

Transit Pulse: Utilizing Social Media as a Source for Customer Feedback and Information Extraction with Large Language Model

Jiahao Wang^a, Amer Shalaby^a

^aTransit Analytics Lab (TAL), Department of Civil and Mineral Engineering (Transportation Engineering, University of Toronto, 35 St. George Street, Toronto, Ontario M5S 1A4 Canada

ARTICLE INFO

Keywords:

Public Transit System
Large Language Model
Customer Service
Social Media
Information Extraction
Semantic Analysis

ABSTRACT

Users of the transit system flood social networks daily with messages that contain valuable insights crucial for improving service quality. These posts help transit agencies quickly identify emerging issues. Parsing topics and sentiments is key to gaining comprehensive insights to foster service excellence. However, the volume of messages makes manual analysis impractical, and standard NLP techniques like Term Frequency-Inverse Document Frequency (TF-IDF) fall short in nuanced interpretation. Traditional sentiment analysis separates topics and sentiments before integrating them, often missing the interaction between them. This incremental approach complicates classification and reduces analytical productivity. To address these challenges, we propose a novel approach to extracting and analyzing transit-related information, including sentiment and sarcasm detection, identification of unusual system problems, and location data from social media. Our method employs Large Language Models (LLM), specifically Llama 3, for a streamlined analysis free from pre-established topic labels. To enhance the model's domain-specific knowledge, we utilize Retrieval-Augmented Generation (RAG), integrating external knowledge sources into the information extraction pipeline. We validated our method through extensive experiments comparing its performance with traditional NLP approaches on user tweet data from the real world transit system. Our results demonstrate the potential of LLMs to transform social media data analysis in the public transit domain, providing actionable insights and enhancing transit agencies' responsiveness by extracting a broader range of information.


1. Introduction

Understanding user feedback is crucial for public transit agencies. Collecting and analyzing this feedback improves service quality, improves operational efficiency, builds community trust, and supports informed, data-driven decisions. Through user feedback, agencies gain insights into passengers' perceptions of the transit service, identifying satisfactory and unsatisfactory aspects, pinpointing when and where issues occur, and understanding overall needs. This process enables agencies to adapt to changing demands, address safety concerns, and maintain transparency and accountability. Ultimately, this leads to increased ridership and higher customer satisfaction.

However, collecting user feedback is a challenging process. Traditional methods, such as physical or online surveys, while effective, have several limitations. These methods have a limited reach and can be costly, requiring significant investments in money, time, and labor. In addition, respondents often experience survey fatigue, which can reduce the quality of feedback. Feedback collection through surveys is typically delayed, often conducted annually or subannually, which is useful for long-term analysis and decision making but impractical for understanding user perceptions in real-time [21].

Periodic surveys tend to capture opinions about the overall system or a broad network, making it difficult for agencies to obtain specific information about when, where, and what issues occurred. This lack of timely, detailed feedback hampers the ability of transit agencies to address immediate concerns and make prompt improvements to their services.

In contrast, collecting feedback through social media posts offers a more dynamic and cost-effective and relatively more real-time alternative [23]. Social media platforms allow agencies to reach a broader audience and collect crowd-sourced information from diverse users. This method is not only cost-effective, but also provides near-real-time feedback, which is crucial for timely responses and adjustments. Furthermore, social media offers access to a massive dataset, which can be seen as the foundation of applying complex analysis tools [31].

 jhope.wang@mail.utoronto.ca (J. Wang); amer.shalaby@utoronto.ca (A. Shalaby)

ORCID(s):

However, using social media data for user feedback analysis presents several challenges [15]. From a data point of view, social networks often contain irrelevant or low-quality information. The linguistic structure is complex, with frequent use of sarcasm, slang, and informal language, making accurate interpretation difficult. Furthermore, social media data is less structured compared to traditional survey data, which lacks preset questions and predefined problem categories. Consequently, social media data related to public transit can be less specific, posing challenges to extract actionable information and increasing data labeling costs. Moreover, the high volume of data generated on social media platforms requires effective data management strategies to handle and analyze it efficiently.

From a methodological tool perspective, traditional Natural Language Processing (NLP) analysis tools for customer satisfaction often rely on fixed lexicons, limiting their ability to accurately interpret diverse language use on social media. As addressed in [18], social media posts, like tweets, are typically brief and lack context, making it difficult for traditional NLP methods to understand the full scope of an issue. Misinterpretations are common, especially when dealing with sarcasm, irony, or domain-specific information, without additional context.

In this paper, our objective is to improve the use of social media data for public transit user feedback by developing an advanced information extraction tool utilizing a Large Language Model (LLM). This tool offers improved sentiment analysis and expands the scope of feedback by moving beyond predefined topic categories, allowing for a more comprehensive understanding of user experiences. In addition, the tool is designed to extract geographical data from tweets, providing valuable location-specific insights. By increasing the ability to efficiently identify and address customer needs, this tool aims to enable more timely and effective service improvements, ultimately enhancing public transit services and boosting customer satisfaction.

The rest of this paper is organized as follows. The next section provides a general introduction to LLMs, discussing their capabilities and advancements compared with traditional NLP models. The following section explores state-of-the-art applications in the domain of social media analysis. Subsequently, the methodology and framework for using LLMs in the analysis of transit posts on social media are presented. The following section discusses the experimental results of using LLM for social media analysis. The last section addresses limitations, explores potential solutions, and presents possible future works.

2. Literature Review

This section provides a brief overview of widely used methods for information extraction from social media.

2.1. Manual Identification

Manual identification is a common method used to extract useful information from social media. For instance, authors in [26] and [24] manually identified geo-information from tweets mentioning station or route names. Although this method guarantees accuracy, it is inefficient for large volumes of tweets. In [6], 1,624 tweets were identified for route information extraction, while [26] involved manually labeling 3,454 tweets for spatial information. While human labeling is valuable for building benchmark datasets or lexicons, it is highly costly and impractical for frequent data analysis on fast-updating social media datasets.

2.2. Lexicon-Based Approaches

One of the most prominent methods for information extraction is the lexicon-based approach. This method identifies information, such as sentiment, using a collection of tokens with predefined scores. The overall sentiment of a sentence is determined based on the scores of these tokens [17]. Lexicon-based approaches are common in sentiment analysis, where tokens are assigned scores of -1, 0, or 1, representing negative, neutral, or positive sentiments, respectively. The sentence score is then calculated accordingly [14, 7]. To enhance the performance of lexicon-based methods, it is often necessary to generate specific lexicons tailored to particular classification tasks. For instance, in [16], a lexicon-based method was used to identify nine problematic topics in a tweet dataset related to the Calgary transit system. Similarly, [26] employed 435 predefined terms to classify tweets into four topics: punctuality, comfort, breakdowns, and overcrowding. In another study, [3] introduced a human-involved lexicon-building process to identify tweet topics, such as bus-related issues, sentiment, and sarcasm.

Despite its popularity, the lexicon-based approach has inherent drawbacks for social media data analysis. Firstly, dividing sentences into tokens can lead to the loss of contextual information. For example, irony or sarcasm can render individual tokens' meanings opposite to their intended sentiment in the sentence. Additionally, the informal writing habits prevalent on social media make it challenging to build comprehensive lexicon datasets that accommodate informal word usage.

2.3. Machine Learning-Based Approaches

Another notable method for information extraction is the machine learning (ML)-based approach. In this method, ML or deep learning (DL) models, such as logistic regression (LR), support vector machines (SVM), and decision trees (DT) [38], are trained on well-structured datasets for specific topic classification tasks. These models are commonly used in sentiment analysis. Additionally, transformer models with large structures, such as BERT, have been applied to topic classification tasks for social media analysis [26]. A well-constructed dataset with relevant labels and similar content is crucial for ML-based methods. Transfer learning can alleviate this limitation by using models trained on datasets built for similar tasks in different domains, then applying them to current tasks [25]. For example, in [19], the author used BERT to build a model for extracting transit topics from social media. Similarly, in [9], the BERT model was used for classifying pedestrian maneuver types. The model was trained on an online customer feedback dataset with 11 different topics and then used for further classification tasks related to the Washington Metropolitan Area Transit Authority.

However, even with transfer learning, ML-based approaches are still constrained by the scope of the training dataset and cannot identify topics not present in the dataset. Moreover, ML-based classification tasks are typically single-aspect. For instance, a model trained for sentiment analysis only performs sentiment analysis. This approach becomes costly when multiple aspects of information need to be extracted from the target text.

