



# EMERALDS

Project Title	Extreme-scale Urban Mobility Data Analytics as a Service
Project Acronym	EMERALDS
Grant Agreement No.	101093051
Start Date of Project	2023-01-01
Duration of Project	36 months
Project Website	<a href="https://emeralds-horizon.eu/">https://emeralds-horizon.eu/</a>

## D5.1 – Use Cases Scoping Document

Work Package	<b>WP 5, Use Cases Experimentation &amp; Assessment Cycles</b>
Lead Author (Org)	<b>Sascha Hoogendoorn-Lanser (TUD)</b>
Contributing Author(s) (Org)	<b>Winnie Daamen (TUD), Kristine Eglite (G93), Foivos Galatoulas (INLE), Serge Hoogendoorn (TUD), Sascha Hoogendoorn-Lanser (TUD), Jurijs Kondratenko (G93), Erik-Sander Smits (ARA), Jeroen Steenbakkens (ARG), Yannis Theodoridis (UPRC)</b>
Due Date	<b>29.03.2024</b>
Date	<b>29.03.2024</b>
Version	<b>V2.0</b>

### Dissemination Level

- PU: Public
- SEN: Sensitive, only for members of the consortium (including the European Commission)



## Versioning and contribution history

Version	Date	Author	Notes &/or Reason
0.1	23/02/2023	Sascha Hoogendoorn-Lanser (TUD)	TOC and V0.1
0.2	18/03/2023	Jeroen Steenbakkens (ARG), Erik-Sander Smits (ARA), Serge Hoogendoorn-Lanser (TUD), Jurijs Kondratenko (G93)	Brief use case requirement description, Chapters 3,4,5
0.3	20/04/2023	ARG, ARA, TUD, G93	Updated Use Case Descriptions
0.4	14/05/2023	Winnie Daamen (TUD)	Chapter 2
0.5	20/05/2023	Foivos Galatoulas (INLE)	Chapter 1 Contributions
0.6	04/06/2023	Jeroen Steenbakkens (ARG)	Altered based on new proposed table of content
0.7	05/06/2023	Winnie Daamen (TUD)	Structure elements added based on the reviewed Table of Contents document
0.8	07/06/2023	Sascha Hoogendoorn-Lanser (TUD)	Detailing UC1, Chapter 6
0.9	27/06/2023	Jeroen Steenbakkens (ARG), Kristine Eglite (G93), Jurijs Kondratenko	Added datasets and mapping, UC3 Description
0.91	28/06/2023	Sascha Hoogendoorn-Lanser (TUD)	Make changes requested by reviewers
1.0	01/03/2024	All authors	Updated version based on progressive insights in use cases and their assessment
2.0	15/03/2024	All authors	Addressed Review Comments

## Quality Control (includes peer & quality reviewing)

Version	Date	Name (Organisation)	Role & Scope
0.4	25/06/2023	Foivos Galatoulas (INLE)	First Reviewer Comments
0.4	26/06/2023	Christos Doulkeridis (UPRC)	Second Reviewer Comments
0.9	28/06/2023	Foivos Galatoulas (INLE) & Christos Doulkeridis (UPRC)	Second Round of Internal Review
0.91	28/06/2023	Yannis Theodoridis (UPRC)	Approval by Scientific and Technical Manager
1.0	29/06/2023	Foivos Galatoulas (INLE)	Final review by Coordinator
2.0	23/03/2024	Yannis Theodoridis (UPRC)	Approval by Scientific and Technical Manager
2.0	24/03/2024	Foivos Galatoulas (INLE)	Final Review for resubmission



**Funded by  
the European Union**

*This project has received funding from the European Union's Horizon Europe research and innovation programme under Grant Agreement No 101093051*

#### **Disclaimer**

EMERALDS - This project has received funding from the Horizon Europe R&I programme under the GA No. 101093051. The information in this document reflects only the author's views and the European Community is not liable for any use that may be made of the information contained therein. The information in this document is provided "as is" without guarantee or warranty of any kind, express or implied, including but not limited to the fitness of the information for a particular purpose. The user thereof uses the information at his/ her sole risk and liability.

#### **Copyright message**

©EMERALDS Consortium. This deliverable contains original unpublished work except where clearly indicated otherwise. Acknowledgement of previously published material and of the work of others has been made through appropriate citation, quotation, or both. Reproduction is authorized provided the source is acknowledged.

## Table of Contents

1	Introduction.....	12
1.1	Purpose and Scope of the Document .....	12
1.2	Relation to Work Packages, Deliverables and Activities .....	13
1.3	Contribution to WP5 and Project Objectives.....	13
1.4	Structure of Use Case Scoping Document .....	14
2	Assessment Approach .....	15
2.1	Background on Assessment Methods.....	15
2.2	Focus of emeralds Assessment .....	16
2.3	General Assessment Approach .....	18
2.3.1	UC specific Software Performance Indicators.....	18
2.3.2	UC Business KPIs.....	18
3	UC1: Risk Assessment, Prediction & Forecasting during Events .....	19
3.1	Background – Context – Challenges .....	19
3.1.1	Need for Risk Identification, Prediction and Forecasting.....	20
3.1.2	The Hague UC Background and Developments so far .....	22
3.1.3	Extreme Scale Data Challenges and Need for Emeralds .....	26
3.2	SMART Definition and UC Objectives .....	26
3.2.1	User story 1: Planning for safe and efficient operations during beach days .....	30
3.2.2	User story 2: Crowd Management Planning of an event.....	30
3.2.3	User story 3: Real-time monitoring and decision-making during an event .....	30
3.3	Computational and Data Resources .....	30
3.3.1	Data Resources.....	30
3.4	Assessment and Corresponding KPIs .....	31
3.5	Mapping to EMERALDS Toolset .....	33
3.6	Implementation Plan .....	34
3.7	Key Stakeholders.....	34
4	UC2: Multi-modal integrated Traffic Management.....	35
4.1	Background – Context – Challenges .....	35
4.1.1	The value of data in multi-modal traffic management.....	35
4.1.2	Challenges in data processing and analytics posed from extreme scale .....	37
4.1.3	Rotterdam use case in the context of EMERALDS .....	38
4.2	SMART Definition and UC Objectives .....	40
4.2.1	User-story 1: using off-line data for design of multi-modal measures .....	47
4.2.2	User-story 2: Integrated network management (automated traffic control).....	48



