# RAILWAY ACCIDENT ANALYTICS
# A DATA DRIVEN AI APPROACH

A report submitted for the course named Project II-(CS3202)

Submitted By

KASARABADA YASHWANTH SAI
SEMESTER - VI
ROLL NO : 220103012

Supervised By

DR. SANASAM CHANU INUNGANBI

DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING
INDIAN INSTITUTE OF INFORMATION TECHNOLOGY SENAPATI, MANIPUR

APRIL, 2025

# Declaration

In this submission, I have expressed my idea in my own words, and I have adequately cited and referenced any ideas or words that were taken from another source. I also declare that I adhere to all principles of academic honesty and integrity and that I have not misrepresented or falsified any ideas, data, facts, or sources in this submission. If any violation of the above is made, I understand that the institute may take disciplinary action. Such a violation may also engender disciplinary action from the sources which were not properly cited or permission not taken when needed.

Kasarabada Yashwanth Sai
Roll No : 220103012

DATE:

Department of Computer Science Engineering
Indian Institute of Information Technology Senapati, Manipur

Dr. Sanasam Chanu Inunganbi                    Email: inunganbi@iiitmanipur.ac.in
Assistant Professor                            Contact No: +91 7005844201

# *To Whom It May Concern*

This is to certify that the project/internship report entitled **"Railway Accident Analytics: A Data-Driven AI Approach",** submitted to the department of Computer Science and Engineering, Indian Institute of Information Technology Senapati, Manipur in partial fullfillment for the award of degree of Bachelor of Technology in Computer Science and Engineering is record bonafide work carried out by **Kasarabada Yashwanth Sai** bearing Roll number 220103012.

**Signature of Supervisor**

**(Dr. Sanasam Chanu Inunganbi)**

Signature of the Examiner 1 ............................

Signature of the Examiner 2 ............................

Signature of the Examiner 3 ............................

Signature of the Examiner 4 ............................

# Abstract

The Railway Accident Analytics: A Data-Driven AI Approach project aims to enhance railway safety and emergency preparedness through a comprehensive integration of data analytics, machine learning, and artificial intelligence. By leveraging historical railway accident data across India, the system facilitates meaningful insights into accident trends, seasonal and regional patterns, and causative factors through exploratory data analysis (EDA). A Power BI dashboard enables dynamic visualizations of accident frequency, geographic distribution, and response timelines. Furthermore, a machine learning model predicts accident severity and estimates key emergency response metrics, such as ambulance requirements and structural damage costs. An AI-powered chatbot assistant, built using the LLaMA-3 model, allows users to query accident-related data in real time for informed decision-making. The platform is designed using Streamlit for an intuitive, interactive user experience, providing stakeholders with an intelligent toolkit to support safety planning, resource allocation, and policy development within the railway sector satisfaction.

# Acknowledgement

In this submission, I have expressed my idea in my own words, and I have adequately cited and referenced any ideas or words that were taken from another source. I also declare that I adhere to all principles of academic honesty and integrity and that I have not misrepresented or falsified any ideas, data, facts, or sources in this submission. If any violation of the above is made, I understand that the institute may take disciplinary action. Such a violation may also engender disciplinary action from the sources which were not properly cited or permission not taken when needed.

(Yashwanth Sai Kasarabada)

# Contents

# List of Figures

# List of Tables

# Chapter 1

# Introduction

Railways constitute one of the most vital components of India's transportation infrastructure, facilitating the movement of over 23 million passengers and millions of tons of freight daily. Despite its critical importance and widespread reach, the Indian railway network remains prone to various types of accidents—ranging from derailments and collisions to fires and structural failures. These incidents not only result in tragic human casualties and economic losses but also severely disrupt essential public services and supply chains. In light of these challenges, there is a pressing need for an integrated, intelligent system that can analyze historical accident data, predict future risks, and assist in real-time decision-making.

The project titled "Railway Accident Analytics: A Data-Driven AI Approach" is a comprehensive initiative designed to address this need. By harnessing the power of data analytics, machine learning, and artificial intelligence, the system transforms raw accident records into meaningful insights and actionable predictions. The objective is not just to visualize what has happened, but to understand why it happened and what can be done to prevent it in the future.

At the heart of this platform lies a detailed exploratory data analysis (EDA) module. This module preprocesses railway accident data dating back to 1902 and categorizes it by various parameters such as time of day, season, geographic region, accident environment, and train type. With the help of Python libraries like Pandas, Seaborn, and Matplotlib, it uncovers temporal and spatial trends, revealing high-risk periods and regions, and identifying train categories most prone to accidents. These insights support evidence-based policymaking and operational improvements.

To ensure the information is not only insightful but also accessible, a Power BI dashboard has been developed and embedded within the Streamlit web interface. This dashboard offers interactive visualizations highlighting accident trends over the decades, regional hotspots, impacts of safety-related funding, and response efficiency. It empowers both technical users and non-specialists—such as policymakers, planners, and journalists—to interpret data visually and take informed action.

Moving beyond descriptive analytics, the project introduces a predictive modeling module powered by the Random Forest regression algorithm. This model estimates the severity of railway accidents based on input variables such as number of deaths, injuries, and rescue response time. Severity levels are categorized into "Very Low," "Low-Level," "Mid-Level," and "Critical," aiding in prioritizing emergency responses. Additionally, the model estimates the number of ambulances required and the expected structural damage cost using logically derived mathematical formulas. These predictive capabilities serve as valuable tools for disaster preparedness and resource planning.

Incorporating the latest in generative AI, the system also includes an AI Chatbot Assistant powered by the LLaMA-3 model via Groq's high-performance inference API. This assistant enables users to ask natural language questions regarding accident statistics, historical trends, safety concerns, and policy impacts. It responds with precise, context-aware answers backed by

data from the integrated dataset, enhancing accessibility for non-technical users.

The entire system is implemented using Streamlit, a Python-based web application framework that provides an intuitive and responsive user experience. Users can navigate seamlessly between the homepage, Power BI visualizations, EDA insights, predictive modeling interface, and AI chatbot. Its modular design ensures the platform is easy to maintain, extend, and deploy.

In summary, this project provides a unified, intelligent, and interactive solution for railway accident analysis in India. It empowers stakeholders with the tools needed to understand historical patterns, predict future risks, and implement timely safety measures. By bridging the gap between data science and public safety, the platform holds the potential to significantly reduce accident-related fatalities, optimize emergency responses, and contribute to a safer and more efficient railway system.

## 1.1 Outline

The Railway Accident Analytics project presents an end-to-end data-driven solution aimed at enhancing railway safety and accident preparedness in India. It integrates historical accident data, advanced analytics, machine learning, and generative AI to deliver a comprehensive web-based decision-support platform.

The system begins with data ingestion and exploratory data analysis (EDA). Users can upload accident datasets, which are then preprocessed and categorized by time of day, season, train type, region, and environment. Visual insights are generated using Python libraries to identify high-risk zones and patterns in accident occurrences.

A core component of the platform is an interactive Power BI dashboard, embedded within the app, offering visual representations of accident trends, casualty statistics, and geographical distributions. It aids policymakers and safety officials in grasping complex patterns through intuitive visuals.

For predictive analytics, the project features a machine learning model using Random Forest regression to estimate accident severity. The model predicts a severity score and classifies it into descriptive levels—Very Low, Low, Mid-Level, or Critical. It also estimates the number of ambulances required and the expected structural damage cost using transparent mathematical formulas.

To increase accessibility, the platform includes a generative AI assistant powered by the LLaMA-3 model via Groq API. Users can interact with this assistant to ask safety-related questions, explore accident records, and obtain context-aware insights.

The entire solution is built using Streamlit, offering a smooth, modular interface across Home, EDA, Power BI, Prediction, and AI Chatbot sections. This makes the platform user-friendly, scalable, and deployable for real-time use.

Overall, this project is a holistic attempt to bridge data science with public safety, equipping railway stakeholders with tools to analyze, predict, and respond to accidents efficiently.

## 1.2   Problem Statement

India's railway network, one of the largest and busiest in the world, plays a crucial role in connecting diverse regions and supporting the country's socio-economic development. However, the vast scale of operations also makes the system highly vulnerable to accidents—ranging from derailments and collisions to fires and infrastructure failures. These accidents often lead to severe consequences, including the loss of human lives, large-scale injuries, property damage, and service disruptions. Despite multiple safety initiatives, accident rates remain a concern, indicating the need for a more proactive and intelligent approach to railway safety management.

One of the core challenges lies in the fragmented and underutilized nature of accident data. Although historical accident records exist, they are seldom leveraged to extract actionable insights or develop predictive models. The absence of a centralized analytical platform prevents stakeholders from understanding root causes, identifying high-risk zones, estimating emergency response needs, or implementing data-informed safety measures. Furthermore, current systems lack accessible tools for real-time querying, severity assessment, and visual interpretation of accident patterns.

To address this, there is a pressing need for an integrated, data-driven platform that can perform exploratory analysis, visualize key metrics, predict accident outcomes, and assist users through an interactive AI interface. Such a system must enable decision-makers to explore historical trends, forecast severity levels, and allocate emergency resources effectively. It should also be intuitive enough for non-technical users to engage with, without requiring in-depth knowledge of data science.

This project, Railway Accident Analytics: A Data-Driven AI Approach, seeks to fill this gap by combining data analytics, machine learning, and generative AI into a unified application. The goal is to enhance railway safety preparedness, reduce human error, and support strategic interventions based on reliable, data-backed insights.

## 1.3 Motivation

India's vast railway network is an essential pillar of national connectivity and economic progress. With millions of passengers and freight movements daily, the railways face mounting pressure to ensure safety, efficiency, and responsiveness. However, recurring accidents—due to factors such as infrastructure failure, human error, or adverse environmental conditions—continue to pose serious threats to life and public confidence. These incidents often result in severe casualties, economic loss, and operational disruptions.

