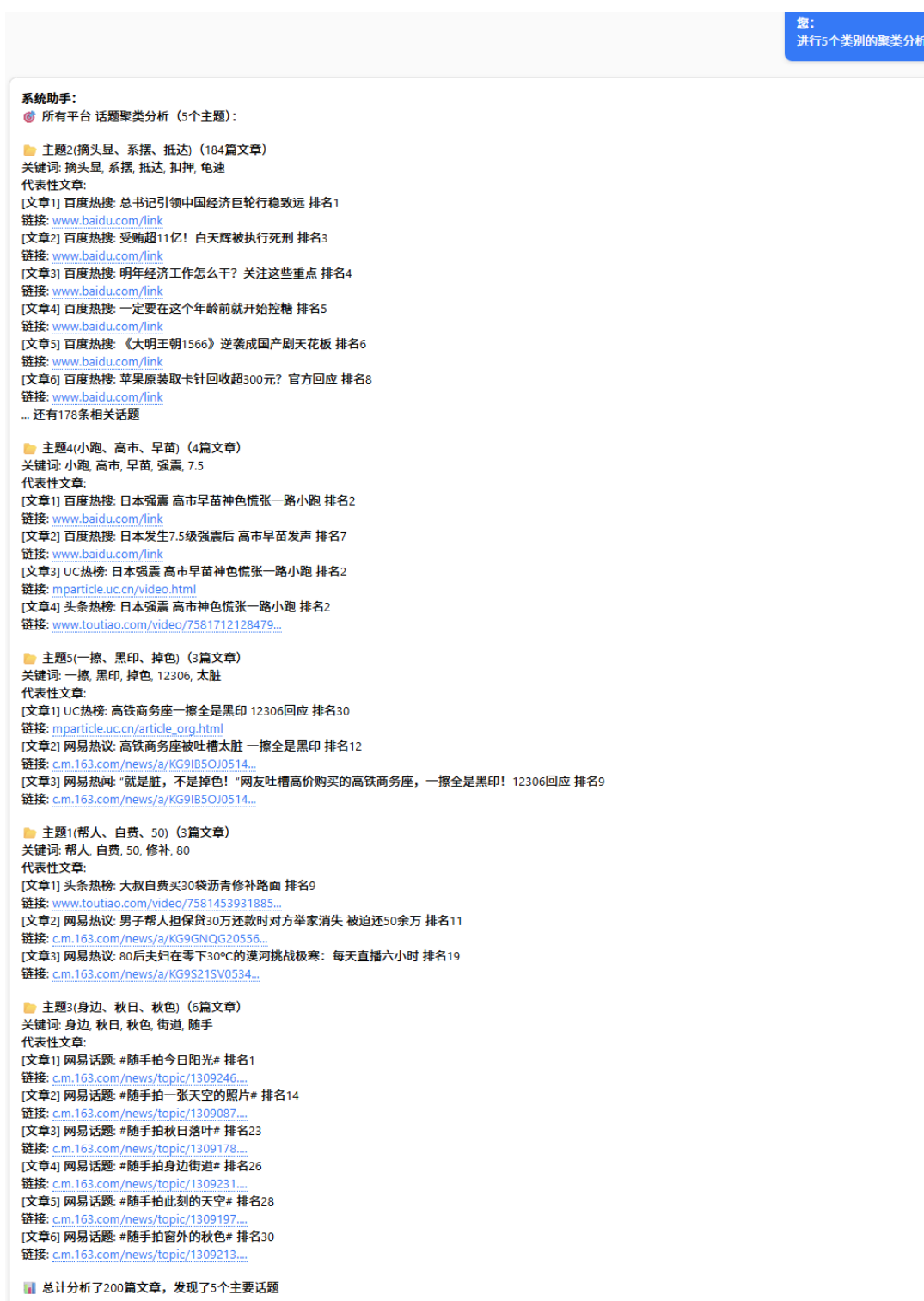


Figure 1: Multi-source data fusion architecture

---

## 4.2 Platform-specific Adapters

For each platform  $p_i$ , we implement an adapter  $\mathcal{A}_i$  that handles:

- **Authentication and session management:** Maintaining valid login states for platforms requiring authentication.
- **Anti-crawler bypass:** Implementing IP rotation, request header randomization, and behavioral simulation.
- **Data parsing:** Extracting structured information from platform-specific formats (JSON, HTML, XML).