2.4. Large Language Model-Based Approaches

Large Language Models (LLMs) are a type of deep learning model built using multi-layer Transformer architectures [36], containing vast numbers of parameters and typically pre-trained on large-scale corpora [?]. Through pre-training, LLMs acquire general knowledge, common sense, and the ability to understand and generate text. Due to their exceptional capabilities, LLMs are increasingly adopted for information extraction tasks. Beyond traditional natural language processing (NLP) tasks such as sentiment analysis [43] and semantic analysis [41], LLMs are widely applied in downstream tasks like knowledge reasoning, question answering, relation extraction, and event extraction, often outperforming traditional NLP models [39].

In the transportation domain, LLMs have been used to answer transportation-related questions such as those concerning transportation economics and driver characteristics [32]. Additionally, LLMs are applied in areas such as transportation infrastructure planning and design [30], project management [1], operations and maintenance [29], and safety control [13]. Multimodal LLMs are also used for tasks like object detection in transportation [4].

As noted in [27, 35], high-quality data ecosystems are essential for effective LLM applications in specific domains. For instance, efforts to benchmark LLM performance in transportation have resulted in datasets like TransportBench [32], designed to assess reasoning abilities in transportation problems, and various question-answer datasets [42], which evaluate decision-making, complex event causality reasoning, and human driving exam performance. Moreover, [28] provides an image dataset for evaluating LLMs' ability to detect transportation-related issues like cracks or congestion. The GTFS-related dataset in [12] evaluates LLM performance in semantic understanding and information retrieval.

However, as discussed in [13], while LLMs are powerful, they face challenges when applied to domain-specific tasks due to knowledge gaps. Fine-tuning is a potential solution, as demonstrated in [37], where an open-source multimodal LLM, VisualGLM, was fine-tuned with 12.5 million textual tokens to enhance its transportation domain knowledge. Similarly, in [44], LLaMA 3 was fine-tuned with a traffic safety dataset to improve its performance in that domain.

Fine-tuning pre-trained LLMs requires large amounts of high-quality training data and significant computational resources. A lightweight alternative is the Retrieval-Augmented Generation (RAG) system, which retrieves domain-specific information from external databases, as proposed in [40]. However, this framework remains largely conceptual, with limited real-world applications.

To address the gap in applying LLMs to transportation information extraction, particularly in user feedback analysis, we introduce Transit Pulse. This is one of the first attempts to provide a lightweight solution for semi-automatic information extraction from social media posts, offering insights for transit monitoring and control.

3. Methodology

This paper addresses two user feedback analysis tasks: traditional classification and information extraction & summarizing. For the traditional multi-class classification task, we cover sentiment classification, sarcasm detection, and transit problem topic classification. This section introduces the classification and information extraction pipelines used in our experiments.

3.1. Traditional Classification Task

As shown in Fig. 1, the traditional classification process involves data processing, vectorization, model training or fine-tuning, and model evaluation. We use both traditional NLP methods—Term Frequency-Inverse Document Frequency (TF-IDF) with machine learning (ML)-based classification—and large language model (LLM)-based classification methods.

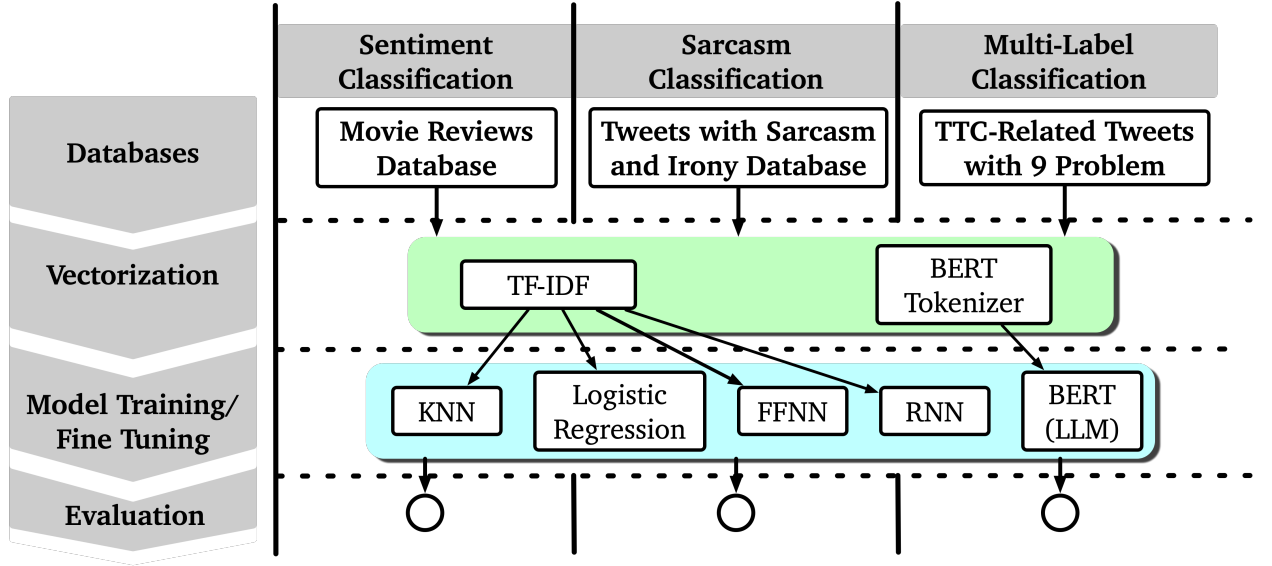


Figure 1: Framework for the Traditional Classification Task, illustrating the steps from data processing to model evaluation.

The initial step involves vectorizing or tokenizing the neutral language input using either TF-IDF or the BERT tokenizer. This converts the input into a multi-dimensional space vector for further training.

3.1.1. TF-IDF

TF-IDF is a statistical measure that evaluates the importance of a word in a document relative to a collection of documents (corpus). It is widely used in information retrieval and text mining.

- **Term Frequency (TF)** measures how frequently a term occurs in a document, normalized to prevent bias towards longer documents:

$$tf(t, d) = \frac{\text{Number of times term } t \text{ appears in document } d}{\text{Total number of terms in document } d}$$

- **Inverse Document Frequency (IDF)** measures the importance of a term across the corpus, reducing the weight of frequently occurring terms:

$$idf(t, D) = \log \left(\frac{N}{|\{d \in D : t \in d\}|} \right)$$

where N is the total number of documents in the corpus D , and $|\{d \in D : t \in d\}|$ is the number of documents in which the term t appears.

- **TF-IDF score** for a term t in a document d is the product of its term frequency and inverse document frequency:

$$tf-idf(t, d, D) = tf(t, d) \times idf(t, D)$$

TF-IDF provides a numerical statistic reflecting the importance of a word to a document in a corpus, enabling effective text analysis and retrieval.

3.1.2. BERT Tokenizer and BERT Model

In contrast to TF-IDF, the BERT tokenizer, part of a pre-trained LLM (the embedding layer in BERT), tokenizes and embeds text into dense vectors with contextual meaning. The BERT tokenizer splits text into tokens using WordPiece, then processes these tokens through bidirectional transformers pre-trained on large corpora to capture contextual information.

The tokenized input is further processed by the pre-trained BERT model [10], which uses self-attention to learn contextual information. The BERT model, thanks to its bidirectional structure, understands word meanings from both directions in a sentence and from distant words. Fine-tuning the BERT model for classification involves adding a dense layer that takes the model's output as a high-level feature vector for final classification. Pre-trained on massive datasets, the BERT model's general understanding of language aids the classification layer in making accurate classifications.

3.2. Information Extraction

Traditional classification methods are easy to deploy and generally reliable but face limitations. They are restricted by the quality and scope of the training data set and can only classify predefined topics. Models trained for specific tasks, such as sentiment classification, cannot be directly used for other tasks like sarcasm detection. Additionally, performance can degrade when applied to scenarios different from the training dataset.

To overcome these limitations, we introduce an information extraction pipeline based on the powerful open-source LLM, Llama 3 by Meta [2]. Llama 3 excels in over 150 benchmark tasks, including answering science / domain-specific questions and common sense reasoning [5]. Its power comes from high-quality training data and a large model structure. Llama 3 was pre-trained with more than 15 trillion tokens, covering all high-quality open data available until December 2023. The model's 70 billion parameters enable it to learn from this massive dataset through supervised fine-tuning (SFT) and reinforcement learning with human feedback (RLHF), resulting in accurate and human-preferred outputs.