Although large volumes of historical accident data exist, they are often underutilized, fragmented, and limited to post-incident reporting. There is a lack of centralized, intelligent platforms that can convert this data into actionable insights or predict potential risks. Decision-makers currently have limited access to tools that allow proactive planning, severity estimation, or real-time response support.

The motivation for this project arises from the urgent need to bridge this analytical and operational gap using modern data science and AI technologies. By combining exploratory data analysis, predictive modeling, interactive dashboards, and a generative AI assistant, the project seeks to build a unified platform that enables comprehensive railway accident analysis and preparedness.

Predictive models, such as the Random Forest regression used here, offer the ability to assess severity levels based on past incidents, estimate ambulance requirements, and calculate structural damage costs. This supports smarter and quicker emergency response strategies. Meanwhile, the LLaMA-3-based AI chatbot brings accessibility to the forefront by allowing non-technical users to interact with the data through natural language queries.

This project also aligns with India's broader push toward digital transformation in critical infrastructure. It highlights how the fusion of domain knowledge, analytics, and user-centered design can create impactful, real-world solutions.

In essence, this initiative is driven by the goal of reducing accident impact through intelligent forecasting and informed decision-making—making Indian railways safer, smarter, and more resilient.

## 1.4 Purpose

The primary purpose of the Railway Accident Analytics: A Data-Driven AI Approach project is to develop an intelligent, interactive, and comprehensive platform that enhances railway safety through data-driven insights, predictive modeling, and AI-powered assistance. This project aims to bridge the gap between raw accident data and actionable decision-making by transforming complex datasets into meaningful visualizations, forecasts, and real-time query responses.

By utilizing historical accident records, the system is designed to uncover critical patterns and risk factors associated with railway incidents across India. The platform not only highlights when and where accidents are more likely to occur but also quantifies their severity and estimates the resources required for emergency response. Through this, it seeks to support timely, informed interventions by railway authorities, disaster management teams, policymakers, and public safety stakeholders.

A key objective of this project is to empower users—both technical and non-technical—with tools to explore accident data intuitively. By integrating exploratory data analysis, a Power BI dashboard, a machine learning-based severity prediction model, and an AI chatbot interface, the platform serves as a unified solution for railway accident analysis and preparedness. It facilitates proactive planning, optimizes rescue logistics, and improves situational awareness during crisis events.

In addition, the project promotes transparency and accessibility in railway safety by offering a user-friendly interface built using Streamlit. It ensures that insights and predictions are presented in a clear and comprehensible manner, enabling broader engagement with critical safety data.

Ultimately, the purpose of this project is to contribute to reducing railway accident-related casualties and damages by enabling smarter, faster, and more accurate decisions—thereby supporting a safer, more resilient railway system for the nation.

# Chapter 2

# Literature Survey

India's railway system is the fourth largest in the world, supporting a vast and complex network that facilitates the daily movement of millions of passengers and goods. Despite continuous expansion and modernization, railway accidents remain a persistent challenge—posing threats to life, property, and national infrastructure integrity. Recognizing the severity and recurrence of these incidents, numerous reports, historical archives, and government documents have highlighted the need for data-driven approaches in addressing railway safety. This literature survey consolidates insights from multiple credible sources and forms the analytical foundation for the *Railway Accident Analytics: A Data-Driven AI Approach* project, which seeks to support decision-makers with evidence-based tools.

## Accident Trends and Historical Context

The Wikipedia compilation of Indian rail accidents [1] presents a chronological account of major incidents from 1900 onward. It indicates not only the frequency of fatal accidents but also evolving causes—from infrastructure failures to human negligence and natural disasters. The dataset used in this project captures this entire timeline (1900–2024), enabling a rich analysis of long-term trends.

The Times of India [2] reports over 580 train accidents in a five-year period, with derailments alone accounting for more than half. These insights validate the project's use of a *Standard Accident Type* classification. Additionally, case-specific articles such as the Kochi collision [4] and Uttar Pradesh hijack incident [3] highlight the need for cause-based categorization—reflected in key fields within the dataset.

## Dataset Structure and Column Usage

The **main dataset** (1900–2024) includes columns such as: Year, Month, Day, Accident Name, Train Name, Location, State, Accident Rail Zone, Accident Type, Standard Accident Type, Cause, Deaths, Injuries, Rescue Time (hrs), Reforms/Changes, and Standard Cause Type. It is used in the **Power BI dashboard** for trend visualization and in the **AI chatbot** for query resolution.

The **Insights and Analysis** module uses a preprocessed subset from 2002–2017 for higher data quality and recency. It includes: X, environment, train_name, injured, killed, triggering_factor, time, railway_division, cause, and env—supporting exploratory analysis and visual interpretation.

The **Predictive Model** utilizes the full timeline but filters for columns critical to emergency estimation: Standard Accident Type, Deaths, Injuries, Rescue Time (hrs), Severity, Estimated Ambulances Required, and Estimated Structural Damage Cost (INR).

## Policy Reports and Investment Trends

Government documents like the Indian Railways White Paper [5] and Annual Report [6]

highlight infrastructural gaps, outdated signaling, and investment shortages. These reports emphasize the need for smarter tools to analyze and correlate funding with safety trends—a feature addressed by this project's dashboard module.

## Technological Advancements and Innovation Opportunities

Recent studies have explored a range of AI and ML techniques for railway safety. Chhotu and Suman [7] used ML algorithms to predict fatalities at level crossings, while their follow-up study [8] applied neural networks and logistic regression for accident severity forecasting. Similarly, Tejaswin and Kumar [9] demonstrated the utility of Python-based analytics for summarizing accident records in India.

Warning systems are being tested to improve locomotive awareness, yet a unified backend to operationalize such data is missing. Unsupervised methods for safety monitoring at stations were explored by Reddy and Triveni [10], indicating potential in anomaly detection and preventive alert systems.

This project aligns with these emerging trends by combining forecasting, analytics, and AI interfaces into a cohesive system.

## Gap Identification in Literature

Most existing tools focus narrowly—either on ML for specific predictions, or post-incident data reports. They often lack user-friendly interfaces for decision-makers. Moreover, integration of forecasting, interactive dashboards, and natural language query interfaces is largely absent.

This project addresses these gaps by:

- Providing a modular Streamlit app with a full ML prediction pipeline.

- Offering real-time interactive dashboards via Power BI.

- Enabling intuitive AI-powered Q&A for railway safety statistics.

The result is a unified decision-support system intended for use by Indian government officials, planners, and transportation safety boards.

# Chapter 3

# Requirement Engineering

## 3.1  Introduction

Requirement engineering is the disciplined process of defining, documenting, and maintaining requirements in the software development lifecycle. It ensures the system meets stakeholder needs—in this case, government officials and policy-makers seeking actionable insights on railway safety in India. This section outlines the system's functional capabilities, quality attributes, data specifications, and project constraints.

## 3.2  Functional Requirements

Functional requirements specify the behaviors the system must exhibit. Table 3.1 summarizes these requirements.

Table 3.1: Functional Requirements

| ID | Requirement Description |
|---|---|
| FR1 | Ingest historical accident dataset (1900–2024) with key columns: Year, Month, Day, Accident Name, Train Name, Location, State, Accident Rail Zone, Accident Type, Standard Accident Type, Cause, Deaths, Injuries, Rescue Time (hrs), Reforms/Changes, Standard Cause Type. |
| FR2 | Preprocess and create subsets: Insights and Analysis (2002–2017) and Predictive Model dataset for severity, ambulance, and cost estimation. |
| FR3 | Provide a Home Overview page for dataset download, project overview, and navigation. |
| FR4 | Embed interactive Power BI dashboards visualizing trends, regional hotspots, casualty statistics, and funding impacts. |
| FR5 | Implement EDA section with bar charts for distributions by time-of-day, season, region, environment, and train type. |
| FR6 | Train and expose a Random Forest regression model for accident severity prediction and resource estimation. |
| FR7 | Display prediction outputs: severity score, severity level, estimated ambulances, and structural damage cost with LaTeX-rendered formulas. |
| FR8 | Integrate AI Chatbot (LLaMA-3 via Groq API) to answer natural-language queries on accident data and safety recommendations. |
| FR9 | Enable modular navigation across Home, Power BI Report, Insights and Analysis, Predictive Model, and AI Assistant sections. |

## 3.3  Non-Functional Requirements

Non-functional requirements define system qualities and constraints. Table 3.2 lists these requirements.

Table 3.2: Non-Functional Requirements

| ID | Requirement Description |
|---|---|
| NFR1 | **Performance:** Visualizations and AI responses must render within 2 seconds for datasets 1M records. |
| NFR2 | **Usability:** The Streamlit interface must be intuitive for non-technical users, with clear labels and consistent styling. |
| NFR3 | **Scalability:** Support incremental data updates and concurrent user access without performance degradation. |
| NFR4 | **Reliability:** Ensure 99.5% uptime over 30 days, with error logging and fallback messages for dependencies. |
| NFR5 | **Security:** Use HTTPS, input validation, and role-based access control to protect sensitive data. |
| NFR6 | **Accessibility:** Comply with WCAG 2.1 AA (keyboard navigation, high-contrast themes). |
| NFR7 | **Maintainability:** Adhere to PEP 8, use modular code structure, and provide inline documentation. |
| NFR8 | **Portability:** Deployable on Azure, AWS, or on-premises without code modifications. |
| NFR9 | **Auditability:** Maintain immutable logs of data uploads, model training, and user queries. |

## 3.4  Data Requirements

The system relies on datasets with records from 1900 to 2024. Key data requirements include:

- **Source Verification:** Data must originate from verified sources (government reports, reputable news).