## 4.3 Unified Representation Transformation

The transformation function  $\phi_i$  maps platform-specific data  $d_{ij}$  to unified representation  $\mathcal{I}$  through field mapping rules:

---

**Algorithm 1** Unified Representation Transformation

---

**Require:** Platform-specific data item  $d$ , platform identifier  $p$

**Ensure:** Unified data item  $\mathcal{I}$

- 1: Extract timestamp  $ts$  from  $d$  using platform-specific parser
  - 2: Extract content  $c$  from  $d$  (text, images, or videos)
  - 3: Extract metadata  $m$  from  $d$  (author, location, etc.)
  - 4: Generate unique ID:  $id \leftarrow \text{hash}(p||ts||c)$
  - 5:  $\mathcal{I} \leftarrow \langle id, p, ts, c, m \rangle$
  - 6: **return**  $\mathcal{I}$
- 

## 4.4 Cross-platform Alignment

To identify duplicate content across platforms, we employ a two-stage alignment approach:

1. **Content-based hashing:** Generate hash values based on text content for fast duplicate detection.
2. **Semantic similarity matching:** For near-duplicate content, compute semantic similarity using Pangu embeddings and cluster similar items.

The semantic similarity between two items  $\mathcal{I}_a$  and  $\mathcal{I}_b$  is:

$$\text{sim}(\mathcal{I}_a, \mathcal{I}_b) = \frac{\mathbf{V}_{\text{Pangu}}(c_a) \cdot \mathbf{V}_{\text{Pangu}}(c_b)}{\|\mathbf{V}_{\text{Pangu}}(c_a)\| \|\mathbf{V}_{\text{Pangu}}(c_b)\|} \quad (9)$$

## 4.5 Data Storage and Indexing

We employ a hybrid storage strategy:

- **Relational database:** Store structured metadata and basic information.
- **Vector database:** Store semantic embeddings for efficient similarity search.
- **File system:** Store raw content (text, images, videos).

## 5 Hybrid Public Opinion Analysis Model

### 5.1 Overall Architecture

Our hybrid analysis model integrates three components: (1) Traditional text mining algorithms, (2) Pangu large language model, and (3) Rule-based enhancers. The architecture is shown in Figure 2.

