**PROJECT PROPOSAL**

**BLOCK 1D-2023-2024**

**Leveraging CNN and MLP for Predictive Accident Prevention in Breda**

**Alexi Kehayias, Musaed Al-Farah, Michon Goddijn & Benjamin Spiertz**

**Team 24**

**24th of May 2024**

# Introduction Contributor: Alexi Kehayias

## Overview

The project leverages deep learning to analyze driving behavior and road condition data.

It aims to reduce high-risk areas for accidents in Breda by identifying high-risk driving patterns and correlating them with road conditions.

## Relevance and Interest

Road safety impacts public health, economic costs, and quality of life.

Predicting and preventing traffic incidents through data-driven insights is crucial.

This project aligns with smart city initiatives and modern trends in data analytics

## Significance of Deep Learning

Deep learning offers tools to analyse large datasets and uncover patterns.

In this project, models predict high-risk areas for accidents based on driving speed, road conditions, and historical data.

This capability enables proactive measures to enhance safety, leading to safer urban environments

# Purpose and Objectives

1. **Develop CNN and MLP Models to Predict Accident Severity**
   - **Convolutional Neural Network (CNN)-** Construct a CNN model to predict the severity of accidents based on data from ANWB Safe Driving, KNMI, Argaleo Greenery, and Breda Road datasets. The model will utilize features such as latitude, longitude, speed, and road conditions to predict accident severity categories.
   - **Multi-Layer Perceptron (MLP)-** Build an MLP model to classify accident severity using the same datasets. The model will incorporate features like driving speed, weather conditions, and temporal data to enhance its predictive capability.
2. **Identify High-Risk Accident Locations and Times**
   - Use the CNN and MLP models to analyse and identify high-risk areas and periods by examining spatial data (latitude, longitude) and temporal data (event start and end times). This analysis will help pinpoint where and when accidents are most likely to occur.
   - Employ geographic and temporal data analysis to detect patterns and trends in accident occurrences, enabling targeted interventions.
3. **Correlate Road and Environmental Conditions with Accident Data**
   - Integrate and analyse data on road and environmental conditions, such as precipitation, temperature, and greenery, with accident records to understand their impact on accident severity.
   - Use the models to explore the relationships between different conditions and accident data, providing insights into how these factors influence accident risk.
4. **Provide Insights for Traffic Management and Insurance Optimization**
   - Utilize findings from the CNN and MLP models to generate actionable insights for local authorities to improve traffic management and enhance road safety.
   - Offer recommendations to insurance companies for optimizing premium structures based on accurate risk assessments derived from the models' predictions.
   - Provide comprehensive insights to guide policy decisions and implement safety measures aimed at reducing accident rates and improving overall traffic safety.

# Significance of Deep Learning in our project

Deep learning plays a crucial role in our project by processing vast amounts of data to accurately predict accident severity and identify high-risk areas. The significance of incorporating deep learning models like CNN's and MLP's in our approach includes:

1. **Pattern Recognition and Prediction**
   a. Deep learning models excel at identifying complex patterns in large datasets, allowing for more accurate predictions of accident severity and risk areas.
   b. By analysing diverse features such as driving behaviour, environmental conditions, and spatial data, these models provide robust predictions that are essential for proactive traffic management and safety interventions.

2. **Enhanced Safety Measures and Insights**
   a. The models help pinpoint high-risk locations and times, enabling targeted interventions to prevent accidents.
   b. Insights derived from the model predictions can inform local authorities and policymakers in making data-driven decisions to enhance road safety.

3. **Dynamic and Real-Time Applications**
   a. Deep learning models offer the potential for real-time feedback and dynamic adjustments to traffic management and insurance pricing.
   b. This capability supports infrastructure improvements and adaptive strategies to mitigate accident risks as conditions change.

4. **Optimized Insurance Premiums**
   a. The ability to predict and quantify accident risks allows insurance companies to tailor premiums more accurately, reflecting the true risk levels associated with different driving behaviours and conditions.
   b. This dynamic approach to pricing can lead to fairer and more competitive insurance premiums for customers, based on real-time risk assessments.

In summary, the application of deep learning models in our project provides significant advantages in predictive accuracy, actionable insights, and adaptive strategies for improving road safety and optimizing insurance practices.

# Problem Statement Contributors: Musaed Al-Fareh & Alexi Kehayias

## Problem Definition

Traffic accidents in high-risk areas of Breda pose significant safety and economic challenges.

Despite existing measures, there remains a pressing need for more effective strategies to predict and mitigate these incidents.

Our project aims to leverage deep learning models to predict high-risk accident zones, providing actionable insights that can lead to enhanced road safety and optimized traffic management.