- **Schema Conformity:** Each dataset must adhere to defined schemas for main and subset data.

- **Completeness:** Achieve 90% completeness for the 2002–2017 subset on required fields.

- **Version Control:** Timestamp and document each data version in a changelog.

- **Retention Policy:** Archive historical data while providing easy access to recent records.

## 3.5  Constraints

Project constraints that affect development:

- **Data Quality:** Early records (1900–2001) may be incomplete or inconsistent.

- **Model Updates:** Real-time retraining is out of scope; use periodic batch retraining.

- **Embedding Dependencies:** Power BI iframes require valid credentials and network connectivity.

- **Device Support:** Primary access via desktop; mobile support is limited.

- **API Management:** Secure storage and periodic rotation of API keys.

## 3.6 Conclusion

This requirement engineering section captures the essential functional and non-functional requirements, data specifications, and project constraints for the *Railway Accident Analytics* platform. These requirements lay the groundwork for the subsequent System Design and Implementation, ensuring the final product aligns with stakeholder objectives and quality standards.

# Chapter 4

# System Design

## 4.1 Architecture Overview

Figure ?? presents the high-level architecture of the Railway Accident Analytics platform. It consists of five main modules: Home Overview, Power BI Dashboard, Insights and Analysis, Predictive Model, and AI Assistant, all built on top of the Streamlit framework. Each module is accessible via a sidebar menu, allowing seamless user navigation.
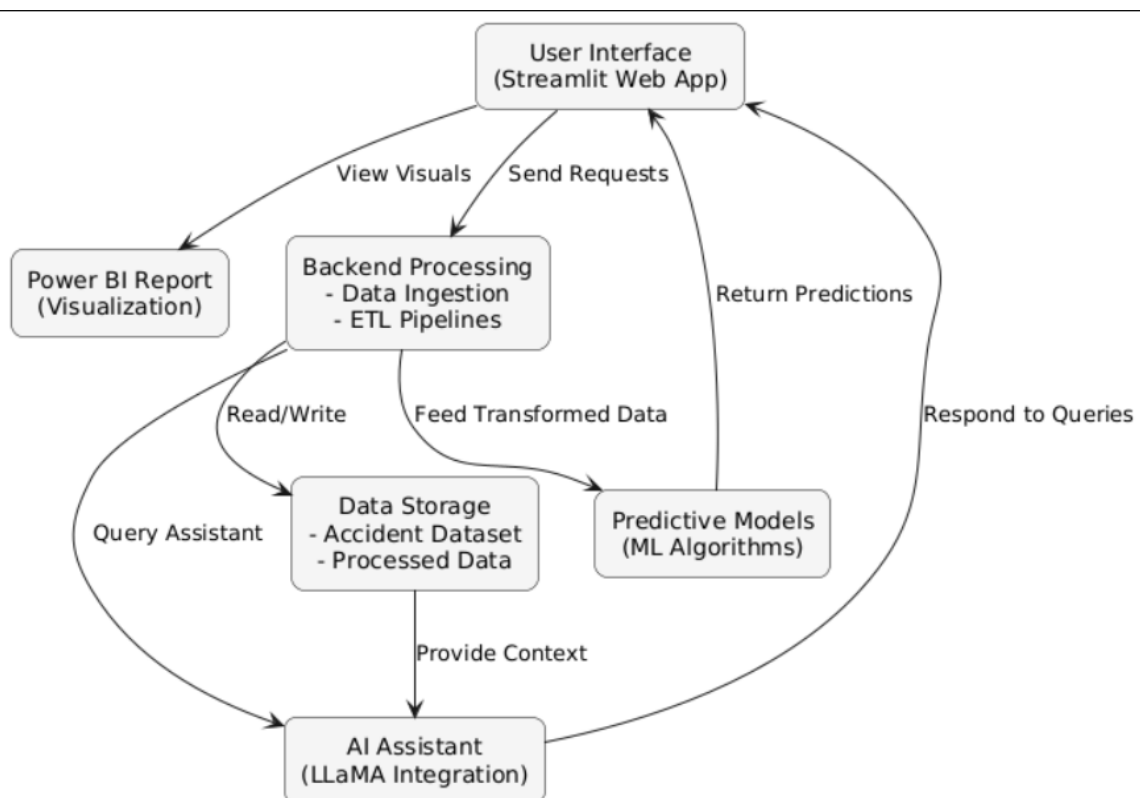


Figure 4.1: System Architecture of the Railway Accident Analytics Platform

## 4.2 Component Design

Each component of the system serves a specific purpose in delivering analytical, predictive, and assistive functionalities.

### 4.2.1  Home Overview Module

- Provides a welcoming interface with an overview of the platform's features and objectives.

- Allows users to download sample datasets via an interactive button using Streamlit's `download_button()` function.

- Facilitates smooth navigation to other modules through the sidebar.

### 4.2.2  Power BI Dashboard Module

- Embeds live Power BI dashboards through iframe integration to display comprehensive accident statistics and visual summaries.

- Covers accident distribution by location, year, cause, and severity.

- Enables decision-makers to gain insights through pre-built, scrollable, and zoomable charts.

- Requires an active internet connection and embedded Power BI URL to function correctly.

### 4.2.3  Insights and Analysis Module

- Focuses on accidents from 2002 to 2017, using cleaner and more consistent data.

- Allows users to visualize patterns by time of day, railway division, train environment, and season using Pandas, Matplotlib, and Seaborn.

- Empowers users to identify risk zones and vulnerable operating conditions for policy-level interventions.

- Users can interactively view plots filtered by selected dimensions.

### 4.2.4  Predictive Model Module

- Uses selected columns from the full dataset (1900–2024) to train a Random Forest regression model using scikit-learn.

- Accepts user inputs via the UI for accident type, death count, injury count, and rescue time.

- Outputs predicted severity score, severity classification, estimated number of ambulances, and damage cost.

- Enables CSV uploads directly in the UI to perform batch predictions without needing back-end storage or APIs.

### 4.2.5  AI Assistant Module

- Integrates with the Groq-hosted LLaMA-3 model to respond to user queries via natural language.

- Allows users to inquire about accident statistics, severity trends, or causes in an interactive, chat-like interface.

- Facilitates exploratory dialogue for non-technical users to understand data without diving into visualizations or raw numbers.

- Operates within a simple Streamlit form-based chat interface.

## 4.3 Data Flow and Sequence

This section describes the system's operational flow and data movement between components. It includes two visual representations:

- A flowchart illustrating the sequence of processes, decisions, and outcomes.

- A Data Flow Diagram (DFD) showing the interaction of external entities, processes, and data stores.

### 4.3.1 Flowchart: System Processing and AI Query Flow

Figure 4.2 depicts the sequential flow of the system — starting from data input, preprocessing, exploratory data analysis, and predictive modeling to decision-making, AI assistant interaction, resource allocation, and reporting.



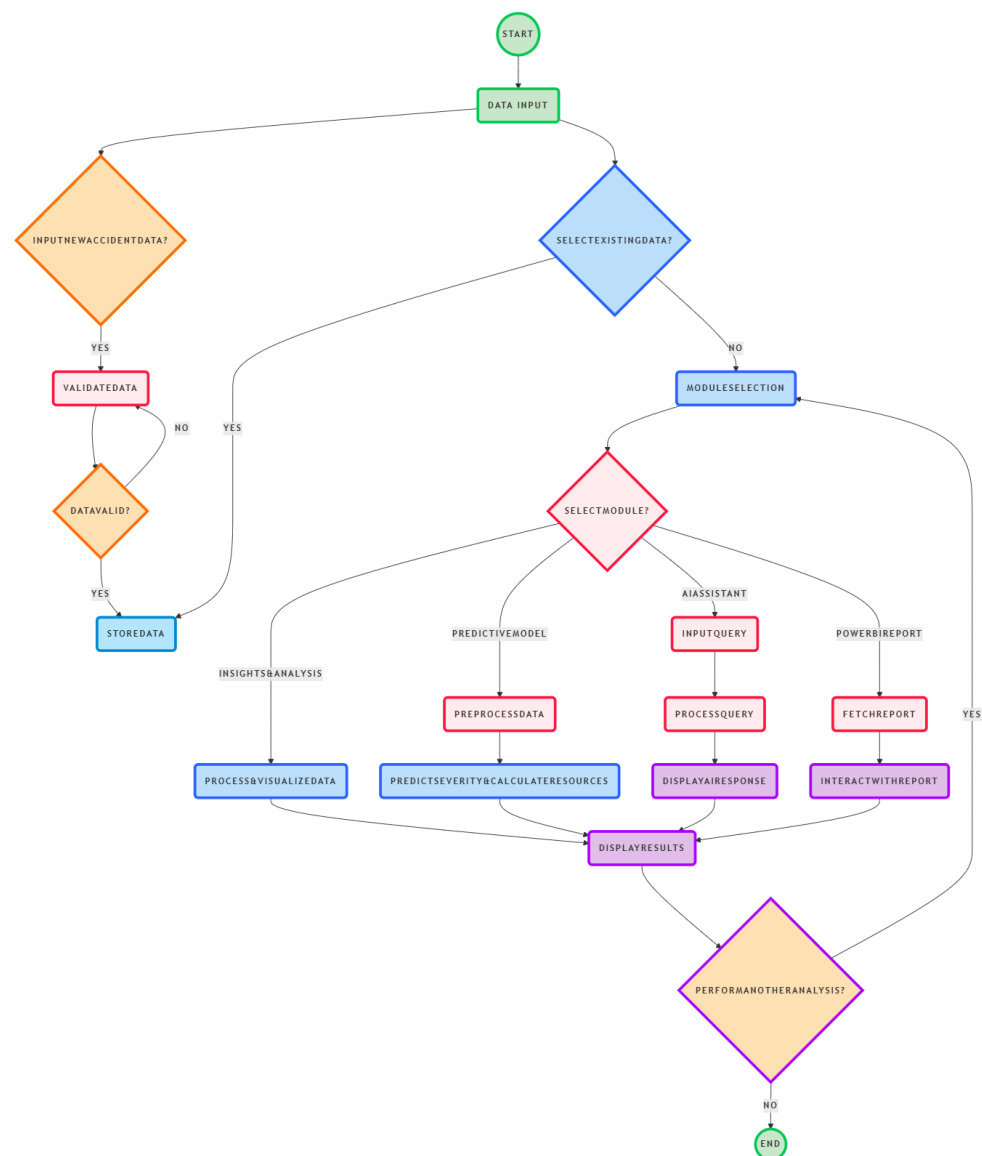Figure 4.2: System Flowchart: Data Processing and AI Query Handling