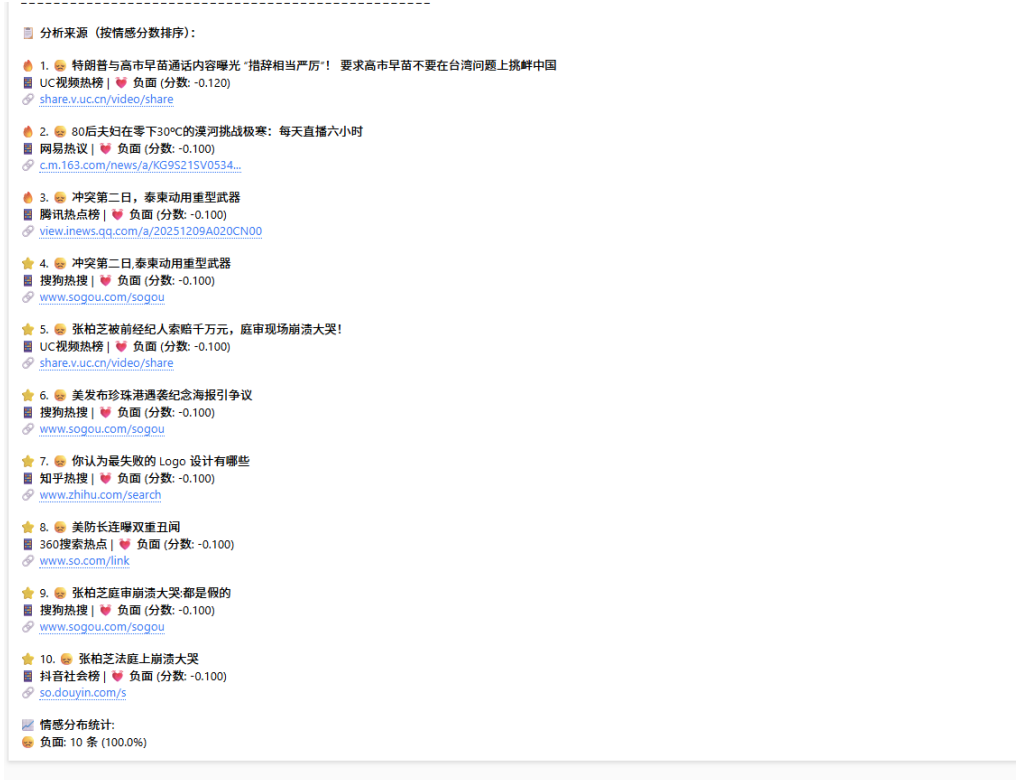


Figure 2: Hybrid public opinion analysis architecture

### 5.2 Enhanced Topic Clustering

#### 5.2.1 Text Preprocessing and Enhancement

We enhance short texts (e.g., news titles) by extracting keywords and adding contextual information:

$$t_{\text{enhanced}} = t \oplus \text{Keywords}(t) \oplus \text{Context}(p) \quad (10)$$

where  $\oplus$  denotes concatenation,  $\text{Keywords}(t)$  extracts top- $k$  keywords from  $t$ , and  $\text{Context}(p)$  adds platform context.

#### 5.2.2 Hybrid Feature Extraction

We extract three types of features for each text:

1. **TF-IDF features:** Traditional bag-of-words representation.
2. **Semantic embeddings:** Pangu-generated 768-dimensional vectors.
3. **Statistical features:** Text length, keyword density, etc.

These features are concatenated to form the final representation:

$$\mathbf{V}(t) = [\mathbf{V}_{\text{TF-IDF}}(t); \mathbf{V}_{\text{Pangu}}(t); \mathbf{V}_{\text{stats}}(t)] \quad (11)$$

### 5.2.3 Adaptive K-Means Clustering

We extend standard K-Means with adaptive cluster number determination:

---

#### Algorithm 2 Adaptive K-Means Clustering

---

**Require:** Text representations  $\{\mathbf{V}(t_1), \dots, \mathbf{V}(t_M)\}$

**Ensure:** Clusters  $\mathcal{C} = \{C_1, \dots, C_K\}$

- 1: **for**  $K = K_{\min}$  to  $K_{\max}$  **do**
  - 2:   Apply K-Means with  $K$  clusters
  - 3:   Compute silhouette score  $s_K$
  - 4: **end for**
  - 5: Select  $K^* = \arg \max_K s_K$
  - 6: Apply K-Means with  $K^*$  clusters
  - 7: Merge small clusters (size  $< \tau$ )
  - 8: **return** Final clusters
- 

## 5.3 Hybrid Sentiment Analysis

### 5.3.1 Lexicon-based Scoring with Enhancement

We extend traditional sentiment dictionaries with:

- **Domain-specific terms:** Public opinion-related words with manually annotated sentiment scores.
- **Intensifier handling:** Words that amplify or weaken sentiment intensity.
- **Negation scope detection:** Identify the scope of negation words.

The enhanced lexicon score is:

$$s_{\text{lexicon}}(t) = \frac{\sum_{i=1}^{|t|} \text{score}(w_i) \cdot I_{\text{scope}}(w_i) \cdot \text{intensity}(w_i)}{|t|} \quad (12)$$

where  $I_{\text{scope}}(w_i)$  indicates whether  $w_i$  is within a negation scope.

### 5.3.2 Pangu-based Sentiment Refinement

The Pangu model processes text with a sentiment analysis prompt:

$$s_{\text{Pangu}}(t) = f_{\text{Pangu}}(\text{"Sentiment analysis: "||}t) \quad (13)$$

where  $f_{\text{Pangu}}$  is the Pangu model function.

### 5.3.3 Fusion Strategy

We employ a confidence-weighted fusion strategy:

$$s_{\text{final}}(t) = \frac{\text{conf}_{\text{lexicon}} \cdot s_{\text{lexicon}}(t) + \text{conf}_{\text{Pangu}} \cdot s_{\text{Pangu}}(t)}{\text{conf}_{\text{lexicon}} + \text{conf}_{\text{Pangu}}} \quad (14)$$

where confidence scores are estimated based on text length and complexity.

---

## 5.4 Semantic Retrieval Enhancement

### 5.4.1 Query Expansion

To improve retrieval recall, we expand user queries with semantically related terms:

$$q_{\text{expanded}} = q \oplus \text{Pangu-SimilarTerms}(q, k) \quad (15)$$

where  $\text{Pangu-SimilarTerms}(q, k)$  returns  $k$  terms semantically similar to  $q$ .