## Context and Background

Breda, a vibrant city known for its bustling traffic and dynamic road networks, faces ongoing challenges related to traffic congestion, road safety, and infrastructure maintenance.

High-risk driving behaviours, such as speeding, sudden braking, and distracted driving, often result in traffic accidents, particularly in certain hotspots.

Traditional methods of traffic management and accident prevention have proven insufficient in addressing the complexities of these issues.

By integrating diverse datasets, such as the ANWB Safe Driving data, KNMI weather data, Argaleo Greenery data, and Breda Road data, our approach aims to uncover patterns and correlations in these high-risk areas that conventional analysis might overlook.

- **Driving Behaviour-** Evaluating the role of driving speed, sudden stops, and other behaviours in contributing to accidents.

# Importance and Benefits

Addressing the problem of traffic accidents in high-risk areas is crucial for several reasons:

1. ## Safety
   - **Life-Saving Potential-** Reducing traffic accidents directly translates to saving lives and minimizing injuries. Predictive insights allow for proactive measures in high-risk zones, potentially preventing accidents before they occur.
   - **Enhanced Emergency Response-** Accurate predictions of high-risk areas can streamline emergency response efforts, ensuring faster and more efficient assistance when accidents happen.

2. ## Economic Impact
   - **Cost Reduction-** Fewer accidents lead to lower medical expenses and vehicle repair costs. Effective management of high-risk areas can also reduce the financial burden on public services and insurance systems.
   - **Efficient Traffic Management-** Improved traffic flow and reduced congestion in high-risk areas help lower fuel consumption and emissions. This not only saves money for drivers but also contributes to environmental sustainability.
   - **Infrastructure Planning-** Insights from predictive models can inform better infrastructure investments and maintenance, optimizing the allocation of resources for road improvements.

3. ## Quality of Life
   - **Stress Reduction-** Less time spent in traffic and fewer accidents lead to a more pleasant driving experience, reducing stress for drivers and commuters.
   - **Community Well-Being-** Safer roads and better traffic conditions enhance the overall quality of life for residents, making communities more liveable and attractive.
   - **Public Trust-** Demonstrating effective use of technology to improve road safety can increase public trust in local authorities and their ability to manage urban challenges.

By focusing on high-risk areas, we can target interventions more effectively, maximizing the impact of safety measures and resources.

Our deep learning models provide a data-driven foundation for these interventions, offering precise, actionable insights that traditional methods might miss.

# AI Business Model Canvas Overview <span style="font-size:smaller">Contributor:</span>

## Musaed Alfareh

Our project leverages advanced deep learning techniques to analyse a comprehensive range of data from the ANWB Safe Driving, Breda Road, KNMI temperature, and Argaleo Greenery datasets. The aim is to revolutionize traffic management and enhance road safety in Breda by offering predictive insights and actionable intelligence.

## Project Goals

1. **Predict High-Risk Accident Areas**
   - **Objective-** Utilize a variety of incident data, including harsh braking, sharp cornering, and speeding, to forecast zones with a higher probability of accidents.
   - **Approach-** By employing sophisticated neural network models, we can detect patterns and trends in driving behaviours that correlate with increased accident risks.
2. **Identify High-Risk Locations and Times**
   - **Objective-** Conduct detailed geographic and temporal analysis to identify specific areas and time frames where the likelihood of accidents is elevated.
   - **Approach-** Integrate spatial and time-series data to pinpoint critical hotspots and periods, allowing for proactive measures to enhance safety.
3. **Inform Traffic Management and Insurance**
   - **Objective-** Provide insightful data to local authorities to improve traffic control and help insurance companies adjust premiums dynamically based on real-time risk assessments.
   - **Approach-** Use predictive models to generate actionable insights that can inform policy decisions, optimize resource allocation, and tailor insurance premiums to reflect actual driving risks.

## Overall Project Vision

Our project aims to seamlessly integrate predictive analytics into Breda's traffic management framework, enabling both immediate and strategic responses to potential risks. By harnessing the power of deep learning:

- **Real-Time Interventions-** Authorities can implement live, data-driven interventions to mitigate risks as they develop, enhancing road safety on a moment-to-moment basis.
- **Strategic Planning-** Long-term planning and infrastructure development can be informed by accurate predictions and insights into high-risk areas and behaviours, fostering a safer and more efficient traffic ecosystem.

## Implementing LSTM in Accident Prevention in Breda

### Data

1.ANWB Safe Driving Dataset
Use: To analyze driving behavior and identify correlations with high-risk incidents.

2. Breda Road Dataset
Use: To correlate road conditions and functionalities with incident data, providing contextual understanding essential for predictive modelling.
3. KNMI Precipitation & Temperature
Use: Analyze weather impact on driving incidents.
4. Argaleo Greenery
Use: Assess greenery's effect on driving conditions.