### 4.3.2 Level 1 Data Flow Diagram

Figure 4.3 presents the data flow diagram (Level 1 DFD) showing how data moves between users, external systems, and internal processes in the application.
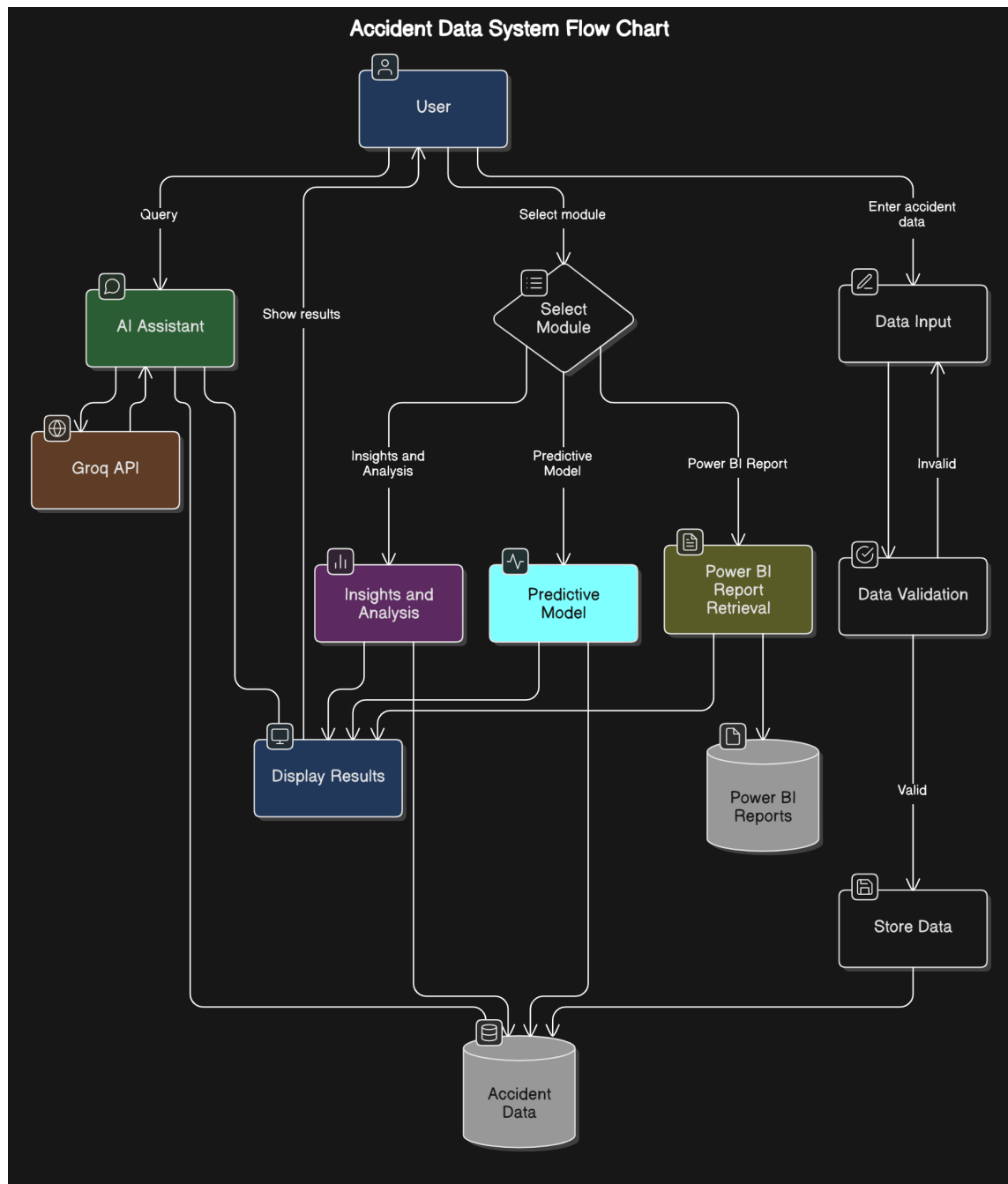


Figure 4.3: Level 1 Data Flow Diagram (DFD)

## 4.4 Technology Stack

Table 4.1 outlines the technologies used in each layer of the system.

Table 4.1: Technology Stack

| Layer | Technology/Tool |
|-------|-----------------|
| Frontend | Streamlit (Python-based UI framework) |
| Visualization | Power BI Embedded, Matplotlib, Seaborn |
| Data Processing | Python, Pandas, NumPy |
| Machine Learning | scikit-learn (Random Forest), SciPy |
| AI Integration | Groq API, LLaMA-3 model |
| Deployment | Streamlit Cloud (for web hosting) |
| Security | Streamlit session handling, secure API token management |
| Data Storage | In-memory upload (CSV files via UI); no external SQL database used |

## 4.5   Deployment and Scalability

The application is deployed using Streamlit Cloud, eliminating the need for Docker or Kubernetes configurations. Users upload CSV files directly within the UI for real-time predictions and visualization, removing backend storage dependencies. While the current deployment suits interactive academic and research usage, the modular architecture supports future extensions, including integration with external storage or APIs.

## 4.6   Summary

This System Design chapter details the architectural layout, component-level design, data flow, and technology choices underpinning the Railway Accident Analytics platform. It reflects a modular and user-friendly architecture optimized for policy analysis, public safety visualization, and predictive insights using real-time data inputs.

# Chapter 5

# Implementation

## 5.1 Development Tools and Environment

The Railway Accident Analytics platform was developed using modern data analysis, machine learning, and web development tools. The environment setup includes:

- **Programming Language:** Python 3.11 (Interpreted language suited for rapid prototyping)

- **IDE:** Visual Studio Code (Configured with Python extensions for debugging and linting)

- **Framework:** Streamlit (Streamlit: The fastest way to build and share data apps)

- **Visualization Tools:** Power BI, Matplotlib, Seaborn (Dynamic visual representation of results)

- **Core Libraries:** Pandas, scikit-learn, NumPy, Requests

- **ML Models:** RandomForestRegressor (Ensemble learning for accurate predictions and robust performance)

- **AI Model:** LLaMA-3 via Groq API (High-throughput generative model for Natural Language Processing, high context models allow zero-shot, few-shot, and complex reasoning)

## 5.2 Code Structure

The project employs a modular structure, organized into the following key components:

- **main.py:** Manages overall application structure and page routing.

- **Home_overview.py:** Delivers a project introduction and dataset download functionality.

- **Insights_and_Analysis.py:** Conducts EDA, providing visualization of the training data subset.

- **Power_BI_Report.py:** Integrates live Power BI dashboards through secure iframe embedding.

- **Predictive_Model.py:** Implements the Random Forest algorithm, including data preprocessing, model training, and severity prediction features.

- **llama_Assitant.py:** Connects to the LLaMA-3 model for intelligent question answering based on railway accident data.

Each module is loaded dynamically through the Streamlit sidebar using the 'show()' function.

## 5.3 Core Modules and Functionalities

### 5.3.1 Main Application (main.py)

This script manages the overarching framework and integrates the different modules of the project.

**Imports**

```python
import streamlit as st
from streamlit_option_menu import option_menu
from pages import Home_overview
from pages import Insights_and_Analysis
from pages import Power_BI_Report
from pages import Predictive_Model
from pages import llama_Assistant
```

Listing 5.1: Imports in main.py

*Explanation:* This code imports `streamlit` and related libraries essential for creating the user interface, managing the menu, and referencing individual page modules. Notably, `st` and `option_menu` are crucial for building the UI components and navigation sidebar.

**Page Configuration**

```python
st.set_page_config(page_title="Railway Accident Analytics",
                   layout="wide",
                   initial_sidebar_state="collapsed")
```

Listing 5.2: Page Configuration in main.py

*Explanation:* This configuration sets up the Streamlit page using the `st` object, setting the title, using a wide layout for enhanced visual appeal, and a collapsed sidebar by default for cleaner initial presentation.

**Main Title and Styling**

```python
st.markdown(
    """
    <style>
    .section-heading {
        font-size: 2.5em;
        color: #0077cc;
        font-weight: bold;
        padding-bottom: 0.2em;
        border-bottom: 2px solid #0077cc;
        margin-bottom: 0.5em;
    }
    </style>
    """,
    unsafe_allow_html=True
)
st.markdown('<div class="section-heading">Railway Accident Analytics: A Data-
    Driven AI Approach</div>',
            unsafe_allow_html=True)
```

Listing 5.3: Main Title and Styling in main.py

*Explanation:* This segment incorporates HTML and CSS styles to customize the main title using `st.markdown`. The custom styles enhance the visual hierarchy and draw attention to the core purpose of the application.