- 4.2.3 User-Story 3: network operations..... 48
- 4.3 Computational and Data Resources ..... 49
  - 4.3.1 Data Resources..... 49
- 4.4 Assessment and corresponding KPIs ..... 50
- 4.5 Mapping to EMERALDS Toolset ..... 52
- 4.6 Implementation Plan ..... 54
- 4.7 Key Stakeholders..... 55
- 5 UC3: Public Transport Trip Characteristics Inference and Traffic Flow Data Analytics ..... 56
  - 5.1 Background – Context – Challenges ..... 56
  - 5.2 SMART Definition and Use Case Objectives..... 59
    - 5.2.1 User story 1: Optimisation of existing PT network ..... 64
    - 5.2.2 User story 2: Prioritising public transport in the mobility hierarchy ..... 64
    - 5.2.3 User story 3: Fleet planning and management ..... 65
  - 5.3 Computational and Data Resources ..... 65
    - 5.3.1 Data Resources..... 65
    - 5.3.2 Computational Requirements ..... 66
  - 5.4 Assessment and Corresponding KPIs ..... 67
  - 5.5 Mapping to EMERALDS Toolset ..... 70
  - 5.6 Implementation Plan ..... 73
  - 5.7 Key Stakeholders..... 74
- 6 Conclusions and Next Steps..... 75
- 7 Annex – Ethics Checklist ..... 77
- 8 References ..... 82

## List of Figures

FIGURE 2-1 – DATA VALUE CHAIN OF URBAN MOBILITY USE CASES. ....	17
FIGURE 3-1 MEDIA PUBLICATION REPORTING ON THE SCHEVENINGEN OVERCROWDING INCIDENT IN 2022.....	19
FIGURE 3-2 – EXAMPLE FROM YANG ET AL (2022) SHOWING LINEAR REGRESSION FIT OF THE TUD DATA AND THE RESIDENTIAL POPULATION OF EACH SUBDISTRICT IN THE AREA WITHIN THE OUTER-RING ROAD OF SHANGHAI, CHINA. ....	20
FIGURE 3-3 – DETERMINING CROWD RISK (ADAPTED FROM [20]).....	21
FIGURE 3-4 COMPARISON BETWEEN FORECASTING ACCURACIES FOR LINEAR REGRESSION AND XGBOOST.....	23
FIGURE 3-5 – HOURLY FORECASTS OF VISITS IN THE SCHEVENINGEN AREA. ....	24
FIGURE 3-6 – THE ARG DIGITWIN SHOWING THE CENTRAL BEACH AREA OF SCHEVENINGEN (ON THE LEFT) AND THE NUMBER OF VISITORS AND ITS PREDICTIONS (ON THE RIGHT). ....	24
FIGURE 3-7 – BOW TIE MODEL FOR THE CROWD MANAGEMENT SYSTEM .....	25
FIGURE 3-8 – PIPELINES FOR USE CASE 1.....	33
FIGURE 4-1 – EXPERIMENTAL FRAMEWORK FROM POELMAN ET AL (2020).....	36
FIGURE 4-2 – OFFLINE AND ONLINE TOOL TRAFFIC MANAGEMENT COMPONENTS .....	39
FIGURE 4-3 – NETWORK ANALYTICS OVERVIEW .....	40
FIGURE 4-4 – EXAMPLE OF BICYCLE COUNT DATA MAP.....	50
FIGURE 4-5 – PIPELINE 1, COVERING SCENARIOS 1-1, 1-2, AND 3-2. “E” INDICATES THE EMERALD CODES FOUND IN TABLE 5.....	52
FIGURE 4-6 – PIPELINE 2, COVERING SCENARIOS 2-3, AND 3-1. “E” INDICATES THE EMERALD CODES FOUND IN TABLE 5. ....	53
FIGURE 4-7 – PIPELINE 3, COVERING SCENARIOS 2-1, AND 2-2. “E” INDICATES THE EMERALD CODES FOUND IN TABLE 5. ....	53
FIGURE 5-1 – THE MAP OF PUBLIC TRANSPORT STOPS IN THE CITY OF RIGA [34].....	56
FIGURE 5-2 – VISUALISATION OF THE ORIGIN-DESTINATION ANALYSIS CONDUCTED SO FAR BY G93: PT TRIP FLOWS (LEFT), TRAJECTORY DATA WITH ENTRY AND EXIT SPOTS (RIGHT).....	57
FIGURE 5-3 – TYPICAL STEPS OF OD ESTIMATION PROBLEM [36].....	58
FIGURE 5-4: – DATA GATHERING AND INGESTION PROCESS AS A BLUEPRINT FOR THE UC SCENARIO DATA PIPELINES THAT WILL BE REPORTED IN D5.6 (M24), D5.7(M33).....	67
FIGURE 5-5 – UC3 PIPELINE.....	72
FIGURE 5-6 – UC3 MILESTONES DURING THE 1ST IMPLEMENTATION CYCLE. ....	73
FIGURE 5-7 – DRAFT OF UC3 MILESTONES DURING THE 2ND IMPLEMENTATION CYCLE. ....	74