### Skills

1.Data engineering skills for building pipelines
2.Machine learning skills to build and train models
3.Data visualization skills to visualize data
4.Product Managers and UX Designers To ensure the solution is user-friendly and meets the needs of city officials.

### Output

1.Incident Reduction: Measurable decrease in traffic incidents based on predictive interventions.
2. Predictive Model Accuracy: High accuracy in predicting high-risk incidents

### Value Propositions

1.Enhanced Road Safety- Reducing high-risk incidents through predictive insights.
2.Dynamic Insurance Pricing- Utilising driving data for more accurate insurance premiums.
3.Data-Driven Decision Making- Enabling the local government and ANWB to make informed decisions on traffic management and safety improvements.

### Integration

1. System Compatibility.
2. Real-Time Data Processing.
3. Stakeholder Systems.

### Stakeholders

1.ANWB.
2. Local Government.
3. Drivers and Residents of Breda.

### Customer Segments

1. Local Road Authorities(Primary)- Direct users of traffic management insights.
2.Insurance Providers(Secondary)- Beneficiaries of risk assessment models for premium adjustments.
3.General Public and ANWB Members(Tertiary)- Recipients of safer driving conditions and possibly lower insurance costs.

### Cost Structure

1.Data Management- Costs for data acquisition, storage, and maintenance.
2.Development and Operations- Expenses related to technology development, system integration, and ongoing operations.
3.Outreach and Engagement- Budget for marketing, stakeholder meetings, and public engagement activities.

### Revenue Stream

1.Direct Revenue- If the solution is sold to other cities or municipalities.
2.Indirect Revenue- Cost savings from reduced accidents and insurance claims, increased efficiency in emergency responses, and improved public safety.
3.Supporting Innovation- Potential development of new safety-related products and services in collaboration

# Data Description <span>Contributor: Alexi Kehayias</span>

## Approach

Our project adopts a comprehensive data-driven approach utilizing deep learning techniques to analyse various datasets.

By examining driving behaviour and road conditions, we aim to predict high-risk accident areas and deliver actionable insights.

These insights will enable targeted interventions, enhancing road safety and optimizing traffic management in Breda.

## Steps Involved

1. **Data Collection**
   - **Source Integration-** Compile and integrate data from multiple sources, including the ANWB Safe Driving, Breda Road, KNMI Precipitation, KNMI Temperature, and Argaleo Greenery datasets.
   - **Data Quality Assurance-** Ensure the integrity and reliability of the data by systematically addressing issues such as missing values and duplicates to maintain a high-quality dataset for analysis.
2. **Data Pre-Processing**
   - **Data Cleaning-** Implement robust data cleaning processes to handle missing values and remove any duplicates. This step is crucial to prepare the data for effective analysis.
   - **Standardization-** Standardize date and time formats to ensure consistency across datasets, facilitating accurate temporal analysis.
   - **Feature Engineering-** Create and refine features, particularly for geospatial data and temporal attributes, to enhance the model's ability to capture relevant patterns and insights.
3. **Model Selection**
   - **Algorithm Choice-** Employ advanced deep learning models, including Convolutional Neural Networks (CNN's) and Multi-Layer Perceptron (MLPs), to assess and predict accident risk based on driving behaviour and environmental conditions.
   - **Customization-** Tailor the models to effectively process and learn from the diverse features of the dataset, optimizing them for the specific requirements of accident prediction.
4. **Training and Evaluation**
   - **Data Splitting-** Divide the cleaned and pre-processed data into training and testing subsets to ensure robust evaluation of model performance.

- o **Model Training-** Train the selected deep learning models using the training subset, fine-tuning parameters to achieve optimal performance.
- o **Performance Metrics-** Evaluate the models using a range of metrics, including accuracy, precision, recall, and F1-score, to ensure comprehensive assessment of their predictive capabilities.
5. **Deployment**
   - o **System Integration-** Integrate the trained predictive models into Breda's traffic management system, enabling real-time monitoring and analysis.
   - o **Actionable Insights-** Provide real-time alerts and insights to local authorities and insurance companies, facilitating immediate and strategic responses to high-risk driving behaviours and conditions.
   - o **User Interface-** Develop user-friendly interfaces and dashboards to visualize predictions and insights, aiding decision-makers in implementing effective traffic safety measures.

# Datasets Overview- Features & Their Significance

**Contributor: Alexi Kehayias**

In our project, we utilize a diverse range of datasets to capture various aspects of driving behaviour, road conditions, and environmental factors. Here's an overview of the datasets used, their key features, and their significance in predicting high-risk accident areas in Breda.