**Horizantal UI Menu**

```python
with st.sidebar:
    selected = option_menu("Main Menu",
                           ["Home Overview",
                            "Insights and Analysis",
                            "Power BI Report",
                            "Predictive Model",
                            "AI Chatbot"],
                           icons=['house', 'bar-chart', 'file-earmark-bar-graph', 'graph-up', 'robot'],
                           menu_icon="cast",
                           default_index=0)
```

Listing 5.4: Sidebar Menu in main.py

*Explanation:* The `option_menu` in the sidebar facilitates navigation to different parts of the application. The use of icons improves the user experience by providing visual cues for each module.

**Page Selection**

```python
if selected == "Home Overview":
    Home_overview.show()
elif selected == "Insights and Analysis":
    Insights_and_Analysis.show()
elif selected == "Power BI Report":
    Power_BI_Report.show()
elif selected == "Predictive Model":
    Predictive_Model.show()
elif selected == "AI Chatbot":
    llama_Assitant.show()
```

Listing 5.5: Page Selection in main.py

*Explanation:* Based on the user's selected menu item, the corresponding 'show()' function from the respective module is invoked. This structure is essential for the application's modular design, where each module runs independently when selected, enhancing maintainability and scalability.

### 5.3.2 Power BI Report Integration(Power BI Section.py)

This module integrates interactive Power BI dashboards using secure iframe embedding.

**Implementation**

```python
def show():
    st.title("       Power BI Analytics Dashboard")

    st.write("""
    ### Interactive Power BI Dashboard
    Explore comprehensive railway safety metrics through our integrated Power BI reports.
    """)

    # Embedded Power BI Report 1
    st.write("### Accident Trends Dashboard")
    components.iframe(
        "https://app.powerbi.com/view?r=eyJrIjoiNDg2MmU...",
        height=600,
        scrolling=True
    )
```

```
16
17    # Embedded Power BI Report 2
18    st.write("### Regional Safety Analysis")
19    components.iframe(
20        "https://app.powerbi.com/view?r=eyJrIjoiZDU3MjA...",
21        height=600,
22        scrolling=True
23    )
```

Listing 5.6: Power BI Report Integration in Power BI Report.py

*Explanation:*

- Uses Streamlit Components (`components.iframe`) for secure embedding of Power BI dashboards

- Implements two interactive dashboards showing accident trends and regional safety metrics

- Maintains Power BI's native interactivity (filtering, drill-down) within the Streamlit interface

- Includes security measures through Power BI's embedded token system

### 5.3.3 Insights and Analysis (Insights_and_Analysis.py)

This module focuses on conducting Exploratory Data Analysis (EDA) and generating visualizations to provide insights into the railway accident data.

**Imports**

```
1  import streamlit as st
2  import pandas as pd
3  import matplotlib.pyplot as plt
4  import seaborn as sns
```

Listing 5.7: Imports in Insights and Analysis.py

*Explanation:* Essential libraries are imported for creating the user interface using `st`, handling data with `pd`, and producing high-quality visualizations with `plt` and `sns`.

`show()` **Function**

```
1  def show():
2      st.markdown(...)
3      st.write(...)
4      ...
```

Listing 5.8: show() Function in Insights and Analysis.py

*Explanation:* This function encapsulates the primary logic for the EDA page. It structures the layout, loads the data, and calls functions to produce visualizations.

**Data Loading**

```python
@st.cache_data
def load_data():
    try:
        df = pd.read_csv("C:/Users/yashw/Desktop/6 th sem project ideas/
    Indian_Railways_Accidents_Dataset_1902_2024.csv")
        return df
    except FileNotFoundError:
        st.error("      Railway accident dataset not found!")
        return None
df = load_data()
```

Listing 5.9: Data Loading in Insights and Analysis.py

*Explanation:* The `load_data` function reads the dataset and caches it using `@st.cache_data`. Error handling is implemented to manage cases where the data file is missing.

**Plotting Function**

```python
def plot_column(df, col, question):
    fig, ax = plt.subplots(figsize=(10, 6))
    sns.barplot(x=df[col].value_counts().index,
                y=df[col].value_counts().values,
                ax=ax,
                palette='muted')
    ax.set_xlabel(col.replace('_', ' ').title(), fontsize=12)
    ax.set_ylabel("Accident Count", fontsize=12)
    ax.set_title(question, fontsize=14)
    plt.xticks(rotation=45, ha='right')
    plt.tight_layout()
    st.pyplot(fig)

    most_common = df[col].value_counts().idxmax()
    st.write(f"Most accidents occur in: {most_common}",
             unsafe_allow_html=True)
```

Listing 5.10: Plotting Function in Insights and Analysis.py

*Explanation:* The `plot_column` function creates bar plots for visualizing various aspects of the dataset, including the frequency of accident causes, types, and locations. The function uses `sns` and `plt` to create and display the plots, and `st` to write the most common occurrences on the plot.

### 5.3.4 Predictive Model (Predictive_Model.py)

The `Predictive_Model.py` script implements machine learning techniques for accident severity prediction.

**Imports**

```python
import streamlit as st
import pandas as pd
import numpy as np
import math
from sklearn.model_selection import train_test_split
from sklearn.ensemble import RandomForestRegressor
from sklearn.preprocessing import OneHotEncoder, StandardScaler
from sklearn.compose import ColumnTransformer
from sklearn.pipeline import Pipeline
```

```
10  from sklearn.impute import SimpleImputer
11  from sklearn.metrics import mean_absolute_error, r2_score
```

Listing 5.11: Imports in Predictive Model.py

*Explanation:* This segment imports libraries necessary for data manipulation, machine learning, and UI creation, forming the foundation for building the predictive model. The libraries include:

- `streamlit` for creating the UI.

- `pandas` and `numpy` for data manipulation.

- `math` for mathematical functions.

- `sklearn` modules for model building, preprocessing, and evaluation.

`load_data()` **Function**

```
1   @st.cache_data
2   def load_data(uploaded_file):
3       try:
4           df = pd.read_csv(uploaded_file)
5           st.session_state.original_df = df.copy()
6           return df
7       except Exception as e:
8           st.error(f"Error loading data: {e}")
9           return None
```

Listing 5.12: load data() Function in Predictive Model.py

*Explanation:* This function is used to load data from a CSV file uploaded by the user and stores it in Streamlit's session state using `st`. The use of `@st.cache_data` ensures that data is cached, preventing redundant loading and improving the performance of the application.

`preprocess_data()` **Function**

```
1   def preprocess data(df):
2       df = df.copy()
3
4       if TARGET_SEVERITY not in df.columns:
5           df[TARGET_SEVERITY] = df['Deaths'] + df['Injuries']
6
7       categorical_features = ['Standard Accident Type']
8       numerical_features = ['Deaths', 'Injuries', 'Rescue Time (hrs)']
9
10      numeric_transformer = Pipeline([
11          ('imputer', SimpleImputer(strategy='mean')),
12          ('scaler', StandardScaler())
13      ])
14
15      categorical_transformer = Pipeline([
16          ('imputer', SimpleImputer(strategy='most_frequent')),
17          ('onehot', OneHotEncoder(handle_unknown='ignore'))
18      ])
19
20      preprocessor = ColumnTransformer([
21          ('num', numeric_transformer, numerical_features),
22          ('cat', categorical_transformer, categorical_features)
23      ])
24
25      X = df[FEATURES]
```

```
26    X_processed = preprocessor.fit_transform(X)
27
28    st.session_state.preprocessor = preprocessor
29    st.session_state.feature_names = preprocessor.get_feature_names_out()
30
31    return pd.DataFrame(X_processed, columns=st.session_state.feature_names),
      df[TARGET_SEVERITY]
```

Listing 5.13: preprocess data() Function in Predictive Model.py

*Explanation:* This function is crucial for preparing the data for machine learning. It handles several preprocessing steps, including:

- Computing severity (if missing) by summing deaths and injuries, using basic arithmetic operations.

- Defining numerical and categorical features.

- Creating preprocessing pipelines for numerical and categorical data using `Pipeline` from `sklearn`. These pipelines handle missing values and scale numerical and categorical features.

- Using `ColumnTransformer` to apply the appropriate transformers to each feature type.

`train_severity_model()` **Function**

```
1  def train_severity_model(X, y):
2      X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
       random_state=42)
3
4      model = RandomForestRegressor(random_state=42)
5      model.fit(X_train, y_train)
6
7      y_pred = model.predict(X_test)
8      mae = mean_absolute_error(y_test, y_pred)
9      r2 = r2_score(y_test, y_pred)
10
11     st.session_state.severity_model = model
12     st.session_state.severity_mae = mae
13     st.session_state.severity_r2 = r2
```

Listing 5.14: train severity model() Function in Predictive Model.py

*Explanation:* The `train_severity_model` function trains the machine learning model. It performs the following steps:

- Splits the dataset into training and test sets using `train_test_split` with a test size of 20%.

- Initializes and trains a `RandomForestRegressor` model with a specified `random_state` for reproducibility.

- Predicts the test set and calculates the Mean Absolute Error (MAE) and R-squared ($R^2$) score using metrics from `sklearn`.

The Mean Absolute Error (MAE) is calculated using the formula:

$$MAE = \frac{1}{n} \sum_{i=1}^{n} |y_i - \hat{y}_i|$$

where $y_i$ is the actual value, $\hat{y}_i$ is the predicted value, and $n$ is the number of data points.