## List of Tables

TABLE 1 – TERMINOLOGY.....	7
TABLE 2 – MATRIX OF ALIGNMENT .....	11
TABLE 3 – DESCRIPTION OF AVAILABLE DATA SOURCES AND THEIR CONNECTION TO THE SCENARIOS AND USE CASE OBJECTIVES.....	49
TABLE 4 – KPIS FOR THE DIFFERENT SCENARIOS. ....	51
TABLE 5 – EMERALDS AND THEIR CODES REFERRED TO IN THE PIPELINES.....	53
TABLE 6 – MAPPING OF THE PIPELINES TO THE DATA AND SPECIFIC DATA REQUIREMENTS. ....	54
TABLE 7 – UC3 (RIGA) DATA RESOURCES.....	65
TABLE 8 – MAPPING OF EMERALDS COMPONENTS TO UC3 OBJECTIVES .....	71

## Terminology

---

<b>Terminology/Acronym</b>	<b>Description</b>
<b>AI</b>	<b>Artificial Intelligence</b>
<b>DoA</b>	<b>Description of Action</b>
<b>EC</b>	<b>European Commission</b>
<b>GA</b>	<b>Grant Agreement</b>
<b>GDPR</b>	<b>General Data Protection Regulation</b>
<b>IEA</b>	<b>Independent Ethics Advisor</b>
<b>INM</b>	<b>Integrated Network Management</b>
<b>KPI</b>	<b>Key Performance Indicator</b>
<b>MAaaS</b>	<b>Mobility Analytics as a Service</b>
<b>ML</b>	<b>Machine Learning</b>
<b>PT</b>	<b>Public Transport</b>
<b>RIA</b>	<b>Research and Innovation Action</b>
<b>SMART</b>	<b>Specific, Measurable, Action-oriented, Relevant, Time-bound</b>
<b>UC</b>	<b>Use Case</b>

**Table 1 – Terminology**

## Executive Summary

---

This executive report provides an overview of the use case scoping within the EMERALDS project, specifically focusing on Task 5.1: Use Cases Orchestration & Validation. The EMERALDS project aims to develop Mobility Analytics as a Service (MAaaS) methods to enhance urban mobility and transportation decision-making. This report outlines the purpose, scope, and objectives of the use cases, emphasizing the coordination, stakeholder engagement, and validation processes involved.

The three use cases in the EMERALDS project serve as real-world scenarios to showcase the application of extreme data analytics and Artificial Intelligence (AI) techniques in urban mobility and transportation. The goal is to develop innovative solutions that address complex mobility challenges and contribute to scientific knowledge in the field. The use case scoping document establishes a clear understanding of the research focus, objectives, and challenges associated with urban mobility data. The first use case comprises risk-assessment and forecasting of crowds during events in the beach area of The Hague. The second use case applies traffic network analytics on estimated and predicted multi-modal traffic states to improve traffic management operations in Rotterdam. The third use case aims to investigate public transport network efficiency and passenger trip patterns in Riga by inferring trip characteristics and analysing traffic flows.

MAaaS refers to the provision of analytics capabilities and solutions related to mobility and transportation as a service. The EMERALDS project adopts a modular approach, developing distinct components of MAaaS services in various work packages. T5.1 takes a proactive approach in coordinating the use cases, fostering integration, co-creation, and engagement with relevant stakeholders.

The use case scoping document aligns the use cases with the overall goals of the project and establishes a structured framework for conducting research and development activities. The document also emphasizes knowledge transfer and reproducibility, recognizing the importance of collaboration and sharing across domains.

WP5 aims to explore and evaluate the applicability of extreme-scale data acquisition, processing, management, analytics, and visualization tools in selected use cases. The use cases are designed to address specific research objectives and challenges set by the project's work packages. T5.1, along with the associated deliverable D5.1, coordinates the use cases and facilitates a thorough assessment and generalization of the MAaaS methods developed within the project.

The use case scoping document provides a detailed description of each use case, including background context, SMART problem definitions and objectives, key performance indicators, computational and data resources, mapping to the EMERALDS toolset, implementation plans, and key stakeholders involved. The document is structured to facilitate a comprehensive assessment of the use cases and guide the research and development process.