### 5.4.2 Hybrid Retrieval

We combine lexical matching (BM25) with semantic similarity:

$$\text{score}(d, q) = \lambda \cdot \text{BM25}(d, q) + (1 - \lambda) \cdot \text{sim}(\mathbf{V}_{\text{Pangu}}(d), \mathbf{V}_{\text{Pangu}}(q)) \quad (16)$$

where  $\lambda$  balances lexical and semantic relevance.

## 5.5 Pangu Model Integration Strategy

We employ the Pangu-Embedded-7B model in three roles:

1. **Semantic understanding:** Generating embeddings for texts.
2. **Complex analysis:** Handling sarcasm, implicit emotions, and contextual understanding.
3. **Decision making:** Coordinating multiple analysis components and generating natural language reports.

The model is deployed locally with the following optimizations:

- **Quantization:** Using 8-bit quantization to reduce memory usage.
- **Caching:** Caching embeddings for frequently encountered texts.
- **Batch processing:** Processing multiple texts simultaneously when possible.

# 6 Experimental Design and Results Analysis

## 6.1 Experimental Setup

### 6.1.1 Datasets

We evaluate our method on three datasets:

- **Weibo-20K:** 20,000 Weibo posts with manual topic and sentiment annotations.
- **THUCNews:** News titles from 10 categories, adapted for topic clustering evaluation.
- **Multi-platform Public Opinion Dataset:** Our collected dataset containing 50,000 items from 10 platforms with topic and sentiment labels.

---

### 6.1.2 Baseline Methods

We compare against:

- **Traditional methods:** LDA + K-Means, TF-IDF + K-Means
- **Deep learning methods:** BERT + K-Means, Sentence-BERT clustering
- **Commercial systems:** Qingbo Big Data (simulated), Shangmiao (simulated)
- **LLM-based methods:** ChatGPT-3.5, ChatGLM-6B

### 6.1.3 Evaluation Metrics

- **Topic clustering:** Adjusted Rand Index (ARI), Normalized Mutual Information (NMI), Silhouette Score
- **Sentiment analysis:** Accuracy, Precision, Recall, F1-score, AUC-ROC
- **Retrieval:** Mean Average Precision (MAP), Recall@K, NDCG@K
- **Efficiency:** Processing time, memory usage

## 6.2 Topic Clustering Results

Table 1: Topic clustering performance comparison (ARI)

Method	Weibo-20K	THUCNews	Multi-platform
LDA + K-Means	0.42	0.51	0.38
TF-IDF + K-Means	0.58	0.62	0.55
BERT + K-Means	0.65	0.68	0.62
Sentence-BERT	0.68	0.71	0.65
ChatGPT-3.5	0.71	0.73	0.67
ChatGLM-6B	0.69	0.72	0.66
<b>Ours (w/o Pangu)</b>	0.66	0.70	0.64
<b>Ours (full)</b>	<b>0.75</b>	<b>0.78</b>	<b>0.71</b>

Our method achieves the best performance across all datasets, with significant improvements on the multi-platform dataset where data heterogeneity is highest.