As illustrated in Fig. 2, the information extraction pipeline begins by embedding target tweets into a structured prompt to guide the LLM in extracting the desired information. The prompt defines the LLM's role and tasks, including extracting and summarizing information related to the transit agency tweeted about (in our case the Toronto Transit Commission, TTC for short), such as station name, sentiment, sarcasm, and problem topic. The output is JSON-like text, which is processed through an information aggregation step. This step involves segregating text chunks into key-value pairs, filtering unuseful text with regex, and constructing a structured JSON dataset. The consensus mechanism determines the most common answer from multiple LLM responses to mitigate performance variations.

3.3. Retrieval Augmented Generation (RAG) System

Despite Llama 3's strengths in context understanding and NLP tasks, it has limitations with unfamiliar data or domain-specific knowledge. For instance, it may misinterpret station names in TTC-related tweets, especially when names are abbreviated or misspelled, as shown in Fig. 3.

To address this, we implemented a Retrieval Augmented Generation (RAG) system [20]. The RAG system supplements the LLM with external knowledge to improve accuracy in domain-specific information extraction.

As shown in Fig. 4, the RAG process starts by embedding the external knowledge base using a pre-trained LLM embedding model in a vector space. Each document is represented as a point in this high-dimensional space. When a query (e.g., tweet content) is embedded into the same space, we compare it to the knowledge base to find the closest matches. These matches are relevant external documents selected based on distance metrics like Euclidean Distance, cosine similarity, or maximum inner product (MIP). The retrieval process typically outputs more candidates than needed for re-ranking by a more complex embedding system or LLM-driven prompts. The retrieved information is then added to the information extraction pipeline, enhancing accuracy and comprehensiveness.

In summary, our methodology combines advanced LLM techniques and the RAG system, to effectively classify and extract information from social media data related to public transit.

4. Experiments

This section presents our experiments on using LLM-based pipeline for traditional classification problems and information extraction & summarization. We begin by introducing the experimental setup, including datasets, models, and hardware configuration. Then, we discuss model performance on each dataset and the information extracted by our pipeline. Additionally, we provide a case study showcasing the practical application of our information extraction system for real-time public transit service monitoring.

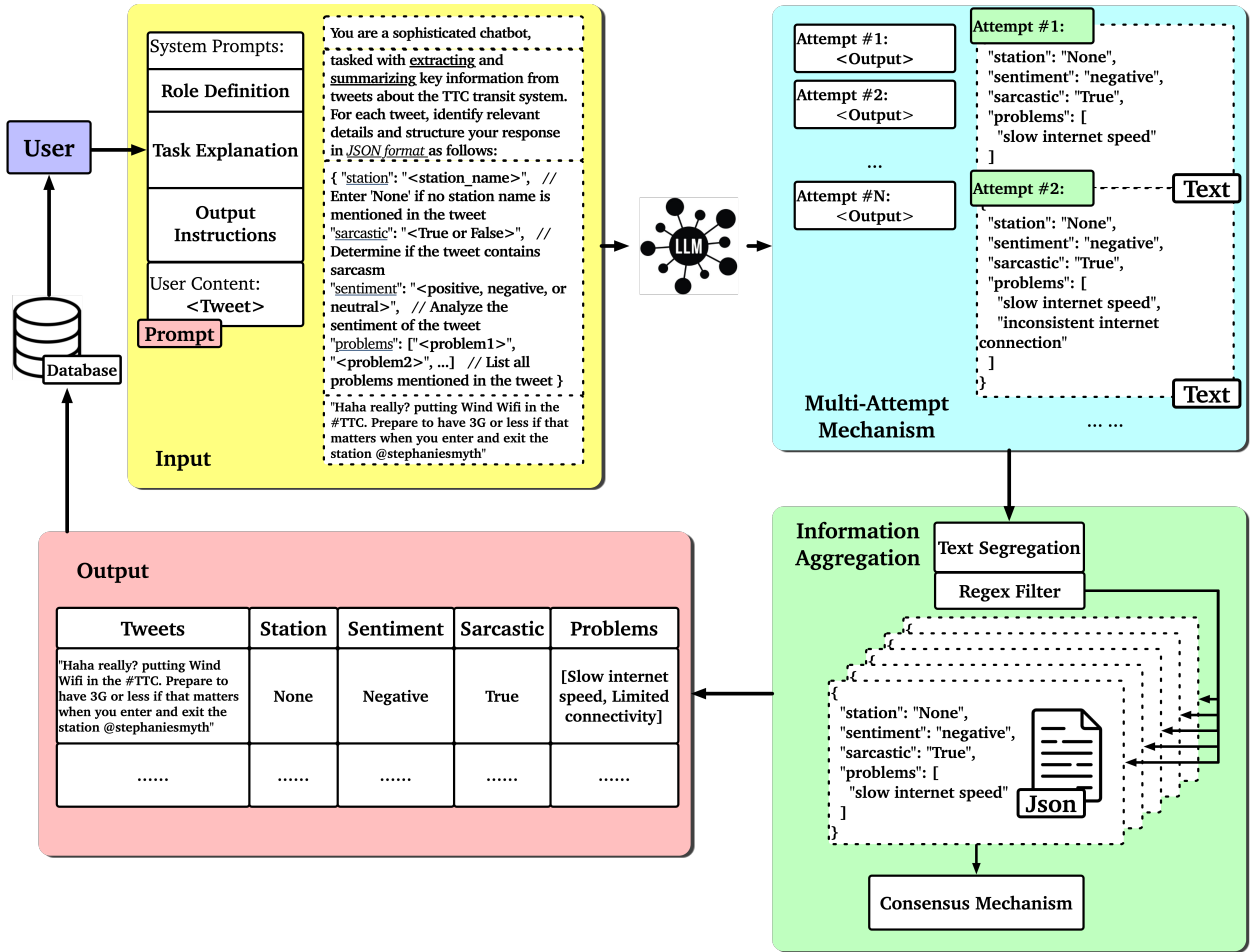


Figure 2: Information Extraction Pipeline with Llama 3, demonstrating the process of embedding tweets into structured prompts and aggregating the extracted information.

- "Haha really? putting Wind Wifi in the #TTC. Prepare to have 3G or less if that matters when you enter and exit the station @stephaniesmyth"

(a) Extracted "TTC" as station name.
- You need better comms, #TTC. Line 1 is in CHAOS and all you're saying is a delay Lawr. to Shep. Try harder. Do better. And fix the signals!

(b) Extracted "Lawr and Shep" as station name.
- I just can't believe! I wait years for the DAMN streetcar to come at Baaaaathurst station

(c) Extracted "Baaaaathurst Station" as station name.

Figure 3: Examples of Station Information Misinterpretation in Tweets.

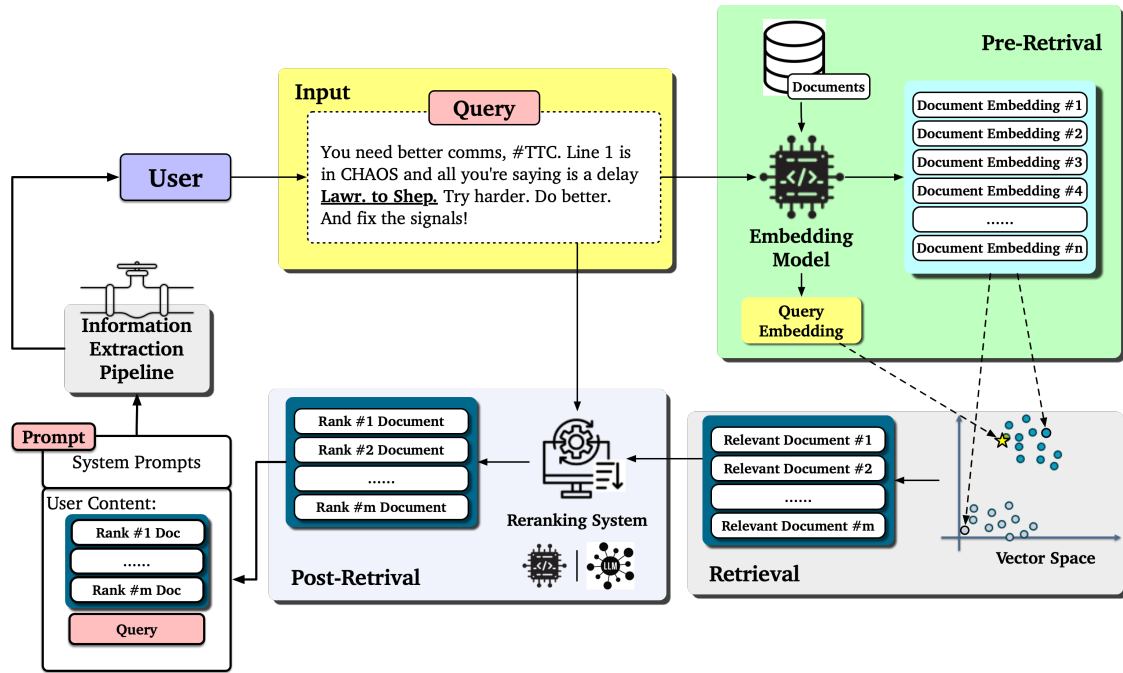


Figure 4: Retrieval Augmented Generation (RAG) System, showing the process of embedding external knowledge and matching it with user queries to enhance information extraction accuracy.