## 1. ANWB Safe Driving Dataset

- **Key Features**
  - o event_start- The start time of the driving event.
  - o event_end- The end time of the driving event.
  - o latitude- The geographic latitude of the event.
  - o longitude-The geographic longitude of the event.
  - o speed_kmh- The speed of the vehicle in kilometers per hour.
  - o Category- Type of driving event (e.g., harsh braking, cornering).
  - o incident_severity- Severity level of the incident (e.g., low, minor, moderate, severe).
- **Reasoning-** This dataset provides crucial insights into driving behavior and incidents. By analyzing events like harsh braking or cornering, and their severity levels, we can identify patterns that correlate with high-risk zones. Features such as speed, location, and event type are essential for understanding and predicting accident patterns.

## 2. Breda Road Dataset

- **Key Features**
  - objectbegintijd- Start time of the road segment or condition.
  - objecteindtijd- End time of the road segment or condition.
  - bgt_functie- Function or usage type of the road (e.g., residential, commercial).
  - wegdeeloptalud- Road segment slope or incline.
- **Reasoning-** This dataset provides context about road conditions and functionalities, which are crucial for understanding how different road types and segments contribute to accident risks. Features like road function and slope help in identifying and contextualizing accident-prone areas.

## 3. KNMI Precipitation Dataset

- **Key Features**
  - dtg- Date and time of the weather observation.
  - latitude- Geographic latitude of the observation point.
  - longitude- Geographic longitude of the observation point.
  - dr_pws_10- Wind speed (10-minute average).
  - ri_pws_10- Wind direction (10-minute average).
- **Reasoning-** Weather conditions significantly impact driving safety. This dataset allows us to assess how wind speed and direction influence driving behavior and incident frequency. Understanding these weather variables helps in predicting how adverse weather can increase accident risks.

## 4. KNMI Temperature Dataset

- **Key Features**
  - dtg- Date and time of the weather observation.
  - latitude- Geographic latitude of the observation point.
  - longitude- Geographic longitude of the observation point.
  - t_dryb_10- Dry bulb temperature (10-minute average).
  - t_dewp_10- Dew point temperature (10-minute average).
- **Reasoning-** Temperature changes affect road conditions and vehicle performance. Analyzing features like dry bulb temperature and dew point helps in understanding how temperature variations contribute to accident risks, particularly in severe weather conditions.

## 5. Argaleo Greenery Dataset

- **Key Features**
  - objectbegintijd- Start time of the greenery observation.
  - objecteindtijd- End time of the greenery observation.
  - bgt_status- Status of the greenery (e.g., maintained, overgrown).
  - bgt_type- Type of greenery (e.g., trees, bushes).
- **Reasoning-** Greenery can influence driving conditions by affecting visibility, road surface conditions, and driver behavior. This dataset helps in incorporating environmental aspects into our analysis. By examining the status and type of greenery, we can determine how these factors correlate with accident likelihood and severity.

## Summary

By integrating these datasets, our approach encompasses a comprehensive view of the factors influencing driving safety in Breda.

From driving behaviors captured by the ANWB Safe Driving data to the impact of road conditions, weather, and environmental features, each dataset contributes unique insights essential for predicting high-risk accident areas.

This multifaceted analysis allows us to develop more accurate and actionable predictions, ultimately enhancing road safety and traffic management.

# Methodology

**Contributors: Musaed Al-Fareh, Alexi Kehayias, & Michon Goddijn**

Our approach utilizes a robust, data-driven methodology with deep learning to analyse driving behavior and road conditions, aiming to predict high-risk accident areas. This methodology provides actionable insights for targeted interventions, ultimately improving road safety and efficiency in traffic management.

## Approach

We leverage comprehensive datasets and advanced machine learning models, specifically Convolutional Neural Networks (CNN) and Multilayer Perceptron (MLP), to assess and predict the risk levels associated with different driving behaviors and conditions. Our steps include:

1. Data collection.

2. Pre-processing.

3. Model selection.

4. Training.

5. Evaluation.

6. Deployment.

## Steps Involved

### 1. Data Collection

- **Data Sources**
    - **ANWB Safe Driving-** Provides data on driving behavior including event start/end times, location, speed, and incident severity.
    - **Breda Road-** Offers information on road conditions and functionalities.
    - **KNMI-** Includes weather data such as precipitation, temperature, and wind speed/direction.
    - **Argaleo Greenery-** Supplies data on the status and type of greenery affecting road conditions.

- **Data Integrity**
  - Ensure data quality by correcting missing values, removing duplicates, and verifying data consistency across all datasets.