## GA Matrix of alignment

GA Component Title	GA Component Outline	Document Chapter(s)	Justification
<b>Deliverable</b>			
D5.1 - Use Cases Scoping Document	Report outlining the use cases, providing the context and objectives for each use case. Moreover, the data resources, the implementation plan and the key stakeholders are provided for each use case. In a separate chapter, the procedure to assess the functionality of the Emeralds is described in more detail.	<b>Chapters 2, 3, 4, 5</b>	This deliverable provides the scoping of the three use cases in the EMERALDS project, covered in Chapters 3, 4, and 5 for use cases 1, 2, and 3, respectively. For each use case a description is provided, including the use case objectives and the data sources that are available in each use case. Chapter 2 describes the assessment approach, that is, a general description of the different aspects of the assessment (technical, functional, etc.).
<b>Tasks</b>			
T5.1 – Use Cases Orchestration & Validation	T5.1 entails the coordination of the case studies T5.2-T5.5. Next to the smooth running of the tasks, we aim to maximise the potential for a thorough assessment and generalisation of the MAaaS methods, looking at the value of the toolset developed in the project. T5.1 aims to maximise the possibilities to identify the added value of the developed innovations and coordinates the use cases in that direction. It will develop an implementation plan that integrates and co-creates innovations engaging all relevant stakeholders. Each use case will document the KPIs and timeline, including intermediate milestones. Furthermore, it will acquire key stakeholders buy-in on the implementation plan, identify key resources and datasets. This task aims to apply an end-to-end validation procedure to evaluate performance and effectiveness of the new services compared to the initial situation (Baseline)	<b>Introduction, chapter 2</b>	The alignment of the use cases is provided in the introduction of this deliverable. Each use case is described in more detail in the respective chapters 3, 4, and 5. The validation is addressed through the assessment approach described in chapter 2 of this deliverable.
T5.2 – Use Case 1 Risk	T5.2 aims to design a risk assessment, prediction and forecasting toolkit for	<b>Chapter 3</b>	This chapter contains a description of use case 1,

<p>Assessment, Prediction and Forecasting in crowded events</p>	<p>regular conditions and events, exploiting extreme data analytics and AI methods. Risk is defined in terms of the possibility that people may be injured, or that damage will be inflicted in other ways. This can be due to crushing (e.g., too many people at the event location), but also due to situations getting out of hand in other ways, e.g., riots. The first step in T5.2 is to identify all relevant data sources (social data – including Instagram images, dedicated sensor data, Resono data, police reports, weather data, event data, etc.); after collecting all relevant data, two types of data analyses are preformed: the first is to look into prediction and forecasting of crowding; the second will be to identify the risk context and which factors are relevant to determine the aggravation circumstances. The second step pertains to the building the tools for prediction and forecasting on the one hand (left side of the bowtie), and to determine the current / predicted / forecasted risk level on the other hand (right side of the bowtie). In the third step, the pilot will be run with the actual stakeholders providing feedback on its performance and usability. The final step involves the assessment, the identification of additional data sources needed to enhance performance, and an overall analysis of broader deployment and business case.</p>	<p>situated in the Hague, consisting of a background, the use case objectives, the overview of data resources, KPIs and a mapping to the EMERALDS toolset. In addition, an overview is provided of the key stakeholders and their buy in to the project activities, as well as an implementation plan.</p>
<p>T5.3 – Use Case 2 Multi-modal traffic network management</p>	<p>T5.3 aims to design a multi-modal traffic modelling and prediction module comprised of tightly coupled processing elements such as extreme scale data collection of big sensory data real-time and long-historical data. Self-identification of underperformance and self-tuning of traffic control algorithms, including multi-modal network decomposition; Mobility AI-driven data fusion and predictive analytics methods for identification, assessment, and solvability of bottlenecks; traffic simulator which boosts the raw sensory data dataset into rich training sequences; traffic prediction model which learns from the training data set; route calculation as a service. Main data input into the system will be the microscopic and aggregated data from the various sensors (radar, loop sensors, etc) a large historical data set of floating car data (FCD). FCD is represented by geo position and the speed of vehicle</p>	<p><b>Chapter 4</b></p> <p>This chapter contains a description of use case 2, situated in Rotterdam, consisting of a background, the use case objectives, the overview of data resources, KPIs and a mapping to the EMERALDS toolset. In addition, an overview is provided of the key stakeholders and their buy in to the project activities, as well as an implementation plan.</p>

	<p>sensed approximately each 5 seconds from navigation devices, that is from millions of devices every day over the period of several years. For proof-of-concept purposes, the component will operate on Rotterdam counting thousands of vehicles daily. In EMERALDS the key processing components of traffic modelling suites will be advanced by incorporating extreme scale data processing and analytics services. The use of efficient AI methods will allow the traffic management functionalities to be extended for multi-modal traffic and to collect and process more data, while addressing all privacy and security concerns.</p>		
<p>T5.4 – Use Case 3 Trip Characteristics Inference and Traffic Flow Data Analytics</p>	<p>Tools advanced within the projects will be used to provide value added data analytics for fostering sustainable mobility in the city of Riga, mainly by drawing important inferences for public transport network optimisation. Riga use case will contribute to defining functionality of the tools/services and then later will be tested. The work in the project will be used to create sustainable mobility and urban planning analytics services. The prototype algorithm for joining Riga Public transport data sets, estimating entry points and trip chaining for exit points will be implemented for better efficiency building on the results of WP3-WP4. G93 has been involved in numerous sustainable mobility planning projects where it made use of heterogenous data.</p>	<p><b>Chapter 5</b></p>	<p>This chapter contains a description of use case 3 situated in Riga, consisting of a background, the use case objectives, the overview of data resources, KPIs and a mapping to the EMERALDS toolset. In addition, an overview is provided of the key stakeholders and their buy in to the project activities, as well as an implementation plan.</p>

**Table 2 – Matrix of Alignment**

# 1 Introduction

---

The work presented in D5.1 is carried out under T5.1 Use Cases Orchestration & Validation. Task 5.1 in the EMERALDS project is responsible for coordinating the case studies outlined in T5.2 to T5.5. Its main objectives go beyond ensuring the smooth execution of these tasks. T5.1 aims to maximize the potential for a comprehensive assessment and generalization of the Mobility Analytics as a Service (MAaaS) methods developed within the project, focusing on identifying the added value of the toolset.