4.1. Datasets

We used three datasets to test model performance on multi-class classification tasks for sentiment analysis, sarcasm detection, and transit problem classification.

For sentiment analysis, we used the Sentiment Analysis on Movie Reviews dataset [8], containing 156,060 training records and 66,292 testing records across five sentiment groups: negative, somewhat negative, neutral, somewhat positive, and positive.

For sarcasm detection, we used the Tweets with Sarcasm and Irony dataset [22], which includes four classes: irony, sarcasm, regular, and figurative. The dataset has 54,618 training records and 7,861 testing records.

These datasets are publicly available and specifically designed for sentiment analysis and sarcasm detection. We did not use our TTC tweets dataset for performance testing due to the lack of human-labeled sentiment and sarcasm information.

For the transit-related problem classification task, we used tweets related to the Toronto Transit Commission (TTC), which also serves as our information extraction dataset. The data was collected from February 5th, 2015, to December 31st, 2015, from two main TTC X (formerly Twitter) accounts (as shown in Fig. 6): TTC Service Alerts and TTC Customer Service. The TTC operates 192 bus routes, 11 streetcar routes, and 3 subway lines, serving over 1.4 million riders daily at its peak in 2023 [33].

The TTC-related tweets dataset includes 631,691 records. After removing duplicates and retweets with minimal additional information, 27,312 tweets remained for analysis. Figure 7 shows the number of tweets posted at different times of the day, peaking during peak hours.

The dataset categorizes tweets into 10 problem categories: maintenance, capacity availability, interaction with staff, travel time, ride quality, winter maintenance, temporal availability, safety and security, accessibility, and communication (Table 1). Initially, a group of keywords for each problem topic was manually identified from a subset of the entire dataset. Subsequently, a lexicon-based method was employed to label the remaining tweets based on the frequency of keyword occurrences, assigning each tweet to the most relevant problem category.

To evaluate model performance with and without the RAG system, we utilized a GTFS-specific question dataset [11] designed to assess the understanding of GTFS standards and the retrieval of information from structured GTFS data. This dataset comprises 195 questions across six categories: term definitions, common reasoning, file structure, attribute

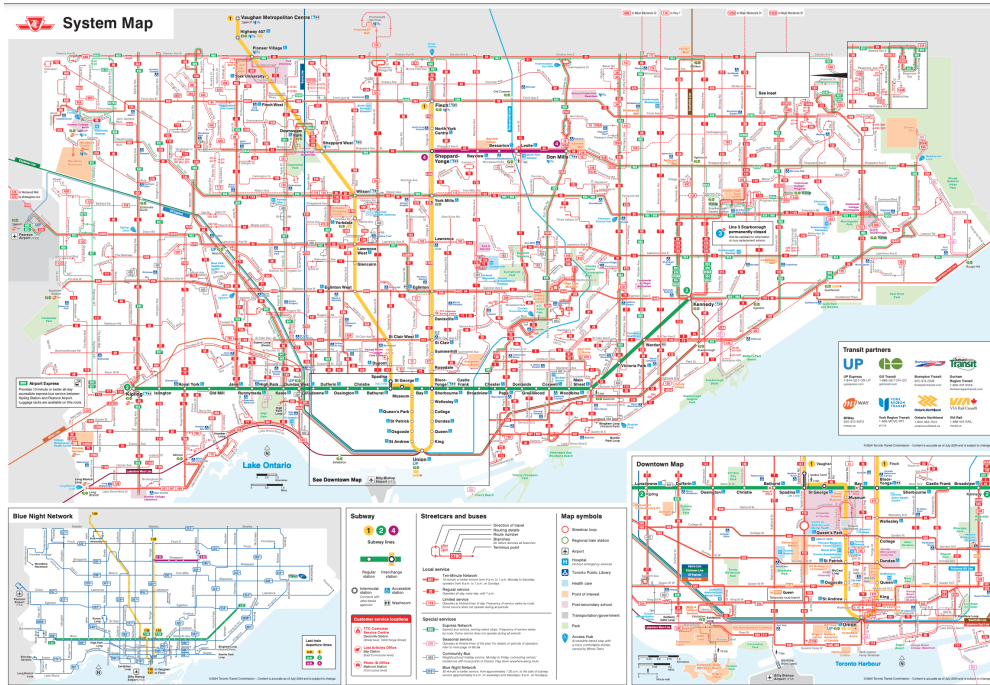
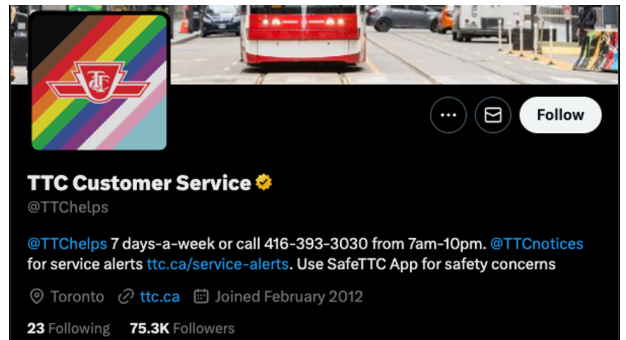


Figure 5: System map of the Toronto Transit Commission (TTC) [34].



(a) TTC Service Alerts, 495.1 thousand followers when taking this screenshot



(b) TTC Customer Service, 75.3 thousand followers when taking this screenshot

Figure 6: Two Main Official X (formerly Twitter) Accounts of TTC (Screenshot taken on June 26, 2024).

mapping, data structure, and categorical mapping. An example of these questions is shown in Fig. 8.

According to [11], Term definitions questions test the model’s ability to understand specific GTFS terms and document structures. Common reasoning questions evaluate basic GTFS knowledge, including abbreviations, usage, and file purposes. File structure questions determine the model’s ability to identify the correct files in given contexts. Attribute mapping questions assess whether the model can correctly associate attributes with their respective files. Data structure questions verify the model’s capability to identify attribute data types accurately. Categorical mapping questions examine the model’s understanding of categorical variables representing data, such as different codes for wheelchair availability on buses.

In addition to these, the dataset includes 87 programming questions aimed at testing the model’s ability to retrieve information from the GTFS dataset effectively. An example of these questions is shown in Fig. 9.

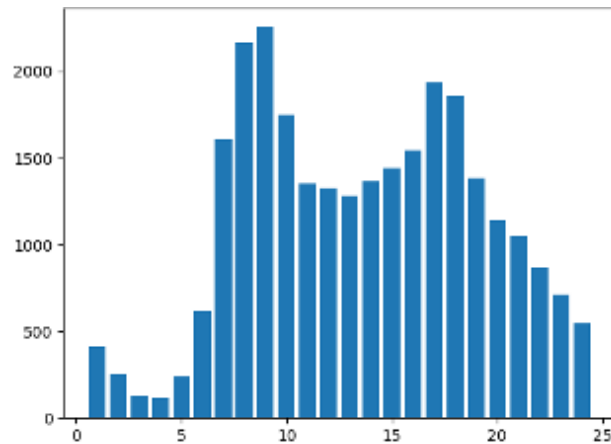


Figure 7: Number of tweets published at different times of the day.

Category	Count
Winter Maintenance	24
Temporal Availability	179
Interaction with Staff	245
Maintenance	396
Capacity Availability	553
Communication	648
Accessibility	1263
Ride Quality	1599
Travel Time	2882
Safety and Security	2943

Table 1

Summary of various transit system aspects and their counts.