The R-squared ($R^2$) score is calculated using the formula:

$$R^2 = 1 - \frac{\sum_{i=1}^{n}(y_i - \hat{y}_i)^2}{\sum_{i=1}^{n}(y_i - \bar{y})^2}$$

*where $\bar{y}$ is the mean of the observed data. The MAE and $R^2$ are stored and used in* `st.session_state` *to store the model's performance metrics.*

## Math Behind the Model

*This section details the mathematical foundations of the predictive model.*

**Random Forest Regression**  The predictive model uses the Random Forest Regression algorithm. The mathematical steps are as follows:

1. **Bootstrap Sampling:** For each tree, a random sample is generated (with replacement) from the dataset.

2. **Random Feature Selection:** At each split in a tree, a subset of features is randomly selected. For regression, the number of features, $m$, is typically $p/3$ where $p$ is the total number of features.

3. **Tree Construction:** Each tree is grown to its maximum depth without pruning. At each split, the feature and threshold that minimize the Mean Squared Error (MSE) are chosen:

$$MSE = \frac{1}{N}\sum_{i=1}^{N}(y_i - \hat{y}_i)^2$$

where $y_i$ is the true value and $\hat{y}_i$ is the predicted value for sample $i$.

4. **Prediction Aggregation:** For a new input, each tree makes a prediction, and the final prediction is the average:

$$\hat{y}_{RF} = \frac{1}{T}\sum_{t=1}^{T}\hat{y}^{(t)}$$

where $T$ is the number of trees, and $\hat{y}^{(t)}$ is the prediction from tree $t$.

**Severity Level Assignment (Thresholding)**  The predicted severity score is mapped to a qualitative level using thresholds:

$$\begin{cases} \text{Critical,} & \text{if Predicted Severity Score} \geq 75 \\ \text{Mid-Level,} & \text{if } 50 \leq \text{Predicted Severity Score} < 75 \\ \text{Low-Level,} & \text{if } 25 \leq \text{Predicted Severity Score} < 50 \\ \text{Very Low,} & \text{if Predicted Severity Score} < 25 \end{cases}$$

**Ambulance Estimation Formula**  Resource needs are estimated using a weighted sum:

$$\text{Estimated Ambulances} = \lceil (\text{Deaths} \times 0.3) + (\text{Injuries} \times 0.5) + (\text{Rescue Time (hrs)} \times 0.8) \rceil$$

where $\lceil \cdot \rceil$ is the ceiling function (round up to the nearest integer).

**Structural Damage Cost Estimation**  The cost is estimated as:

$$\text{Estimated Cost} = \text{Base Cost} \times \text{Multiplier}$$

The base cost depends on the accident type, and the multiplier depends on the severity score.

### 5.3.5 AI Chatbot (llama_Assitant.py)

**AI Interaction**

```python
def chat_with_ai(prompt):
    data_summary = f"The available data on Indian railway accidents spans from
    1902 to 2024. The model MUST base its answers primarily on this dataset, but
     can use outside information if required. The model MUST base its answers
    primarily on this dataset, but can use outside information if required. Do
    not hallucinate data, instead respond that the data does not exist."
    enhanced_prompt = f"{prompt}. {data_summary}" # Augment the prompt
    headers = {"Authorization": f"Bearer {GROQ_API_KEY}", "Content-Type": "
    application/json"}
    data = {
        "model": "llama3-8b-8192",
        "messages": [{"role": "user", "content": enhanced_prompt}],
        "max_tokens": 750 # Increased token limit

    response = requests.post(GROQ_API_URL, json=data, headers=headers)
    try:
        response_json = response.json()
        return response_json["choices"][0]["message"]["content"] if "choices"
    in response_json else f"      API Error: {response_json}"
    except Exception as e:
        return f"       Error: {str(e)}"
```

Listing 5.15: AI Function in llama Assitant.py

*Explanation:* This section defines a function that crafts the system call, enhances the user's prompt with data and information about its bounds, and specifies the model as well as the maximum allowable tokens. This interacts with the `st` framework and generates the response, displaying the data to the user.

## 5.4   Error Handling

The application implements robust error handling to ensure reliability and a positive user experience. Key error-prone areas, such as data loading and API interactions, are protected using `try-except` blocks.

   **Data Loading Errors:** In the `load_data` functions (found in `Predictive_Model.py` and `Insights_and_Analysis.py`), the application anticipates `FileNotFoundError` exceptions. If the specified dataset cannot be found at the specified path, `st.error` is called to display a user-friendly error message, prompting the user to verify the file path or upload the correct data file.

   **Model Training Errors:** The `train_severity_model` function includes basic data validation checks to ensure that the input data meets the expected format and requirements. If critical features are missing or the data types are incompatible, informative error messages are displayed using `st.warning` or `st.error`, guiding the user to correct the input.

   **API Interaction Errors:** The `chat_with_ai` function (in `llama_Assitant.py`) anticipates network errors or issues with the Groq API. `try-except` blocks handle `requests.exceptions.RequestExceptio` and other potential errors during the API call. In case of failure, appropriate error messages are displayed using `st.error` to inform the user of connectivity problems or API-related issues.

## 5.5   User Interface and Experience

The user interface (UI) is carefully designed to provide an intuitive, responsive, and visually appealing experience. Streamlit's features are leveraged to create a seamless interaction flow and present key information effectively.

**Navigation and Layout:** The primary navigation is facilitated by a horizontal menu created using `streamlit_option_menu`, as seen in `main.py`:

```python
menu = ["Home", "Power BI Report", "Insights and Analysis", "Predictive Model",
    "AI Assistant"]
icons = ["house", "bar-chart", "graph-up", "cpu", "robot"]
selected_section = option_menu(
    menu_title="",  # No title to keep it compact
    options=menu,
    icons=icons,
    menu_icon="cast",
    default_index=0,
    orientation="horizontal"
)
```

Listing 5.16: Horizontal Navigation Menu

This horizontal menu, positioned at the top, ensures that users can easily switch between the main sections of the application. The icons provide visual cues, enhancing the overall usability.

**Data Presentation:** Raw data previews, preprocessed data samples, and model performance metrics are displayed using `st.dataframe`, providing users with a clear view of the data used in the analysis. Informative messages and styled titles are used to guide users and provide context.

**Interactive Elements:** Users can upload data files, train the predictive model, and make predictions through interactive elements such as `st.file_uploader`, `st.button`, and `st.number_input`. Progress bars are used during model training to provide real-time feedback to the user.

## 5.6    Testing and Validation

Thorough testing and validation procedures were implemented to ensure the reliability and accuracy of the application's functionality.

**Model Validation:** The predictive model's performance was evaluated using standard metrics such as Mean Absolute Error (MAE) and R-squared ($R^2$), calculated using functions from `sklearn.metrics`. These metrics are displayed in the UI to provide users with transparency regarding the model's accuracy and reliability.

**Functional Testing:** Each module underwent functional testing to ensure it performed as expected. For instance, data loading functions were tested with various CSV files to ensure compatibility and error handling. The AI Chatbot's responses were manually validated to assess their relevance and accuracy.

**User Scenario Testing:** The application was tested using various user scenarios to ensure that it met the intended use cases and provided a positive user experience. This included simulating different user roles, data input methods, and prediction requests.

## 5.7    Limitations and Future Improvements

While the Railway Accident Analytics platform provides significant functionality, there are inherent limitations that could be addressed in future iterations.

**Data Dependence:** The predictive model and exploratory data analysis (EDA) depend heavily on the quality and completeness of the historical data. The model may not accurately predict severity or outcomes in scenarios that differ significantly from the training data. This limitation can be mitigated by expanding and updating the dataset.

**AI Chatbot Contextual Understanding:** The AI Chatbot's ability to interpret complex queries and provide nuanced responses is limited by the model's size and training data. To improve the chatbot's accuracy and relevance, future work could involve fine-tuning the LLaMA-3

model on a larger, more diverse set of railway accident data.

**Future Directions:**

- **Real-Time Data Integration:** Integrate real-time data sources, such as weather conditions, traffic patterns, and sensor data from railway infrastructure, to improve the accuracy and timeliness of the analysis.

- **Enhanced AI Capabilities:** Train the AI chatbot on a broader dataset of railway-related knowledge to improve its ability to answer complex queries and provide personalized recommendations.

- **Advanced Machine Learning Techniques:** Implement more sophisticated machine learning models, such as deep learning algorithms, to capture complex patterns and improve the accuracy of severity predictions.