## 2. Pre-Processing

- **Data Cleaning**
  - Handle missing values by imputing or removing them as necessary.
  - Remove duplicates to maintain the accuracy and quality of the dataset.
  - Standardize date and time formats to ensure uniformity.
- **Feature Engineering**
  - Extract meaningful features from raw data, such as time_of_day, distance_to_landmark, and event_duration.
  - Incorporate spatial features and time-based attributes to enhance predictive capabilities.
- **Data Normalization**
  - Scale and normalize data to prepare it for efficient model training, ensuring that all features contribute equally to the model's learning process.

## 3. Model Selection

- **Convolutional Neural Networks (CNN)**
  - **Rationale-** CNN's are adept at capturing spatial relationships and patterns in data, making them suitable for tasks involving geographic and temporal analysis.
  - **Architecture-** Multi-layered architecture with dense layers, dropout for regularization, and batch normalization to improve training stability.
- **Multilayer Perceptron (MLP)**
  - **Rationale-** MLPs are versatile and effective for structured data where spatial relationships are not as critical.
  - **Architecture-** Simpler compared to CNN's, with several dense layers and dropout to prevent overfitting.

## 4. Training and Evaluation

- **Input Data**
  - Combined features from the ANWB Safe Driving, Breda Road, KNMI, and Argaleo Greenery datasets are used to train the models.

- **Data Splitting**
  - **Training Set (80%)-** Used for training the model, allowing it to learn patterns and relationships in the data.
  - **Testing Set (20%)-** Reserved for evaluating the model's performance on unseen data to ensure it generalizes well.
- **Evaluation Metrics**
  - **Accuracy-** Measures the proportion of correct predictions out of all predictions made.
  - **Precision, Recall, and F1-Score-** Provide a detailed evaluation of the model's performance in classifying different risk levels.
- **Model Training**
  - **CNN and MLP Training-** Models are trained using extensive epochs and adjusted learning rates, with early stopping to prevent overfitting.
  - **Validation-** Ongoing validation during training helps monitor and improve model performance, ensuring the model doesn't overfit to the training data.
- **Learning Curves Analysis**
  - Evaluate learning curves for both models to understand their training dynamics and ability to generalize to new data.

## 5. Deployment

- **Integration into Breda's Traffic System**
  - The trained models are integrated into a real-time traffic management system, enabling dynamic monitoring and prediction of high-risk areas.
  - A user-friendly interface is developed using **Streamlit**, allowing local authorities to visualize and act upon real-time insights. This application provides an accessible way for stakeholders to monitor traffic conditions and predict potential accidents.
- **Real-Time Insights**
  - The deployment ensures that predictions and alerts are provided in real-time. Drivers are alerted if they are at high or low risk of being involved in an accident, allowing for immediate interventions and adjustments in driving behavior and traffic management.
  - Continuous model updates and retraining are planned to maintain accuracy as new data becomes available, ensuring the system adapts to evolving traffic patterns and conditions.

## Summary

Our comprehensive methodology encompasses the entire lifecycle of data-driven prediction, from initial data collection to final deployment of predictive models. By integrating CNN and MLP models with detailed pre-processing and robust training protocols, we aim to deliver high-precision, real-time insights for accident prevention and improved traffic safety in Breda.

# Deep Learning Model

## MLP (Multilayer Perceptron)

Multi-Layer Perceptrons are versatile neural networks that consist of multiple layers of neurons. They are capable of learning complex patterns through deep network structures, making them suitable for a variety of prediction tasks. MLPs are effective for structured data where relationships between features are less hierarchical compared to spatial data

## CNN (Convolutional Neural Network)

Convolutional Neural Networks are designed to capture intricate patterns and features in data, particularly useful for tasks involving spatial or hierarchical data structures. In this context, CNNs excel at identifying complex relationships between various driving features and predicting accident severity categories.

# Model Architecture

## CNN (Convolutional Neural Network) Architecture

### Model Description

- **Input Layer-** Dense layer with 512 neurons and ReLU activation function, regularized with L2 to prevent overfitting.
- **Hidden Layers**
    - Second Dense layer with 256 neurons, ReLU activation, and L2 regularization.
    - Third Dense layer with 128 neurons, ReLU activation, and L2 regularization.
    - Fourth Dense layer with 64 neurons, ReLU activation, L2 regularization, and Batch Normalization to stabilize the learning process.
    - Fifth Dense layer with 32 neurons, ReLU activation, and L2 regularization.
- **Dropout Layers-** Applied after each dense layer to prevent overfitting by randomly setting a fraction of input units to 0.
- **Output Layer-** Dense layer with 5 neurons and softmax activation to match the number of severity categories.
- **Tools and Techniques**
    - **Data Analysis-** pandas, NumPy
    - **Visualization-** matplotlib, seaborn
    - **Deep Learning-** TensorFlow (Sequential model with multiple Dense layers, Dropout, BatchNormalization, Adam optimizer)
    - **Database Management-** Postgres SQL for loading data.