Question: What value is used in the "wheelchair_accessible" field of the "trips.txt" file to indicate that the trip has no specific information regarding wheelchair accessibility?

a) 0 b) 1 c) None of these d) 3

Answer: a

Figure 8: Example of GTFS understanding benchmarking questionnaire.

Question: What is the route_short name for route_id 192?

Answer: 192

Figure 9: Example of GTFS programming question.

Model	Sentiment Analysis	Sarcasm Detection	Problem Topic Classification
Logistic Regression (TF-IDF)	62.90%	74.32%	81.59%
KNN	63.19%	60.80%	84.40%
FNN	64.87%	75.46%	86.83%
RNN	64.71%	75.39%	79.12%
BERT (LLM)	78.02%	82.64%	91.50%

Table 2
Model Performance Comparison

4.2. Models and Environment Setup

We implemented various models for classification tasks, including logistic regression, KNN, FFNN, and RNN, to compare with the relatively small LLM: BERT. As discussed in the methodology section, traditional classification methods using ML/DL are based on TF-IDF for text embedding and trained from scratch. In contrast, BERT is fine-tuned for each task using corresponding datasets, showcasing the power of LLM structures.

For the information extraction pipeline, we used LLama-3 as the core for text understanding, information extraction, and summarization.

The experiments were conducted in the following environment:

- **System:** Pop OS
- **Processor:** 2.4 GHz 8-Core Intel Core i9
- **GPU:** Nvidia 4090
- **Programming Language:** Python
- **Large Language Model (for classification):** BERT
- **Large Language Model (for information extraction):** LLama-3 8B
- **Prompt Framework:** LangChain

4.3. Experiment Results

4.3.1. Classification Tasks

Table 2 presents the accuracy results of various models across three different tasks: Sentiment Analysis, Sarcasm Detection, and Problem Topic Classification.

The LLM model demonstrates the highest accuracy across all three tasks, substantially outperforming other models, especially in Sentiment Analysis and Sarcasm Detection. For Problem Topic Classification, BERT (LLM) achieves 91.5% accuracy, highlighting its robustness in handling complex language understanding tasks compared to traditional and other neural network models.

4.3.2. GTFS Understanding and Retrieval

This section presents the performance of the LLM for GTFS understanding and retrieval tasks, with and without the RAG system. Figure 10 shows the prompts for the LLM answering multiple choice questions on GTFS understanding, both with and without the retrieved information from the RAG system. The external documents in the RAG system were built using the official GTFS documentation¹.

The results of these experiments are shown in Figure 11. For multiple choice questions across six types, the LLM with RAG consistently achieves a higher number of correct answers. Notably, for Categorical Mapping questions, the model with RAG scores 51 correct answers compared to 27 without RAG. Overall, the accuracy of answering GTFS-related multiple choice questions increases from 57.43% to 73.85%. For programming questions, accuracy improves from 24.14% (21/87) to 52.87% (46/87), demonstrating the enhanced capability of RAG in improving LLM performance for domain-specific questions.

¹<https://gtfs.org/schedule/reference/>

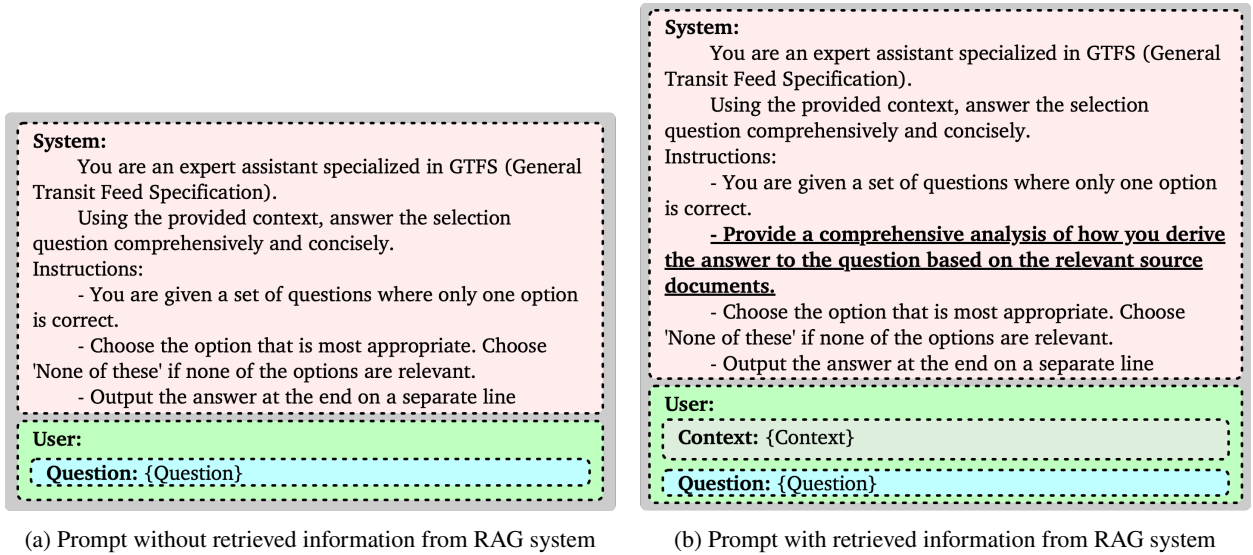


Figure 10: Prompts for LLM answering GTFS understanding multiple choice questions, with and without information from RAG system

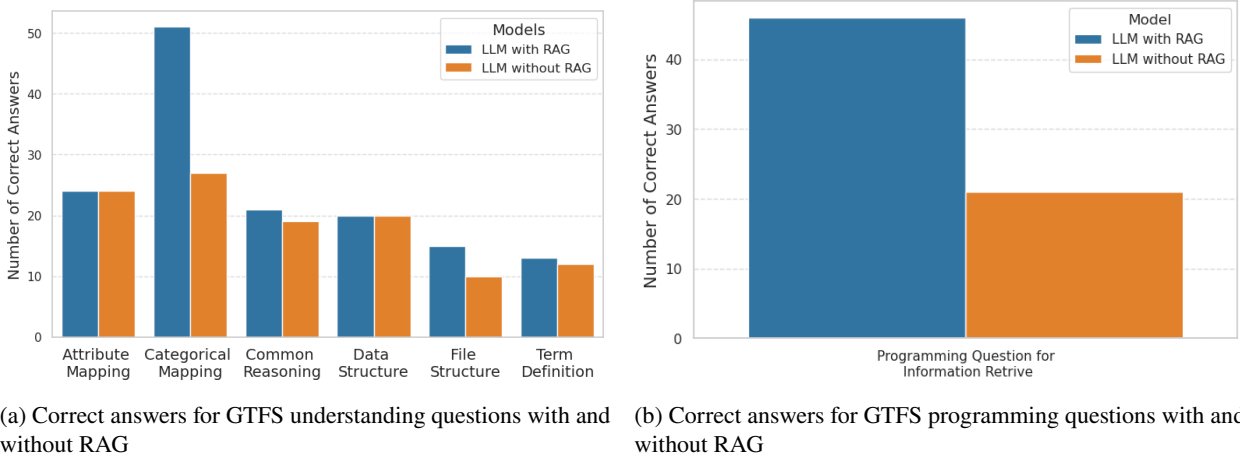


Figure 11: Performance of LLM on GTFS tasks with and without RAG system

4.3.3. Station Name Extraction with RAG

The prompts used in the RAG system during the post-retrieval phase to enhance the accuracy and formality of station name extraction are shown in Fig. 12. To guide the language model in selecting the most relevant station name from the retrieved list, we employed the Chain-of-Thoughts (CoT) and Few-Shot techniques. These methods instruct the model to think step-by-step and provide reference examples for making decisions.

Figure 13 shows three examples of using the RAG system for station extraction. As illustrated, the RAG system enables the language model to extract more formal station information, avoiding issues related to using abbreviations or misspellings. Additionally, the model can better handle instances where "TTC" is incorrectly identified as a station name.

4.3.4. Information Extraction with LLM

This section demonstrates information extraction from the TTC-related tweets dataset using LLM, compared with traditional NLP methods. A notable advantage of the LLM is its ability to generate comprehensive outputs, including station information, sentiment, sarcasm, and problem summarization in a single response.