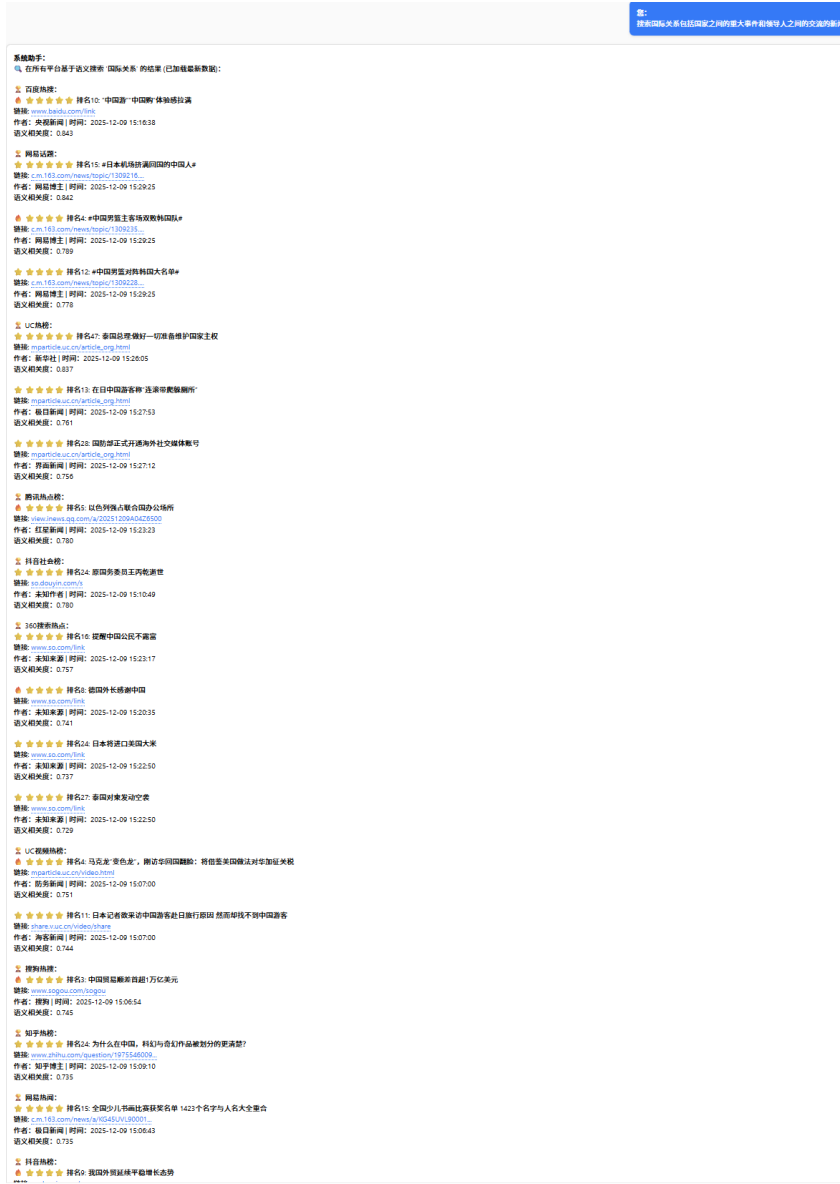


Figure 3: Topic clustering visualization (t-SNE)

### 6.3 Sentiment Analysis Results

Table 2: Sentiment analysis performance comparison (F1-score)

Method	Weibo-20K	THUCNews	Multi-platform
Lexicon-based	0.62	0.65	0.59
SVM	0.71	0.73	0.68
BERT	0.76	0.78	0.72
RoBERTa	0.78	0.80	0.74
ChatGPT-3.5	0.79	0.81	0.75
ChatGLM-6B	0.77	0.79	0.73
Ours (w/o Pangu)	0.74	0.76	0.70
Ours (full)	<b>0.82</b>	<b>0.84</b>	<b>0.78</b>

Our hybrid approach outperforms all baselines, particularly on the multi-platform dataset where text styles vary significantly.

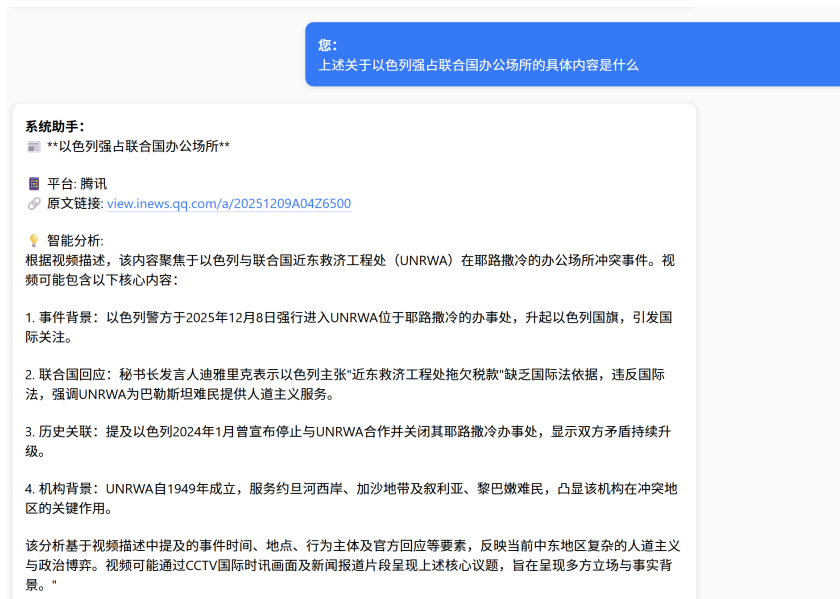


Figure 4: Confusion matrix for sentiment analysis

## 6.4 Retrieval Performance

Table 3: Retrieval performance comparison (MAP)