## MLP (Multi-Layer Perceptron) Architecture

### Model Description

- **Input Layer-** Dense layer with 32 neurons and ReLU activation function.
- **Hidden Layers**
    - First Dense layer with 16 neurons and ReLU activation, regularized with L2.
    - Second Dense layer with 8 neurons and ReLU activation, regularized with L2.
- **Dropout Layers-** Applied after each dense layer to prevent overfitting by randomly setting a fraction of input units to 0.
- **Output Layer-** Dense layer with 5 neurons and softmax activation to match the number of severity categories.

- **Tools and Techniques**
    - **Data Analysis-** pandas, NumPy
    - **Visualization-** matplotlib, seaborn
    - **Deep Learning-** TensorFlow (Sequential model with multiple Dense layers, Dropout, Adam optimizer)
    - **Database Management-** Postgres SQL for loading data.

# Legal Obligations and Approach

## GDPR Compliance

To comply with GPDR laws we will need to maintain detailed records of data processing activities, obtain explicit consent from data subjects, and provide clear data handling policies.

# Transparency

To be transparent with stakeholders and to do this we would have to provide comprehensive documentation, offer regular updates on data usage and model development, and engage stakeholders through consultations and feedback sessions.

## Bias and Fairness

Ensure fairness and mitigate bias.

This would be done by conducting regular audits to identify and address biases, use fairness metrics to evaluate the impact on different demographic groups, and implement strategies to ensure equitable treatment across all segments.

## Ethical AI Use

Adhere to ethical standards.

This would be done by following ethical guidelines for AI development, engage ethics review boards to assess potential impacts, and incorporate ethical considerations into the AI system's design and implementation.

# Risk Assessment

# Model Accuracy

The risk associated here is incorrectly identifying high-risk accident areas.

This would have an impact of misallocation of resources, failure to prevent accidents.

A solution to this would be regular evaluation, feedback loops, and cross-validation.

## Bias and Fairness

A risk posed by this is that the AI system might favour certain demographic groups.

Having an impact on ethical issues, potential discrimination.

To mitigate this problem, we would carry out regular audits, use fairness metrics and implement equitable strategies.

## Operational & Performance Risk

The biggest risk when integrating and deploying the AI system is ensuring real-time accuracy and reliability in predicting high-risk accident areas.

This would lead to Ineffective interventions, disrupted operations.

One way to fix this would be to use reliable deployment strategies such as Pilot Testing, incremental rollout and continuous improvement.

## Overall Risk Level

Based on the AI Act's criteria and the detailed risk assessment, this project is classified as High-Risk due to its significant impact on public safety, potential biases affecting fairness, and operational challenges.

The AI system's primary objective is to predict high-risk accident areas, which directly impacts road safety by preventing accidents and reducing injuries.

However, incorrect predictions can lead to misallocated resources and ineffective interventions.

Additionally, inherent biases in historical data could lead to ethical issues and discrimination, while ensuring real-time accuracy and reliability presents substantial operational challenges.

To mitigate these risks, enhanced compliance measures including strict GDPR adherence and following AI Act guidelines are essential. Transparency and accountability can be maintained through stakeholder engagement and regular audits.

Ethical AI practices, such as bias mitigation and ethics reviews, along with comprehensive risk management strategies like continuous monitoring and phased deployment, are crucial for ensuring the AI system's effectiveness and fairness.

# Project Timeline

## Week 1: Data Preparation And Setup (May 28th to May 31st)
## Tasks

## Alexi

## Week 1

## Task 5.1: Database Connection Setup(0.5 Hours)

- Write and test the code to connect to the SQL database.

## Task 5.2: Initial Data Queries(2 Hours)

- Implement simple queries to access the data, load it into dataframes and merge it.

## Task 5.3 & 5.4: Data Preprocessing Initiation (6 Hours) & Advanced Data Processing (4 Hours)

- Execute SQL joins and perform advanced data preprocessing such as normalization and feature extraction.
- Perform data cleaning and preparation, including encoding and splitting the dataset.
- Inserting the joined, cleansed data into a new table called merged_data.

## Task 5.5: Data Balancing (4 Hours)

- Balance the class distribution and prepare a PDF detailing all preprocessing steps.

# Week 2

## Task 5.5: Creating Views In The Dataset (4 Hours)

- Create views for the  dataset to give us some insight and useability of the data

## Task 6.1: Model Design & Assignment Begins (6 Hours)

- Discuss, decide, and assign model development tasks to team members.