Extract station name from a tweet.
Think it step by step.
 1. First, you should read the tweet and the list of stations provided.
 2. You should determine whether a station name in the list is/are specifically mentioned in the tweet.
 2.0 If there is no station name ****specificly**** mentioned in the tweet, you should only select 'None' from the list of stations provided.
 2.1 If there is, you should only select the station name from the list of stations provided.
 3. Your answer should be in the format of a list e.g. ['station1', 'station2', ...].
Do not infer or guess the station name if it is not clearly mentioned in the tweet.

Example #1:

 Tweet: "#RiderAlert Viva Blue experiencing various delays, SB, of up to 20 minutes due to traffic delays and broken down TTC bus slowing traffic."
 Station List:
 =====
 Bloor St West at Bathurst St
 Bloor St West at Green Lanes East Side
 Bloor St West at Bay St West Side
 =====
 Your selection: ['None']

Example #2: [example_content]
 Example #3: [example_content] ...

User:
 Here is the Station List:
 =====
 {stations}
 =====
 Tweet: "{tweet}"
 Your selection:

Figure 12: Prompts used in the RAG system to improve station name extraction.

Tweets	Station Name Extracted by LLM <u>Without Domain-Specific Knowledge</u>	Station Name Extracted by LLM <u>With RAG</u>
"Haha really? putting Wind Wifi in the #TTC. Prepare to have 3G or less if that matters when you enter and exit the station @stephaniesmyth"	{... "Station": "TTC" }	{... "Station": "None" }
You need better comms, #TTC. Line 1 is in CHAOS and all you're saying is a delay <u>Lawr. to Shep</u> . Try harder. Do better. And fix the signals!	{... "Station": "Lawr and Shep" }	{... "Station": ["Lawrence", "Sheppard"], }
I just can't believe! I wait years for the DAMN streetcar to come at <u>Baaaaathurst station</u>	{... "Station": "Baaaaathurst" }	{... "Station": "Bathurst" }

Figure 13: Comparison of station name extraction by Llama 3 with and without the RAG system.

Figure 14 compares sentiment and sarcasm labels extracted by traditional NLP methods and LLM. The left heatmap, using traditional NLP, shows an implausible balance between sentiment and sarcasm labels. Conversely, the LLM-labeled heatmap reveals that most sarcastic tweets are labeled with negative sentiment, demonstrating the LLM’s advanced capability in detecting and interpreting sarcasm, as well as sentiment identification. Figure 15 illustrates examples where traditional NLP methods identified the posts as positive. However, our proposed methods correctly identified them as negative and sarcastic.

Figure 16 shows examples of the most challenging cases where the language model identified posts as sarcastic but still positive. These examples demonstrate that even for humans, it is difficult to discern the writer’s true opinion

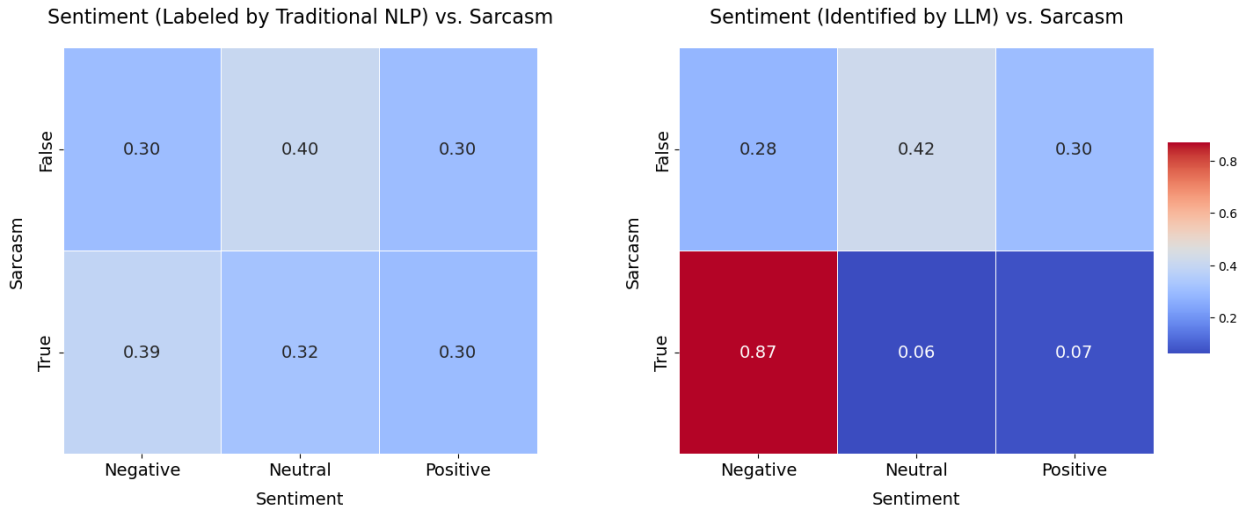


Figure 14: Comparison of Sentiment and Sarcasm Labels: Traditional NLP (left) vs. LLM (right)

without additional context.

Thank you TTC. I just love when the train decides to go out of service halfway to its destination
 Gotta love taking the streetcar down Queen and Jameson with people who just punch people randomly and brawl at 9am #ttc
 When I asked driver when next 504 EB to Broadview, driver of 3rd Roncey car just shrugged. Great customer service, TTC! #NotTheBetterWay
 Witnessed a #TTC driver allowing a guy on for free. I guess that no asking for fare thing has started already eh? @TTChelps #workingclass

Figure 15: Examples where traditional NLP methods identified posts as positive, but our proposed methods identified them as negative and sarcastic.

I swear TTC is holding casting calls for their streetcar drivers, why you all so good-looking?! My driver said he'd see me later
 Not all ttc bus drivers are rude woah
 I know the tide is turning when I'm actually enjoying a TTC bus ride.
 I made my bus. Thank u ttc gods

Figure 16: Examples of challenging sentiment analysis where the LLM identified posts as sarcastic but positive.

For system problem extraction, the LLM summarizes the issues mentioned in tweets. Figure 17 shows the most frequently mentioned keywords in the LLM-generated problem summaries for each problem category. The alignment of these keywords with the problem categories illustrates the LLM’s ability to accurately understand and summarize context.

Despite the dataset being (semi-)manually labeled, it contains errors due to the manual labeling process. Figure 18 compares human-labeled and LLM-extracted problem information. The number of records with problems identified by humans is significantly lower than those identified by the LLM, highlighting the limitations of manual labeling.

As shown in Fig. 18a, negative tweets, whether estimated by traditional NLP or LLM, are almost evenly found in the 'No problem detected' and 'Problem detected' categories when those categories are labeled by humans, revealing inconsistencies. In contrast, Fig. 18b shows that using the LLM method for labelling the categories, negative tweets are more accurately identified, aligning better with common sense and behavioral logic.

4.3.5. Case Study: System Monitoring with LLM Powered Information Extraction Pipeline

This section demonstrates how the LLM-powered information extraction pipeline can help transit agencies respond promptly and efficiently to system issues.



Figure 17: Word Clouds for Different Problem Categories Summarized by LLM

Starting with station information, geo-specific data allows for station-level performance analysis. Figure 19 shows tweet counts for the five most-mentioned stations from 7 AM to 11 AM. Union, Yonge, and Eglinton stations are frequently mentioned, with Bloor station peaking from 9 AM to 10 AM, indicating a possible incident.