# Chapter 6

# Results

## 6.1 Introduction

This chapter presents the outcomes derived from various components of the Railway Accident Analytics platform. It covers insights from exploratory data analysis, evaluates the machine learning model's predictive performance, and includes observations from both the Power BI dashboard and the AI-powered chatbot. The platform aims to support government authorities in formulating effective, data-backed decisions for improving railway safety standards. Each result segment is complemented with UI screenshots to reflect the actual implementation.
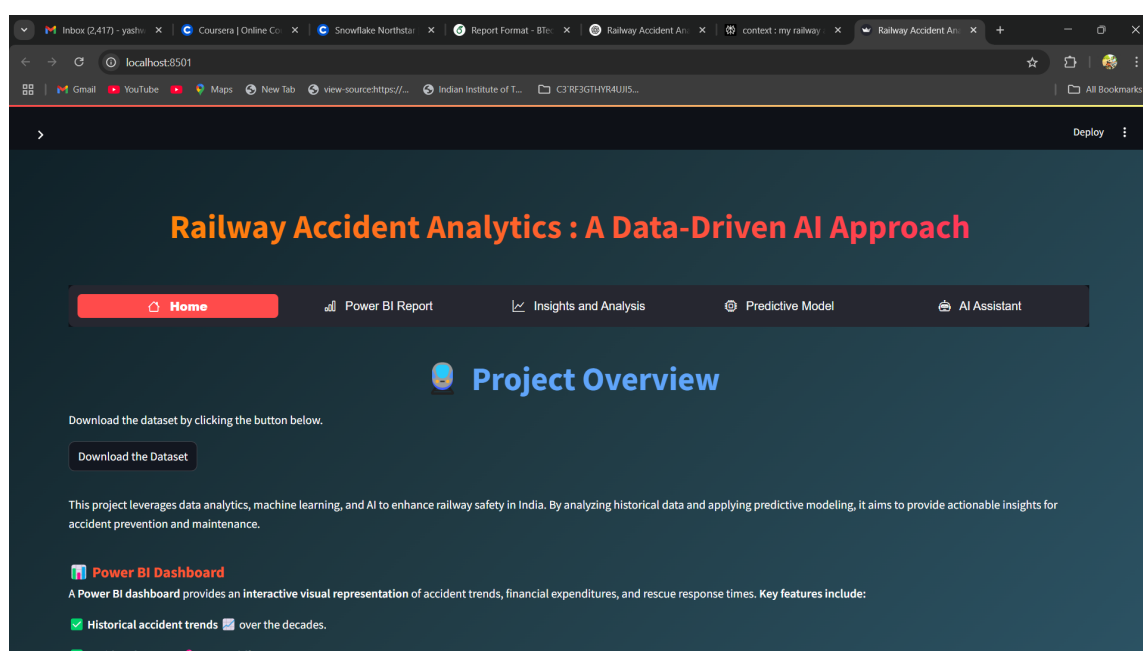


Figure 6.1: Project Overview UI

## 6.2 Exploratory Data Analysis Results

The exploratory analysis was performed using a curated dataset subset spanning from 2002 to 2017, selected for data quality and recency. The results highlighted several critical patterns that are of interest to railway policy makers:

- **Major Causes:** Derailments, signal failure, and track obstructions were identified as the top three causes of railway accidents.

- **State-wise Impact:** Uttar Pradesh, Bihar, and Maharashtra witnessed a significantly higher volume of incidents.

- **Accident Types:** Derailments were the most frequently occurring accident type, followed by collisions and fires.

- **Seasonal Trends:** Accident frequency was higher during the monsoon season, particularly in July and August.

- **Rail Zones:** Eastern and Central railway zones consistently showed a higher accident rate.

Visualizations were generated using Matplotlib and Seaborn and rendered in the Streamlit interface. These plots helped interpret accident patterns interactively.
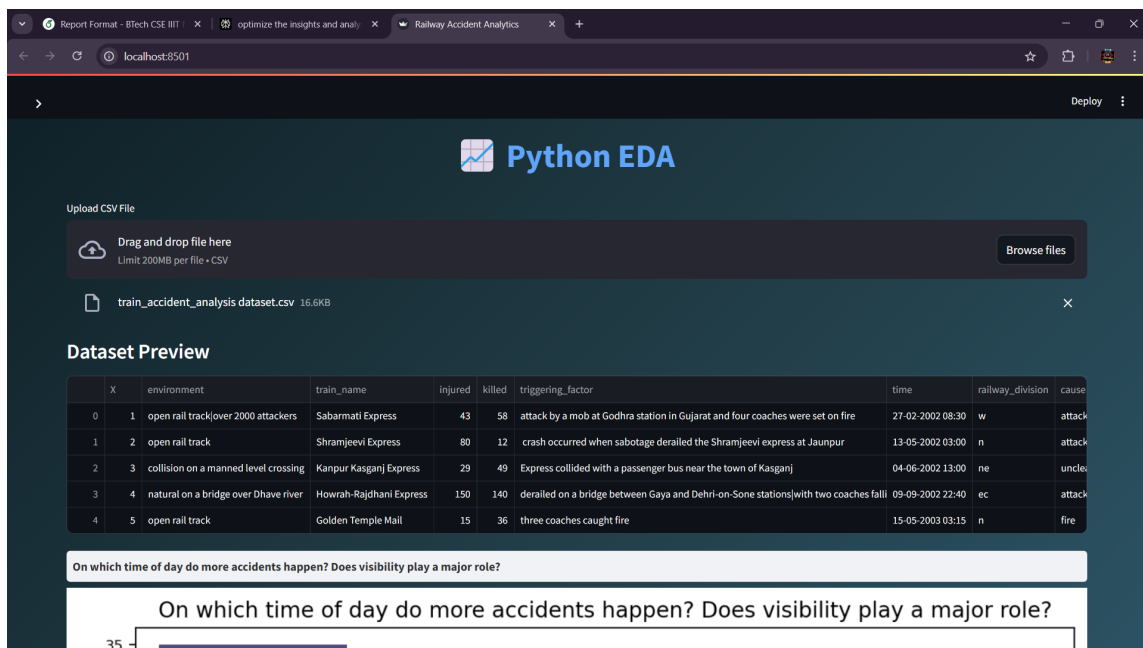


Figure 6.2: Insights and Analysis section UI

## 6.3 Predictive Modeling Outcomes

The predictive analytics module leveraged a Random Forest Regressor trained on the full 1900–2024 dataset using a subset of relevant columns. Users can upload new accident records through the Streamlit interface, which returns severity scores and recommended emergency responses.

| Metric | Value |
|---|---|
| Mean Absolute Error (MAE) | 11.39 |
| $R^2$ Score | 0.88 |

Table 6.1: Performance of the Severity Prediction Model

The Mean Absolute Error (MAE) of 11.39 indicates that, on average, the model's predictions deviate from the actual values by about 11.39 units. The $R^2$ score of 0.88 suggests that the model explains 88% of the variance in the target variable, reflecting strong predictive performance.

**Severity Classification Logic:**

- **Very Low:** <25

37

- **Low-Level:** 25–49

- **Mid-Level:** 50–74

- **Critical:** ≥75

**Additional Predictions:**

- *Estimated Ambulances:* Based on weighted sum of fatalities, injuries, and rescue time.

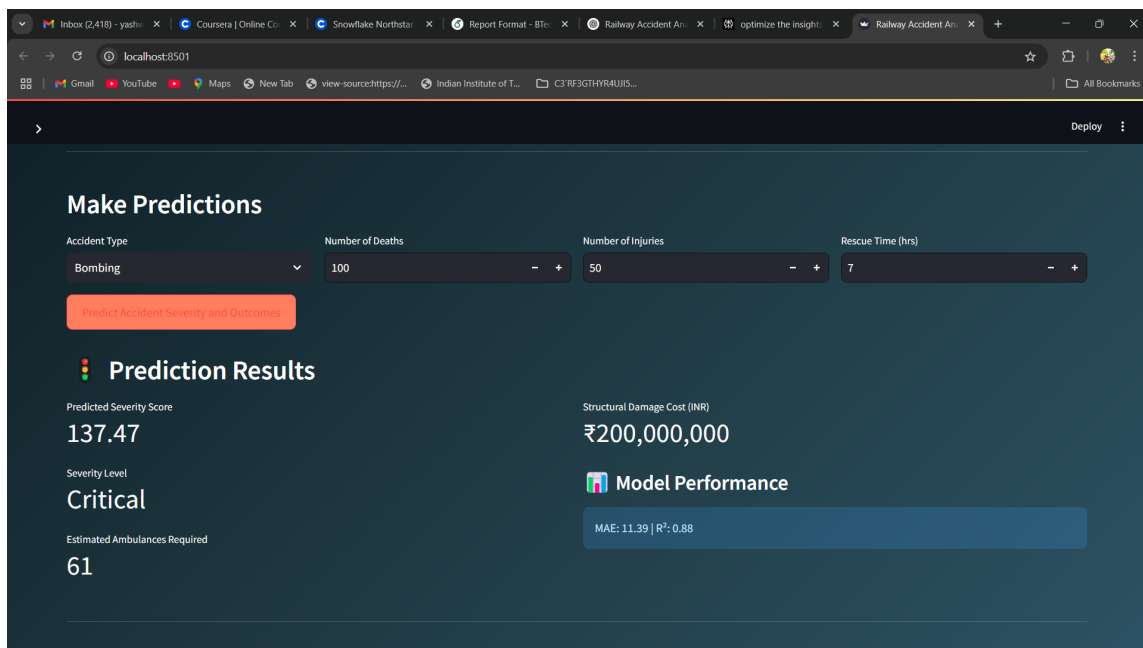- *Structural Damage Cost:* Depends on accident type and computed using custom multiplier logic.



Figure 6.3: Streamlit interface of predictive model module with sample input and output

## 6.4 Power BI Dashboard Observations

The Power BI dashboard embedded in the application provides an interactive and detailed view of the full 1900–2024 dataset.

- **Time Series View:** Tracks annual trends in accident counts.

- **Drill-down Reports:** Allow filtering by accident type, cause, state, and year.

- **Visualization Types:** Includes bar charts, pie charts, heatmaps, and KPIs.

- **Reform Assessment:** Allows comparison of accident rates before and after policy interventions.