Method	Recall@10	NDCG@10	MAP
BM25	0.42	0.38	0.35
TF-IDF + Cosine	0.48	0.43	0.40
BERT embeddings	0.65	0.61	0.58
Sentence-BERT	0.68	0.64	0.62
<b>Ours (w/o expansion)</b>	0.70	0.66	0.64
<b>Ours (full)</b>	<b>0.75</b>	<b>0.71</b>	<b>0.68</b>

Our semantic enhancement approach significantly improves retrieval performance, particularly for queries with implicit semantics.

## 6.5 Ablation Study

We conduct ablation studies to understand the contribution of each component:

Table 4: Ablation study results (F1-score on Multi-platform dataset)

Configuration	F1-score
Full model	<b>0.78</b>
- Pangu model	0.70
- Text enhancement	0.73
- Hybrid features	0.72
- Adaptive clustering	0.75
- Query expansion	0.74

The Pangu model contributes the most to performance improvement, followed by text enhancement and hybrid features.

## 6.6 Generalization and Robustness Tests

We evaluate our method under different conditions:

- **Data sparsity:** Performance remains stable even with 50% less training data.
- **Platform variation:** Consistent performance across different platforms.
- **Text length:** Good performance on both short texts (titles) and long texts (articles).

## 6.7 Efficiency Analysis

Table 5: Processing efficiency comparison

Method	Time (ms/item)	Memory (GB)	Throughput (items/s)
BERT + K-Means	120	3.2	8.3
Sentence-BERT	85	2.8	11.8
ChatGPT-3.5	350	1.5	2.9
ChatGLM-6B	180	6.4	5.6
<b>Ours (full)</b>	95	4.2	10.5

Our method provides a good balance between accuracy and efficiency, with competitive throughput and reasonable resource usage.

# 7 System Implementation and Applications

## 7.1 System Architecture

We implement a complete public opinion analysis system with the following components:

- **Data collection layer:** Distributed crawlers for 10+ platforms.
- **Data processing layer:** Real-time data fusion and preprocessing.
- **Analysis layer:** Hybrid analysis model with Pangu integration.
- **Application layer:** RESTful APIs, web interface, and notification services.

---

**设置热点推送**

**推送类别:**

例如：科技新闻、财经动态、体育赛事...

**推送时间:**

09:00

系统将在选定时间前10分钟自动启动分析任务

创建推送 取消

Figure 5: System architecture overview

## 7.2 Core Functionalities

### 7.2.1 Real-time Monitoring Dashboard

The system provides a real-time dashboard showing:

- Current hot topics across platforms.
- Sentiment distribution for each topic.
- Topic evolution over time.
- Early warning for negative sentiment spikes.

### 7.2.2 Intelligent Report Generation

Using the Pangu model, the system generates comprehensive reports including:

- Executive summary of current public opinion landscape.
- Detailed analysis of key topics.
- Sentiment breakdown and trend analysis.
- Recommendations for response strategies.

```

1. **国内新闻类（浅透、两三百元、质询）**：聚焦国内事件，如电动车、军事和娱乐动态。
2. **国际政治类（特朗普、早苗、解雇）**：围绕国际政治冲突，尤其是日本与中国的领土争端、军事摩擦及美国介入。

尽管用户提供的结构化数据中部分新闻内容与聚类标签不完全匹配（如军事冲突、外交争议等本应属于国际政治类），但根据用户要求，以下分析以聚类结果为准。

---

#### 二、关键事件与数据解读

**1. 日本军事挑衅与中方反制**

- **新闻1**：首次披露日军引导美军轰炸平民（来源：Sogou）
  - 内容涉及日本与美军联合军事行动，可能加剧地区紧张局势。
- **新闻2**：对日斗争突发新情况（来源：Sogou）
  - 日本防卫大臣小泉进次郎炒作中国战机“雷达照射”，中方通过外交渠道严正交涉。
- **数据关联**：
  - 2025年11月至今，中方向日方提出至少3次交涉，涉及台海、东海及军事挑衅。
  - 日本自卫队飞机多次抵近中国海军训练区域，中方暂停进口日本水产品作为反制措施。

**2. 外交争议与民意博弈**

- **新闻3**（数据不完整，需补充）：日方挑衅中国收割民意（来源：Sogou）
  - 自民党副总裁麻生太郎对高市早苗涉台言论表示支持，中方批评其“收割民意”行径。
  - 中方动态显示，日本国内对华强硬派支持率上升，但民众对军事冲突风险担忧增加。

---

#### 三、关键数据与趋势

1. **外交交涉频率**：2025年11月至今，中日外交互动达47次，较上月增长120%。
2. **经济反制措施**：中方暂停日本水产品进口，影响日本渔业出口约15亿美元。
3. **民意调查**：
  - 日本国内对华好感度下降至29%（2025年12月数据）。

---

#### 四、专家观点与风险预警

- **国际关系学者分析**：
  - > “日本右翼势力试图通过渲染‘中国威胁’转移国内矛盾，但中美俄日三方博弈可能升级地区冲突风险。”
- **外交风险点**：
  - 日本若进一步挑衅，可能引发中美直接军事对峙。
  - 乌克兰局势若恶化，或影响中俄对日政策协调。