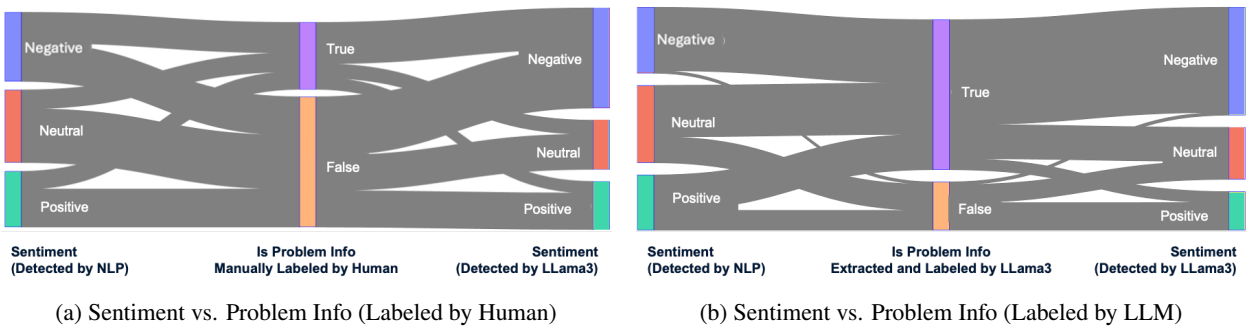


Figure 18: Comparison of Sentiment Labels with Problem Information Extraction

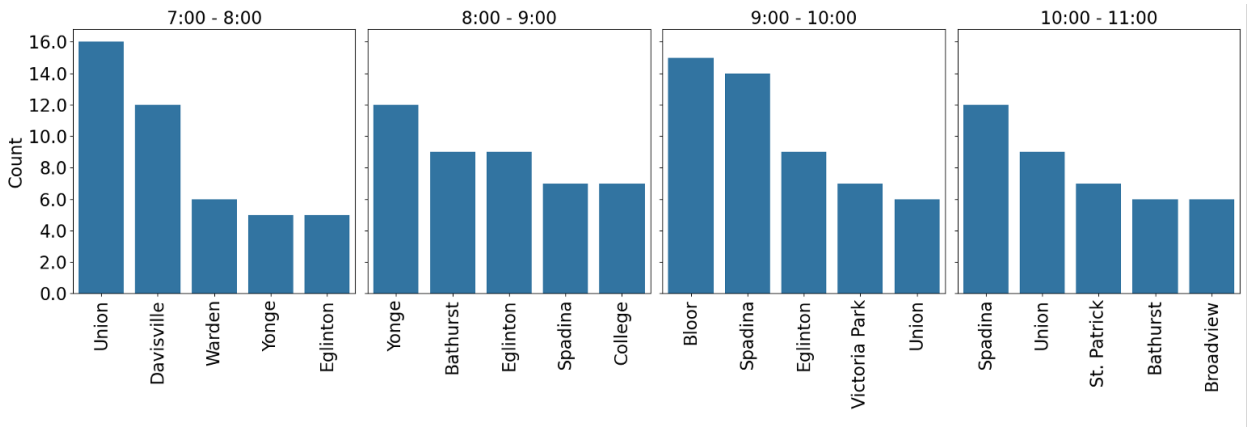


Figure 19: Tweet counts for the five most-mentioned stations from 7 AM to 11 AM

Examining negative tweets related to Bloor station between 9 AM and 10 AM (Figure 20), we identify significant issues. Figure 21 compares keywords extracted directly from tweets and LLM summaries. Traditional NLP methods fall short in clarity, whereas LLM summarization clearly identifies the main issue: unusually long lines for shuttle buses, indicating a capacity problem during the morning peak hours. With this information, transit authorities can further investigate using onsite cameras and implement better solutions to address the passenger overflow.

'Line up to get back into bloor Station #TTC #commuterlife4ever #somuchforECON',
 'Longer than normal wait times for trains on Line 1 (Yonge), between Union and Bloor, due to signal related problems at Bloor Station. #TTC',
 'Man this commute. Left at 8:15 am and I'm walking from bloor. Shuttle bus line way to long #ttc', 'The line up to get a SB #ttc train at',
 'This-&#gt;@kjkfid Why is the Bloor subway line closure & delays big news for TTC? This is how the King Streetcar line operates on a daily basis.',
 'WOW TTC YOU F**KED UP. THE F**KING BLOOR LINE UP TO GET ON A SHUTTLE BUS UNTIL BAY STREET. THIS IS A JOKE',
 'Why is the Bloor subway line closure & delays big news for TTC? This is how the King Streetcar line operates on a daily basis.'

Figure 20: Negative tweets related to Bloor station from 9 to 10 AM

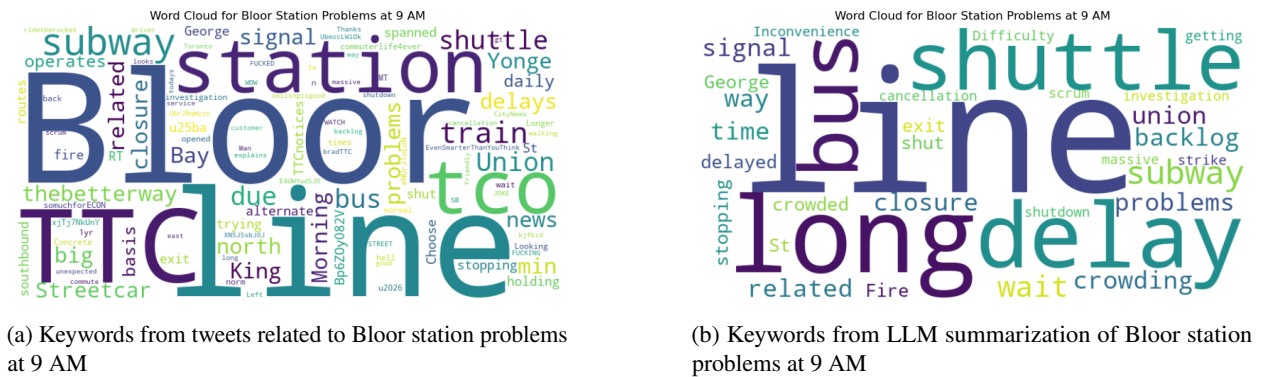


Figure 21: Comparison of keywords related to Bloor station problems at 9 AM extracted directly from tweets and from LLM summarization

This example shows how our advanced techniques provide actionable insights for real-time problem-solving.

5. Conclusion and Future Work

The LLM-powered information extraction pipeline has demonstrated significant advantages in analyzing social media data for transit systems. This pipeline enhances sentiment analysis by simultaneously detecting sentiment and sarcasm, providing a more accurate understanding of public opinions.

Additionally, the pipeline facilitates actionable insights for transit agencies by transitioning from system-level to station-level analysis. By extracting station-specific information and summarizing potential issues mentioned in social media posts, it enables targeted interventions and improvements.

Moreover, using this LLM-based information extraction pipeline reduces the dependency on pre-identified labels in datasets. By adjusting the prompts and incorporating relevant external guidance documents, the pipeline can extract more useful information with less human effort, broadening the scope of analysis.

However, there are limitations and areas for future research. The TTC-related tweets dataset used in this study lacks comprehensive human annotations for sentiment classification and sarcasm detection. Additionally, station information was not pre-extracted, and there was no human review of labeling results. Building a more robust dataset with high-quality annotations is crucial. This would provide ground truth for classification and extraction tasks and benefit future LLM studies in the transit domain by offering a reliable dataset for training and fine-tuning models, serving as a benchmark for LLM applications in public transit.

Furthermore, refining the prompts used in the LLM is necessary to ensure a more consistent output format, which would enhance the information aggregation process and minimize the loss of relevant information.

Lastly, the LLM used for information extraction, Llama 3, is a high-performance model with 8 billion parameters, requiring substantial computing resources. Future work will focus on employing techniques such as Knowledge

Distillation to transfer knowledge from larger models to smaller ones, thereby improving performance and efficiency for handling a relatively limited range of tasks.