MAaaS refers to the provision of analytics capabilities and solutions related to mobility and transportation as a service. It entails leveraging data, technologies, and algorithms to extract insights, optimize operations, and improve mobility decision-making. Leveraging MAaaS solutions, organizations and stakeholders in the mobility sector can access a range of analytics services and tools without the need for significant investments in infrastructure, expertise, and resources. In EMERALDS, a modular approach is pursued where the distinct components of MAaaS services are developed in WP3, WP4 and Industry Partners from WP2, WP4 and WP5 combine this innovative software in their centralized platform or system where users can upload their mobility data, apply various analytical techniques, and obtain valuable insights and actionable recommendations. The result consists of a modular functional spectrum of enabling applications that are expected to deliver improved data analytics performance in all environments constituting of multiple resources in a computing continuum.

This document sets the stage for rigorous research and collaboration, aiming to develop innovative solutions that address complex urban mobility challenges and contribute to scientific knowledge in the field of urban mobility data analytics as well as mobility and transportation planning. Selection of the use cases has been performed based on showcasing a variety of analytics tasks and employing different combinations of the assets to successfully execute data workflows.

It should be noted that further refinement of the use case descriptions is anticipated under T5.2, T5.3 and T5.4, since the project is following an agile approach (steering 2 design sprints) but also due to the fact that additional computational and data resources will become available to the consortium. The UC tasks will further engage with different stakeholders from their local ecosystems over 3 cities (The Hague, Rotterdam, and Riga) to best determine business requirements, approaches, and the most effective way to present the results to system operators.

## 1.1 Purpose and Scope of the Document

---

The purpose of this document is to provide a comprehensive overview of the use case scoping for the application of extreme data analytics and Artificial Intelligence (AI) techniques in the context of enhancing urban mobility and transportation over a set of selected use cases tied with urging needs of city authorities. It outlines the scope, objectives, and expected outcomes of each individual use case (T5.2-T5.4), serving as a guide for the research and development activities to be undertaken.

Importantly, D5.1 aims to establish a clear understanding of the scientific goals and research focus of each use case, highlighting the specific challenges and opportunities associated with urban mobility data, whilst setting the ground for harnessing the potential of the data to deal with real world problems. It provides a structured framework for conducting in-depth investigations, developing innovative solutions, and evaluating the effectiveness of the EMERALDS toolset in addressing these challenges. In each use case description, SMART (Specific, Measurable, Action-oriented, Relevant, Time-bound) problem definitions have been incorporated with the aim to guide the research and development process towards specific, measurable, achievable, relevant, and time-bound objectives.

Once the first version of the EMERALDS toolset will be released in M18 (D2.2), testing of the tools within the use cases will be conducted with the involvement of stakeholders who will provide



feedback on its performance and usability (featured in D5.2, D5.3 and D5.4). This feedback will inform further enhancements and refinements to the EMERALDS toolset, leading to the revised version of the toolset in M33 (D2.3) and Second Validation Cycle Deliverables (D5.5, D5.6 and D5.7). Together with WP2, specifically D2.1 and D5.1 constitutes the verification means for project milestone 1 MS1 (*Reference Architecture & Performance KPIs*) reached on M9.

## 1.2 Relation to Work Packages, Deliverables and Activities

---

D5.1 Use Case Scoping Document serves as a scientific and business roadmap for the project, hence has significant value for the EMERALDS work plan, providing a comprehensive understanding of the research focus, objectives, data requirements, and implementation plans for each use case. Through the UC to EMERALDS toolset mapping, the identification of how each use case aligns with and supports the scientific, technical, and wider goals of the project is achieved, ensuring that the use cases are designed to address the research objectives and challenges set-up by WP3, WP4 and WP2. Each use case has been allocated a task within WP5, specifically UC1 is carried out in T5.2, UC2 in T5.3, and UC3 in T5.4.

Given that each use case brings together different stakeholders, WP5 has a strong interest in strengthening the stakeholder engagement and interest in testing of the EMERALDS innovations, therefore actively collaborates with WP6 in all stages of the use case experimentation plan, undertaking a series of workshops and communication activities.

The independent ethics advisor (IEA) attached to the project collaborates with UC stakeholders to develop an ethical framework that guides the design and implementation of the UCs. In this direction, potential ethical risks and implications are identified and timely mitigated ensuring that UC activities align with ethical standards. The ethical compliance monitoring procedure is continual and comprises a concrete oversight mechanism promoting ethical awareness, accountability, and responsible AI practices throughout the project lifecycle.

In addition to establishing an interdisciplinary approach, by combining expertise from different domains, such as mobility studies, transportation, AI, data science, and urban planning, D5.1 Use Case Scoping Document also emphasizes the importance of knowledge transfer and reproducibility. The report recognises that collaboration and knowledge sharing across domains are essential for achieving meaningful and impactful outcomes throughout the duration of the EMERALDS project.

## 1.3 Contribution to WP5 and Project Objectives

---

WP5 examines and assesses the applicability and potential of the developed Extreme Scale Data acquisition, processing, management, analytics, and visualization tools within the context of urban mobility. This will be achieved by conducting evaluations on the defined set of use cases, (with the work performed in T5.2, T5.3 and T5.4), allowing for a comprehensive exploration of the capabilities and effectiveness of the tools (developed in WP3, WP4) in addressing urban mobility data processing and analytics challenges. Selection of the use cases has been based on showcasing a variety of analytics tasks and employing different combinations of the assets.