## Task 6.2: Initial Model Training & Evaluation (4 Hours)

- Train and evaluate the first models using f1-score, MSE, RMSE, MAPE, Accuracy, Precision and Recall.

## Task 6.3: Model Optimisation (5 Hours)

- Improve model performance through various adjustments and preprocessing enhancements

### Task 6.4: Model Comparison (3 Hours)

- Compare different machine learning models to assess their complexity and performance.

### Task 6.5: Advanced Model Training (5 Hours)

- Develop and train a complex neural network or other suitable deep learning model.

# Week 3

### Task 7.1: Code Documentation & Review  (2 Hours)

- Review and document all Python scripts, ensuring proper commenting and structure.

### Task 7.2 & 7.3: Virtual Environment Serup, Code Structuring & Commenting  (5 Hours)

- Apply linting, logging and docstrings to all project code.

### Task 7.4: Unit Testing (2 Hours)

- Develop and run unit tests to achieve at least 30% code coverage.

### Task 7.5: Deployment Planning  (3 Hours)

- Detail the deployment strategy using frameworks such as Streamlit, Flask, Django, or FastAPI.

### Task 7.5: Building Application  (5 Hours)

- Apply the deployment strategy and build the application, test it and publish it.

## Preparing Final Presentation & Final Ethics Write Up (6 Hours)

- Conduct a final rehearsal and make last-minute adjustments to the presentation.

### Musaed
### Week 1

## Exploratory Data Analysis (Combined 11 hours over three days)

- Explore the data overall, analyse patterns and anomalies the data may show.

## Research Into Types of Model To Use Based Off Of EDA ( 5 Hours)

- The research here allows us to see which models are the most suitable to work with, compare and select the best model of the chosen models.

# Week 2

## Task 6.1: Model Design & Assignment Begins (6 Hours)

- Discuss, decide, and assign model development tasks to team members.

## Task 6.2: Initial Model Training & Evaluation (4 Hours)

- Train and evaluate the first models using f1-score, MSE, RMSE, MAPE, Accuracy, Precision and Recall.

## Task 6.3: Model Optimisation (5 Hours)

- Improve model performance through various adjustments and preprocessing enhancements

## Task 6.4: Model Comparison (3 Hours)

- Compare different machine learning models to assess their complexity and performance.

## Task 6.5: Advanced Model Training (5 Hours)

- Develop and train a complex neural network or other suitable deep learning model.

# Week 3

## Task 7.1: Code Documentation & Review  (2 Hours)

- Review and document all Python scripts, ensuring proper commenting and structure.

## Task 7.2 & 7.3: Virtual Environment Serup, Code Structuring & Commenting  (5 Hours)

- Apply linting, logging and docstrings to all project code.

## Task 7.4: Unit Testing (2 Hours)

- Develop and run unit tests to achieve at least 30% code coverage.

## Task 7.5: Deployment Planning  (3 Hours)

- Detail the deployment strategy using frameworks such as Streamlit, Flask, Django, or FastAPI.

## Task 7.5: Building Application  (5 Hours)

- Apply the deployment strategy and build the application, test it and publish it.

## Preparing Final Presentation & Final Ethics Write Up (6 Hours)

- Conduct a final rehearsal and make last-minute adjustments to the presentation.

## Task 7.5: Finalising Application (5 Hours)

- Testing deployability, debugging any errors and planning on adding some final features to the application.

# Benjamin

## Week 1

## Task 3.4: Risk Analysis & Ethics Start (3 Hours)

- Document the level of risk associated with the developed AI system.

## Task 3.4: Risk Analysis & Ethics Completed (3 Hours)

- Document the level of risk associated with the developed AI system.

## Task 3.5: Legal & Ethical Requirements Mapping Start (3 Hours)

- Identify essential and non-essential legal requirements according to the AI Act.

## Task 3.5: Legal & Ethical Requirements Mapping Finalised (3 Hours)

- Identify essential and non-essential legal requirements according to the AI Act.

## Exploratory Data Analysis (4 Hours)

- Explore the data overall, analyse patterns and anomalies the data may show.

# Week 2

## Task 3.6: Legal Requirement Relevance (4 Hours)

- Explain the relevance and reasoning behind the inclusion or exclusion of specific legal requirements.

## Task 6.1: Model Design & Assignment Begins (6 Hours)

- Discuss, decide, and assign model development tasks to team members.

## Task 6.2: Initial Model Training & Evaluation (4 Hours)

- Train and evaluate the first models using f1-score, MSE, RMSE, MAPE, Accuracy, Precision and Recall.

## Task 6.3: Model Optimisation (5 Hours)

- Improve model performance through various adjustments and preprocessing enhancements

## Task 6.4: Model Comparison (3 Hours)

- Compare different machine learning models to assess their complexity and performance.