References

- [1] Abbas, M., Ferrari, A., Shatnawi, A., Enoiu, E., Saadatmand, M., Sundmark, D., 2023. On the relationship between similar requirements and similar software: A case study in the railway domain. *Requirements Engineering* 28, 23–47.
- [2] AI@Meta, 2024. Llama 3 model card URL: https://github.com/meta-llama/llama3/blob/main/MODEL_CARD.md.
- [3] Al-Sahar, R., Klumpenhouwer, W., Shalaby, A., El-Diraby, T., 2024. Using twitter to gauge customer satisfaction response to a major transit service change in calgary, canada. *Transportation Research Record* 2678, 190–206.
- [4] Ashqar, H.I., Jaber, A., Alhadidi, T.I., Elhenawy, M., 2024. Advancing object detection in transportation with multimodal large language models (mllms): A comprehensive review and empirical testing. *arXiv preprint arXiv:2409.18286*.
- [5] Beeching, E., Fourrier, C., Habib, N., Han, S., Lambert, N., Rajani, N., Sanseviero, O., Tunstall, L., Wolf, T., 2023. Open llm leaderboard. https://huggingface.co/spaces/open-llm-leaderboard-old/open_llm_leaderboard.
- [6] Casas, I., Delmelle, E.C., 2017. Tweeting about public transit—gleaning public perceptions from a social media microblog. *Case Studies on Transport Policy* 5, 634–642.
- [7] Collins, C., Hasan, S., Ukkusuri, S.V., 2013. A novel transit rider satisfaction metric: Rider sentiments measured from online social media data. *Journal of Public Transportation* 16, 21–45.
- [8] Cukierski, W., 2014. Sentiment analysis on movie reviews. URL: <https://kaggle.com/competitions/sentiment-analysis-on-movie-reviews>.
- [9] Das, S., Oliiae, A.H., Le, M., Pratt, M.P., Wu, J., 2023. Classifying pedestrian maneuver types using the advanced language model. *Transportation research record* 2677, 599–611.
- [10] 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*.
- [11] Devunuri, S., Qiam, S., Lehe, L., 2023. Chatgpt for gtfs: From words to information. *arXiv preprint arXiv:2308.02618*.
- [12] Devunuri, S., Qiam, S., Lehe, L.J., 2024. Chatgpt for gtfs: benchmarking llms on gtfs semantics... and retrieval. *Public Transport*, 1–25.
- [13] Du, Y., 2024. Large models in transportation infrastructure: a perspective. *Intelligent Transportation Infrastructure* 3, liae007.
- [14] El-Diraby, T., Shalaby, A., Hosseini, M., 2019. Linking social, semantic and sentiment analyses to support modeling transit customers' satisfaction: Towards formal study of opinion dynamics. *Sustainable Cities and Society* 49, 101578.
- [15] Grant-Muller, S.M., Gal-Tzur, A., Minkov, E., Nocera, S., Kuflik, T., Shoor, I., 2015. Enhancing transport data collection through social media sources: methods, challenges and opportunities for textual data. *IET Intelligent Transport Systems* 9, 407–417.
- [16] Hosseini, M., El-Diraby, T., Shalaby, A., 2018. Supporting sustainable system adoption: Socio-semantic analysis of transit rider debates on social media. *Sustainable cities and society* 38, 123–136.
- [17] Kiritchenko, S., Zhu, X., Mohammad, S.M., 2014. Sentiment analysis of short informal texts. *Journal of Artificial Intelligence Research* 50, 723–762.
- [18] Kuflik, T., Minkov, E., Nocera, S., Grant-Muller, S., Gal-Tzur, A., Shoor, I., 2017. Automating a framework to extract and analyse transport related social media content: The potential and the challenges. *Transportation Research Part C: Emerging Technologies* 77, 275–291.
- [19] Leong, M., Abdelhalim, A., Ha, J., Patterson, D., Pincus, G.L., Harris, A.B., Eichler, M., Zhao, J., 2024. Metroberta: Leveraging traditional customer relationship management data to develop a transit-topic-aware language model. *Transportation Research Record*, 03611981231225655.
- [20] Lewis, P., Perez, E., Piktus, A., Petroni, F., Karpukhin, V., Goyal, N., Küttler, H., Lewis, M., Yih, W.t., Rocktäschel, T., et al., 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in Neural Information Processing Systems* 33, 9459–9474.
- [21] Li, C.L., 2016. Better, quicker, together: enabling public transport service quality co-monitoring through a smartphone-based platform. Ph.D. thesis. Massachusetts Institute of Technology.
- [22] Ling, J., Klinger, R., 2016. An empirical, quantitative analysis of the differences between sarcasm and irony, in: *The Semantic Web: ESWC 2016 Satellite Events, Heraklion, Crete, Greece, May 29–June 2, 2016, Revised Selected Papers* 13, Springer. pp. 203–216.
- [23] Liu, J.H.I., Ban, X.J., et al., 2017. Measuring the impacts of social media on advancing public transit. Technical Report. National Institute for Transportation and Communities.
- [24] Lucini, F.R., Tonetto, L.M., Fogliatto, F.S., Anzanello, M.J., 2020. Text mining approach to explore dimensions of airline customer satisfaction using online customer reviews. *Journal of Air Transport Management* 83, 101760.
- [25] Meng, J., Long, Y., Yu, Y., Zhao, D., Liu, S., 2019. Cross-domain text sentiment analysis based on cnn_ft method. *Information* 10, 162.
- [26] Osorio-Arjona, J., Horak, J., Svoboda, R., García-Ruíz, Y., 2021. Social media semantic perceptions on madrid metro system: Using twitter data to link complaints to space. *Sustainable Cities and Society* 64, 102530.
- [27] Papageorgiou, E., Chronis, C., Varlamis, I., Himeur, Y., 2024. A survey on the use of large language models (llms) in fake news. *Future Internet* 16, 298.
- [28] Prajapati, S., Singh, T., Hegde, C., Chakraborty, P., 2024. Evaluation and comparison of visual language models for transportation engineering problems. *arXiv preprint arXiv:2409.02278*.
- [29] Rane, N.L., 2023. Multidisciplinary collaboration: key players in successful implementation of chatgpt and similar generative artificial intelligence in manufacturing, finance, retail, transportation, and construction industry.
- [30] Roberts, J., Lüddecke, T., Das, S., Han, K., Albanie, S., 2023. Gpt4geo: How a language model sees the world's geography. *arXiv preprint arXiv:2306.00020*.
- [31] Schweitzer, L., 2012. How are we doing? opinion mining customer sentiment in us transit agencies and airlines via twitter. Technical Report.
- [32] Syed, U., Light, E., Guo, X., Zhang, H., Qin, L., Ouyang, Y., Hu, B., 2024. Benchmarking the capabilities of large language models in transportation system engineering: Accuracy, consistency, and reasoning behaviors. *arXiv preprint arXiv:2408.08302*.

- [33] (TTC), T.T.C., 2023. 2022 operating statistics. URL: <https://www.ttc.ca/transparency-and-accountability/Operating-Statistics/Operating-Statistics---2023>.
- [34] (TTC), T.T.C., 2024. Ttc system map. URL: https://cdn.ttc.ca/-/media/Project/TTC/DevProto/Images/Home/Routes-and-Schedules/Landing-page-pdfs/TTC_SystemMap_2024-03.pdf?rev=606dcf47e1cf4020856cbf68d34ec872.
- [35] Tupayachi, J., Xu, H., Omitaomu, O.A., Camur, M.C., Sharmin, A., Li, X., 2024. Towards next-generation urban decision support systems through ai-powered construction of scientific ontology using large language models—a case in optimizing intermodal freight transportation. *Smart Cities* 7, 2392–2421.
- [36] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, Ł., Polosukhin, I., 2017. Attention is all you need. *Advances in neural information processing systems* 30.
- [37] Wang, P., Wei, X., Hu, F., Han, W., 2024. Transgpt: Multi-modal generative pre-trained transformer for transportation. arXiv preprint arXiv:2402.07233 .
- [38] Wankhade, M., Rao, A.C.S., Kulkarni, C., 2022. A survey on sentiment analysis methods, applications, and challenges. *Artificial Intelligence Review* 55, 5731–5780.
- [39] Xu, D., Chen, W., Peng, W., Zhang, C., Xu, T., Zhao, X., Wu, X., Zheng, Y., Chen, E., 2023. Large language models for generative information extraction: A survey. arXiv preprint arXiv:2312.17617 .
- [40] Xu, H., Yuan, J., Zhou, A., Xu, G., Li, W., Ye, X., et al., 2024a. Genai-powered multi-agent paradigm for smart urban mobility: Opportunities and challenges for integrating large language models (llms) and retrieval-augmented generation (rag) with intelligent transportation systems. arXiv preprint arXiv:2409.00494 .
- [41] Xu, S., Wu, Z., Zhao, H., Shu, P., Liu, Z., Liao, W., Li, S., Sikora, A., Liu, T., Li, X., 2024b. Reasoning before comparison: Llm-enhanced semantic similarity metrics for domain specialized text analysis. arXiv preprint arXiv:2402.11398 .
- [42] Zhang, J., Ilievski, F., Ma, K., Kollaa, A., Francis, J., Oltramari, A., 2023a. A study of situational reasoning for traffic understanding, in: *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pp. 3262–3272.
- [43] Zhang, W., Deng, Y., Liu, B., Pan, S.J., Bing, L., 2023b. Sentiment analysis in the era of large language models: A reality check. arXiv preprint arXiv:2305.15005 .
- [44] Zheng, O., Abdel-Aty, M., Wang, D., Wang, C., Ding, S., 2023. Trafficsafetygpt: Tuning a pre-trained large language model to a domain-specific expert in transportation safety. arXiv preprint arXiv:2307.15311 .