The overall goal of WP5 is validating the solutions' advancements/improvements in handling data-intensive processes in an iterative procedure, where developments are fine-tuned with use case leaders, directly contributing to project objective O4 ***Demonstrate and measure the efficiencies of the novel Extreme Scale Analytics services through three pilot use cases and validate the concepts and tools usefulness as well as overall improvements in extreme data workflows through two early adoption applications.***

T5.1 - and the associated D5.1 - focus on the coordination of the case studies and the subsequent early adoption demonstrators T5.2-T5.5. In a further extent, T5.1 aims to maximise the potential for a

thorough assessment and replication of the achieved improvements from the EMERALDS tools and to multiply the possibilities to identify the added value of the developed innovations (e.g., Mobility Analytics as a Service) and directs the use cases in that direction. Each use case description documents the KPIs and timeline of software experimentation, including intermediate milestones and maturity checkpoints fostering the iterative assessment. In this regard, the use cases serve as a maturity testbed for the developed tools in order to affirm the targeted TRL improvement using real world data and creating conditions similar to a real-world scenario. The testing and validation process in EMERALDS is organized through two agile development sprints, where software developers, system integrators and use case owners collaboratively approach the urban mobility data research questions and tackle the pertinent challenges. D5.1 contributes to project milestone 1 (*Reference Architecture & Performance KPIs MS1 -M9*) setting the ground for the application of an end-to-end validation procedure to evaluate performance and effectiveness of the new services compared to the baseline performance of UC assets and processes.

## 1.4 Structure of Use Case Scoping Document

---

D5.1 provides a comprehensive and detailed account of the evaluation and implementation process, highlighting the KPIs, time plans, resources, datasets, and technical solutions foreseen in each use case. It will serve as a reference document for assessing the performance, effectiveness, and technical implementation of the developed tools as it captures the baseline situation of existing tools and approaches in the use cases. Each of the use cases descriptions provided is addressing the following topics:

- Background- Context – Challenges
- SMART problem definition and Objectives
- Key Performance Indicators
- Computational and Data Resources
- Mapping to EMERALDS toolset, developments and learnings addressed
- Initial Implementation Plan, and
- Key Stakeholder roles

The Use Case Scoping Document is structured as follows: Chapter 1 introduces the report, the purpose of the document, highlighting the relation to other work packages and the project's objectives. Chapter 2 describes the approach adopted for assessing the use cases. It provides background information on the assessment methods used in the EMERALDS project and explains the focus of the assessment. It further details the technical, legal, functional, user-acceptance, impact, and socio-economic assessment approaches employed. Chapter 3 focuses on the first UC Risk-assessment, prediction & forecasting during events, Chapter 4 on the second UC Multi-modal integrated traffic management and Chapter 5 on the third UC Public Transport Trip Characteristics Inference and Traffic Flow Data Analytics. A brief summary and conclusion of the report is given in Chapter 6. Last, Chapter 6 describes the next steps to be taken based on the information presented in the document.

Ethical considerations related to D5.1 are addressed in the Ethics Checklist document included in the Annex.

The deliverable and UC requirements, descriptions and implementation plans are the fruit of several one-to-one meetings and numerous conference telcos organized with all involved partners.

## 2 Assessment Approach

---

The aim of the EMERALDS project is to design, develop and create an urban data-oriented Mobility Analytics as a Service (MAaaS) toolset, consisting of the proclaimed EMERALDS services. These services will be applied (demonstrated) in the EMERALDS use cases, which provide the collection and management of large volumes of urban mobility data from various sources such as sensors, GPS devices, proprietary data, and transportation networks.

The EMERALDS services can be used in a variety of applications, so the aim of the UCs is to provide a diverse range of applications. In UC 1 and 2, these applications mostly consist of traffic or crowd management applications, in which information and/or traffic signals are provided to the traffic and crowd managers, who, if necessary, may deploy a variety of traffic / crowd management measures. The (end-)users (general public, i.e., travellers, *and* the decision makers) will react (by changing their travel decisions, their driving behaviour, the interventions they deploy, etc.) to these measures, and with their changed behaviour the overall traffic and transport system operations will be impacted [1,2]. In UC 3, the Emeralds will be used to optimise the public transport network. This optimisation will also lead to a user response (e.g., more people using public transport), with the corresponding impacts on the traffic and transport system.

In the Grant Agreement (GA), the UCs have been defined on a conceptual level, and also a first idea has been given on the data availability for each UC. The UC owners have translated these conceptual GA descriptions to concise business, data, software, infrastructure, functional and non-functional requirements and SMART objectives, as described later in this document. In this translation process, the UCs have been further adapted to UC end-user needs, based on which user stories have been defined. These user stories show the added functionality of the emeralds for an end-user of the UC system. They also present the requirements of the UCs to the technical partners, who have documented their draft tool design specifications and opted functionalities. The fruition of this dialogue is encapsulated on the UC-EMERALDS mapping, signifying the agreement on pursuing an innovation transfer procedure, entailing specific integration steps and the generation of the Data Workflow per objective in pipelines. The UC systems are configured to execute pre-determined tasks of the workflow in the necessary environment, but in the end deliver the outputs through dashboards, interfaces, data-streams or files to an end-user.

The EMERALDS services are available through open software repositories, the ATOS platform and certain functionalities are accessed through CARTO services. In addition, UC owners offer in-house components. Each UC uses different EMERALDS services as well as specific parts of the reference architecture as described in D2.1. These will be specified in the description of the respective use cases later in this document.