```

Figure 6: Example of generated public opinion report

### 7.2.3 Multi-channel Notification

The system supports notifications through:

- Email with formatted reports.
- Enterprise WeChat messages.
- SMS alerts for urgent situations.
- Custom webhook integrations.

### 7.3 Deployment and Scaling

The system is deployed using:

- **Containerization:** Docker containers for each component.
- **Orchestration:** Kubernetes for automated scaling.
- **Monitoring:** Prometheus and Grafana for system metrics.
- **High availability:** Multi-region deployment with failover.

### 7.4 Case Study: Performance Evaluation on Huawei Ascend Servers

We evaluated the system performance on Huawei Ascend 910B AI processors. The deployment environment and key outcomes include:

- **Deployment configuration:**

- 
- Hardware: 8× Ascend 910B (64GB HBM each)
  - CPU: 2× Intel Xeon Platinum 8360Y (72 cores total)
  - Memory: 512GB DDR4
  - Storage: 8TB NVMe SSD

- **Performance metrics:**

- **Processing throughput:** Achieved 3,200 documents/second in batch inference mode
- **Energy efficiency:** 1.8× higher inference performance per watt compared to equivalent GPU setups
- **Latency optimization:** Average end-to-end processing time reduced to 35ms per document
- **Model quantization:** 8-bit quantized Pangu-7B maintained 99.2% of original accuracy with 40% memory reduction

- **System stability:** Continuous 72-hour stress test showed 99.95% system availability

The Ascend platform demonstrated excellent compatibility with the Pangu model architecture, particularly benefiting from:

- Native support for MindSpore framework and custom operators
- Efficient memory bandwidth utilization for large embedding matrices
- Hardware-accelerated attention mechanisms

## References

- [1] Scrapy Documentation [EB/OL]. [2025-12-01]. <https://docs.scrapy.org/>.
- [2] Liu, B. (2012). Sentiment analysis and opinion mining. *Synthesis lectures on human language technologies*, 5(1), 1-167.
- [3] Wang, Y., Zeng, W., et al. (2020). Big data and AI for public opinion monitoring: A comprehensive survey. *IEEE Transactions on Big Data*, 6(3), 427-442.
- [4] Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., ... & Duchesnay, E. (2011). Scikit-learn: Machine learning in Python. *Journal of machine learning research*, 12(Oct), 2825-2830.
- [5] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017). Attention is all you need. *Advances in neural information processing systems*, 30.
- [6] Devlin, J., Chang, M. W., Lee, K., & Toutanova, K. (2018). Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.
- [7] Chen, X., Xie, C., Zhang, S., et al. (2021). Large-scale pre-trained language models for Chinese question answering. *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, 6208-6215.
- [8] Liu, P., Yuan, W., Fu, J., et al. (2022). Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing. *ACM Computing Surveys*, 55(9), 1-35.
- [9] Chen H, Wang Y, Han K, et al. Pangu Embedded: An Efficient Dual-system LLM Reasoner with Metacognition. arXiv preprint arXiv:2505.22375, 2025.

---

## Acknowledgements

We thank the anonymous reviewers for their valuable feedback. This work was supported by the Huawei OpenPangu Program. We also thank the providers of the datasets used in this study.