## Task 6.5: Advanced Model Training (5 Hours)

- Develop and train a complex neural network or other suitable deep learning model.

# Week 3

## Task 4.2: Final Presentation Preparation (4 Hours)

- Conduct a final rehearsal and make last-minute adjustments to the presentation.

## Task 3.7: Code Review for Ethical Compliance (4 Hours)

- Document the AI system's compliance with EU regulations.

### Task 7.4: Unit Testing (2 Hours)

- Develop and run unit tests to achieve at least 30% code coverage.

### Preparing Final Presentation & Final Ethics Write Up (Combined 12 hours over 3 days

- Apply final additions to the presentation, adjust or add any remarks and integrate the ethics into the presentation.

### Final Presentation Dry Run With Ethics Write Up (4 Hours)

- Dry run the presentation with the team, go over the whole project one final time before submission.

# Michon

## Week 1

## Task 3.4: Risk Analysis & Ethics Start (3 Hours)

- Document the level of risk associated with the developed AI system.

## Task 3.4: Risk Analysis & Ethics Completed (3 Hours)

- Document the level of risk associated with the developed AI system.

## Task 3.5: Legal & Ethical Requirements Mapping Start (3 Hours)

- Identify essential and non-essential legal requirements according to the AI Act.

## Task 3.5: Legal & Ethical Requirements Mapping Finalised (3 Hours)

- Identify essential and non-essential legal requirements according to the AI Act.

## Exploratory Data Analysis (4 Hours)

- Explore the data overall, analyse patterns and anomalies the data may show.

# Week 2

## Task 3.6: Legal Requirement Relevance (4 Hours)

- Explain the relevance and reasoning behind the inclusion or exclusion of specific legal requirements.

## Task 6.1: Model Design & Assignment Begins (6 Hours)

- Discuss, decide, and assign model development tasks to team members.

## Task 6.2: Initial Model Training & Evaluation (4 Hours)

- Train and evaluate the first models using f1-score, MSE, RMSE, MAPE, Accuracy, Precision and Recall.

## Task 6.3: Model Optimisation (5 Hours)

- Improve model performance through various adjustments and preprocessing enhancements

## Task 6.4: Model Comparison (3 Hours)

- Compare different machine learning models to assess their complexity and performance.

## Task 6.5: Advanced Model Training (5 Hours)

- Develop and train a complex neural network or other suitable deep learning model.

# Week 3

## Task 4.2: Final Presentation Preparation (4 Hours)

- Conduct a final rehearsal and make last-minute adjustments to the presentation.

## Task 3.7: Code Review for Ethical Compliance (4 Hours)

- Document the AI system's compliance with EU regulations.

## Task 7.4: Unit Testing (2 Hours)

- Develop and run unit tests to achieve at least 30% code coverage.

## Preparing Final Presentation & Final Ethics Write Up (Combined 12 hours over 3 days

- Apply final additions to the presentation, adjust or add any remarks and integrate the ethics into the presentation.

## Final Presentation Dry Run With Ethics Write Up (4 Hours)

- Dry run the presentation with the team, go over the whole project one final time before submission.

# Reference

Below you will find references to information we used in our proposal:

GDPR.eu. (n.d.). General Data Protection Regulation (GDPR). Retrieved from https://gdpr.eu/

European Commission. (n.d.). Data protection in the EU. Retrieved from https://ec.europa.eu/info/law/law-topic/data-protection/data-protection-eu_en

European Commission. (n.d.). Ethics guidelines for trustworthy AI. Retrieved from https://ec.europa.eu/futurium/en/ai-alliance-consultation

IEEE. (n.d.). Global initiative on ethics of autonomous and intelligent systems. Retrieved from https://ethicsinaction.ieee.org/

AI Fairness 360. (n.d.). IBM AI Fairness 360 open source toolkit. Retrieved from https://aif360.mybluemix.net/

Google AI. (n.d.). Responsible AI practices. Retrieved from https://ai.google/responsibilities/responsible-ai-practices/

National Institute of Standards and Technology (NIST). (n.d.). Cybersecurity framework. Retrieved from https://www.nist.gov/cyberframework

OWASP Foundation. (n.d.). OWASP top ten. Retrieved from
https://owasp.org/www-project-top-ten/

TensorFlow. (n.d.). TensorFlow documentation. Retrieved from
https://www.tensorflow.org/

Fowler, M. (n.d.). Continuous delivery. Retrieved from
https://martinfowler.com/bliki/ContinuousDelivery.html

Google Cloud. (n.d.). AI platform. Retrieved from
https://cloud.google.com/ai-platform