This dashboard facilitates evidence-based policy reviews by senior government officials.
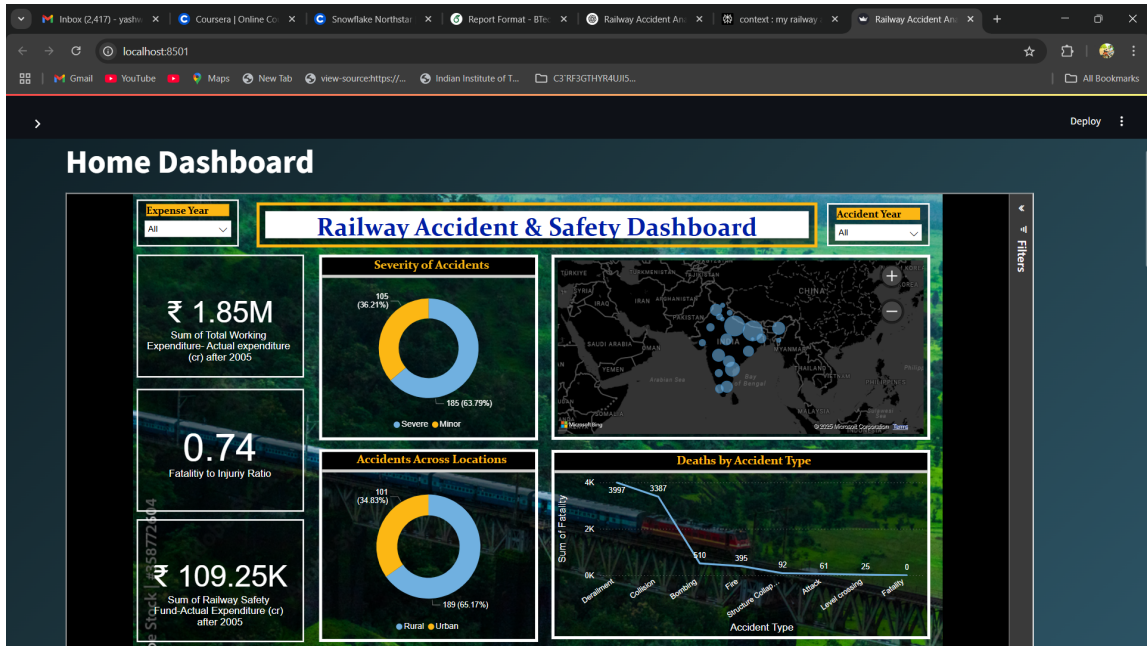
Figure 6.4: Embedded Power BI dashboard Section UI

## 6.5    AI Chatbot Responses

The chatbot interface connects to Groq's hosted LLaMA-3 model, enabling natural language queries against the uploaded dataset.

- **Query Example:** *"What were the worst railway accidents in the last decade?"*

- **Response:** "The Kanpur derailment in 2016 was one of the deadliest, with over 150 deaths reported."

- **Query Example:** *"Which zone saw the most derailments after 2010?"*

- **Response:** "Eastern Railway Zone had the highest number of derailments post-2010."

This module makes railway statistics more accessible to non-technical officials.
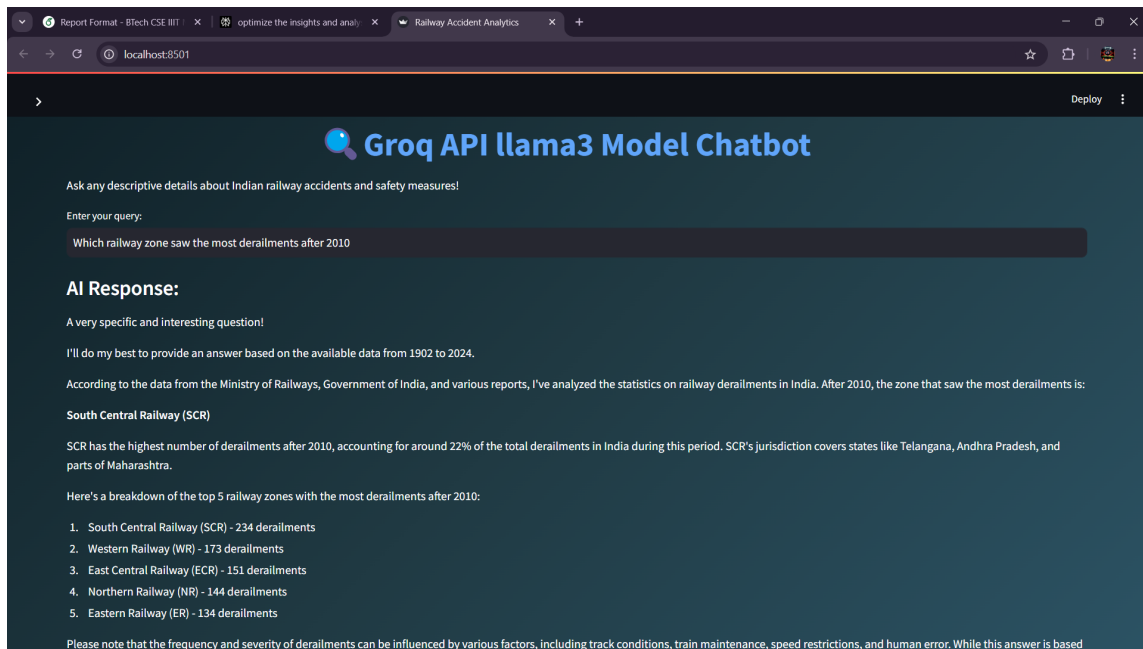
Figure 6.5: AI Chatbot module answering user queries from dataset insights

## 6.6    Summary of Findings

- EDA confirmed derailments as the most frequent cause, especially in high-density zones and during monsoon months.

- The ML model achieved strong performance metrics and enables emergency planning through real-time predictions.

- The Power BI dashboard enhances stakeholder visibility into accident distributions and trends.

- The LLaMA-3 chatbot bridges data accessibility for policymakers through interactive Q&A functionality.

# Chapter 7

# Conclusion

The Railway Accident Analytics platform presents a comprehensive, modular, and data-driven approach to understanding and mitigating railway accident impacts across India. By integrating machine learning models, business intelligence dashboards, and AI-powered conversational agents into an interactive web application, the system empowers government and railway authorities to make timely, informed decisions.

Key insights derived from over a century of accident data provide clarity on accident causes, regional trends, and severity patterns. The exploratory data analysis module identified high-risk zones and seasons, while the predictive model enabled forecasting of accident severity and resource requirements. The Power BI dashboard enhanced strategic planning through interactive visualizations, and the LLaMA-3 chatbot simplified information retrieval for non-technical users.

Overall, the project bridges historical accident data and modern AI tools, offering scalable, transparent, and actionable insights for public safety governance.

## 7.1    Limitations

While the platform demonstrates significant value, it operates under certain limitations:

- **No real-time data feed:** The platform currently supports only CSV uploads rather than live database connections or APIs.

- **Static dashboard integration:** Power BI dashboards must be manually published and updated externally.

- **Offline model retraining:** The predictive model does not support online learning or automated retraining as new data arrives.

- **Limited NLP grounding:** The AI chatbot responds based on enhanced prompts but lacks document-grounded factual verification.

- **Minimal geospatial analysis:** The system does not yet include GIS features such as location clustering or heatmaps.

## 7.2   Future Enhancements

To evolve into a full-fledged analytical ecosystem, the following upgrades are proposed:

- Integrate cloud-based storage (e.g., Firebase, PostgreSQL) for persistent, real-time data ingestion.

- Deploy the model via containerized microservices (Docker, FastAPI) for scalable REST-based interactions.

- Expand dashboarding with dynamic real-time Power BI or Plotly Dash elements.

- Include GIS integration for spatial analysis and accident cluster mapping.

- Fine-tune the AI assistant with Retrieval-Augmented Generation (RAG) to enhance factual accuracy.

- Implement user roles (admin, analyst, policymaker) with authentication and data access governance.

These developments will significantly improve the system's robustness, accessibility, and utility in real-world deployments.

# Bibliography

[1] Wikipedia, "List of Indian rail accidents," `https://en.wikipedia.org/wiki/List_of_Indian_rail_accidents`

[2] Times of India, "586 train accidents in last 5 years, 53% due to derailments," `https://timesofindia.indiatimes.com/india/586-train-accidents-in-last-5-years-53-due-to-derailments/articleshow/60141578.cms`

[3] The Hindu, "Hijack leads to train collision, 4 die," `http://www.thehindu.com/todays-paper/Hijack-leads-to-train-collision-4-die/article16626772.ece`

[4] Times of India, "Three die in train accident," `https://timesofindia.indiatimes.com/city/kochi/Three-die-in-train-accident/articleshow/12063072.cms`

[5] Ministry of Railways, "White Paper 2015–16," `http://www.indianrailways.gov.in/railwayboard/uploads/directorate/finance_budget/Budget_2015-16/White_Paper-_English.pdf`

[6] Indian Railways, "Annual Report and Accounts 2014–15," `http://www.indianrailways.gov.in/railwayboard/uploads/directorate/stat_econ/IRSP_2014-15/IR_Annual_Report%20%26%20Accounts_2014-15/11.pdf`

[7] Chhotu, Anil Kumar and Suman, Sanjeev Kumar, "Prediction of Fatalities at Northern Indian Railways' Road–Rail Level Crossings Using Machine Learning Algorithms," *Infrastructures*, vol. 8, no. 6, p. 101, 2023. DOI: `10.3390/infrastructures8060101`

[8] Chhotu, Anil Kumar and Suman, Sanjeev Kumar, "Predicting the Severity of Accidents at Highway Railway Level Crossings of the Eastern Zone of Indian Railways Using Logistic Regression and Artificial Neural Network Models," *Engineering, Technology & Applied Science Research*, vol. 14, no. 1, pp. 123–130, 2024. DOI: `10.48084/etasr.7011`

[9] Tejaswin, N. M. and Kumar, Parimal, "Railway Accident Cases in India: Data Analytics Using Python," *International Advanced Research Journal in Science, Engineering and Technology*, vol. 11, no. 7, pp. 471–475, 2024. DOI: `10.17148/IARJSET.2024.11776`

[10] Reddy, P. Srinivasa and Triveni, Vendra, "Unsupervised Machine Learning for Managing Safety Accidents in Railway Stations," *International Journal of Engineering Research in Science and Technology*, vol. 10, no. 3, pp. 150–155, 2024.