### 2.1 Background on Assessment Methods

---

Cities engage in competition for employment opportunities and investments. In this respect, the quality of life emerges as a crucial determinant, especially for high value-added businesses. A favourable quality of life gives citizens access to locations that meet their basic needs; it also encompasses access to diverse cultural events and recreational activities. Nevertheless, the essential mobility required to make destinations accessible has major social, environmental and economic impacts [3-5]. In recent years, the notion of smart urban mobility has emerged to tackle various needs, including mitigating the adverse impacts of mobility and negative externalities of urban transportation systems [6,7]. This has resulted in the design of a large variety of systems, such as [8,9], all claiming to have a large, beneficial, impact, but impeding concise cross-comparisons. Various authors have therefore developed assessment methods [10-13] addressing the key concerns of mobility such as

vehicular emissions, transportation cost, travel time, congestion, safety, energy consumption, accessibility, and social equity [14-16].

As nowadays more and more data on mobility become available, in multiple resolutions sourcing from an increasing number of interconnected devices, these data can be used to provide insights in the mobility system, and as such support the development of smart urban mobility in various ways. The emeralds form the linking pin between the data and the development of smart urban mobility. For the assessment of the emeralds in this project, elements of the abovementioned literature will be used, in particular the KPIs. Choosing the appropriate KPIs does however require a profound understanding of the relation between data / information quality (accuracy, reliability, granularity, etc.) and the effectiveness of the considered smart mobility applications. As such these KPIs will be specific for each UC, but the overall assessment approach will be similar for all UCs, as explained in the next section.

## 2.2 Focus of emeralds Assessment

---

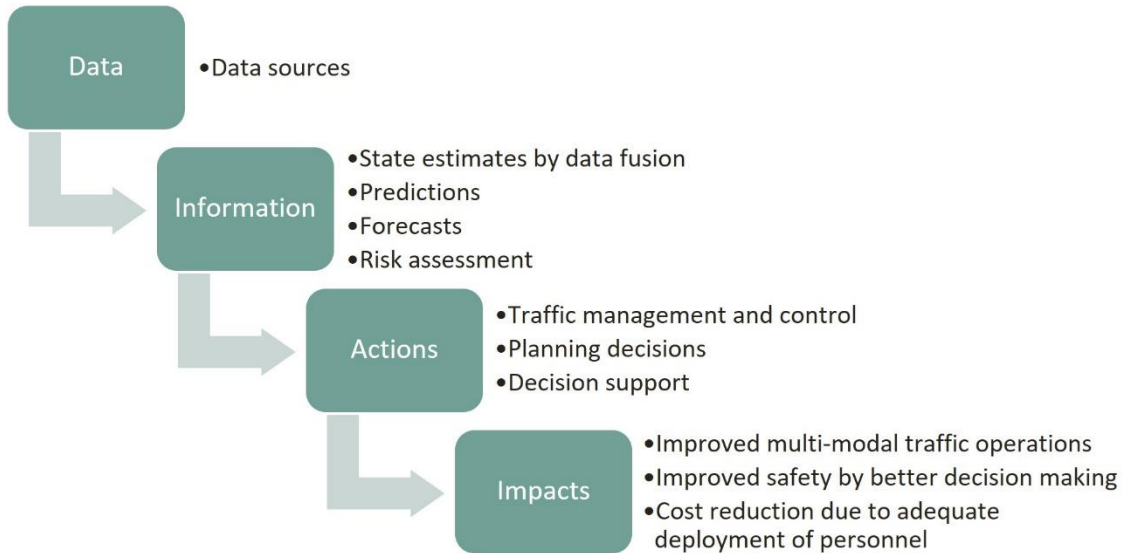
An assessment typically consists of evaluating a variety of key performance indicators (KPIs). These KPIs should cover the different functionalities of the system to be assessed, under different conditions. During the 4<sup>th</sup> steering committee of the EMERALDS project 4 categories of key performance indicators have been identified (below each category some examples are mentioned):

- 1) Software performance (GA, measured by WP3, WP4).
  - a) Storage-querying-analytics backbone can act as a data fabric system reducing data integration time by >20%.
  - b) Porting >30% of analytics pipeline to the edge.
  - c) Perform >50% of sensitive data analytics tasks in-situ.
- 2) UC specific software performance KPIs (measured by WP3 and WP4, D5.1, detailed in Chapters 3,4,5 measured and presented in D5.2-D5.7).
  - a) Percentage of correctly predicted crowdedness (UC1).
  - b) Processing times for real-time traffic data (UC2).
  - c) Accuracy of fusing GPS, GTFS, radar, and other sensor data to create comprehensive trip characteristics (UC3).
- 3) UC business KPIs (D5.1, detailed in Chapters 3,4,5 and reported in D5.2-D5.7).
  - a) Reduction in costs related to hiring staff using better predictions of crowdedness (UC1).
  - b) Reduction in monetized travel time due to more accurate queue length prediction as input for integrated network management (UC2).
  - c) Reduction in passenger delays by optimization of existing PT network using information on mobility hubs (UC3).
- 4) Impact and communication (measured by WP6).
  - a) Scientific Publications
  - b) Press releases.
  - c) Webinars

We emphasize that the four categories are related, as is indicated by the so-called value chain shown in Figure 2-1. The value chain shows how the “software performance” (in the value chain reflected in by the arrow between (raw) data and information) affects the actions taken in the considered smart



mobility applications (which become more timely, more effective, etc.), which in turn affect the impacts of these actions.



**Figure 2-1 – Data Value chain of urban mobility use cases.**

For the assessment of the use cases, the second (UC specific software performance KPIs) and third (UC business KPIs) categories of KPIs are most relevant.

The UC specific software performance indicators are directly related to the EMERALDS services to be developed in the technical work packages (WP3 and WP4). In the corresponding deliverables, these emeralds a technical performance check already takes place. However, in the UCs, these emeralds are functioning in an application that requires a specific accuracy of the emeralds, which justifies a dedicated software performance assessment within each UC. The UC specific software performance assessment of the emeralds is based on indicators that are frequently used to assess AI from a technical / software perspective, that is focusing on the accuracy of the estimations and predictions. This holds both for the assessment by the technical partners and within the UCs. This UC specific software performance assessment or technical assessment is pivotal in the EMERALDS, as we are taking tools from Research, at most TRL-4 and showcase their TRL 5 version within UC partner systems.

For the UC business KPIs the EMERALDS services are assessed based on the consequences of applying the EMERALDS services on the processes (businesses) of the various stakeholders. As the EMERALDS services will not be implemented in these processes, these consequences cannot be measured directly. The assessment will therefore be based on structured interviews with the stakeholders in the UCs, during which simulations are shown of the UC scenarios using the EMERALDS services and the stakeholders can give a qualitative assessment of the consequences.

As each UC uses different emeralds, there is no alignment with respect to the assessment of specific emeralds. However, there will be an alignment in time, where the software specific performance assessment will be done first, and reported in D5.2-5.4, and the business KPIs are assessed in the second part of the project and reported in D5.5-D5.7.

More details on the generic assessment approach of both the use case specific software performance KPIs and the use case business KPIs can be found in section 2.3.



## 2.3 General Assessment Approach

---

In the following, the general assessment approaches are further detailed for the UC specific software performance indicators (section 2.3.1) and for the UC business KPIs (section 2.3.2).

### 2.3.1 UC specific Software Performance Indicators

The UC specific software performance indicators are directly related to the emeralds services to be developed. An essential aspect of evaluating individual software service performance within UC execution environments involves conducting rigorous assessments distinct from mere benchmarking exercises. Their assessment is therefore based on indicators that are frequently used to assess AI. [17] recommends using multiple indicators, preferably with independent characteristics, in order to resolve some of the limitations of individual metrics to ensure the validity of the models. One of the suggested combinations is  $R^2$ , RMSE and the 20-index [17]. The required accuracy needs to be aligned with the needs of the respective stakeholders of each use case. The specific indicators as well as the required accuracy are further specified for each UC, so for UC 1 in section 3.4, for UC 2 in section 4.4 and for UC 3 in section 5.5. In those sections, also the alignment with the technical partners with respect to their assessment of the emeralds is addressed.

### 2.3.2 UC Business KPIs

The performance of the EMERALDS services in relation to the business of the stakeholders for each UC cannot be directly observed, as the EMERALDS services will not be implemented in the systems of the stakeholders. Nevertheless, to get an impression of the added value of the EMERALDS services, we use structured interviews with the stakeholders, and ask them to assess the EMERALDS services using simulations for the scenarios that have been defined for each UC.

The workflow for this assessment is as follows. First, the stakeholders that will participate in the assessment are identified. Together with these stakeholders, in which parts of their business (processes) the emeralds may have an effect is examined. Examples include the planning of staff size to perform crowd management measures (UC 1), the management scheme of a multi-modal intersection (UC 2) and the operations of a public transport operator (UC 3). For each combination of process and stakeholder, the business KPIs are determined. These are the assessment indicators that will identify the added value of the emeralds for the business of the stakeholders. This part of the assessment has been included in this document, for UC 1 in section 3.4, for UC 2 in section 4.4 and for UC 3 in section 5.5.

Each UC has identified scenarios. These scenarios will be quantified in detail, so the impact of the emeralds can be identified for different situations, e.g., under which conditions is the added value of the emeralds expected to be larger or smaller. To this end, historical data are combined with the emeralds for a variety of conditions. Stakeholders are then presented with the results of these scenarios and are asked about the impact of the emeralds on their workflow, addressing the different KPIs. Here, a comparison will be made between the scenarios, and compare with the current situation, so their added value will become apparent.

## 3 UC1: Risk Assessment, Prediction & Forecasting during Events

### 3.1 Background – Context – Challenges

Growing population and urbanization lead to an increased probability of overcrowding and risks of incidents during events. With events we mean the broadest sense of the word. Examples of events are not only organized events such as festivals or sport events, but also an ordinary shopping day in the city, a visit to beach for pleasure or rush hours at a train station. Looking at the Scheveningen beach area in The Hague, in the last years more and more incidents have been reported n due to crowdedness (Figure 3-1 Media publication reporting on the Scheveningen overcrowding incident in 2022).

## Gemeente Den Haag: kom niet meer met auto naar Scheveningen

17 jul. 2022 in BINNENLAND



**DEN HAAG** - Wie naar het strand wil met de auto, kan Scheveningen zondagmiddag beter mijden. De gemeente Den Haag roept mensen op niet meer met de auto naar Scheveningen te komen, want er is geen parkeerplek meer. „Ook de meeste parkeergarages zijn vol”, twittert de gemeente.



© REGIO15

Figure 3-1 Media publication reporting on the Scheveningen overcrowding incident in 2022

More and more data has become available to measure the historical and current crowdedness at locations [18]. However, to properly manage the crowdedness and mitigate risk, proper prediction and forecasting models are needed.

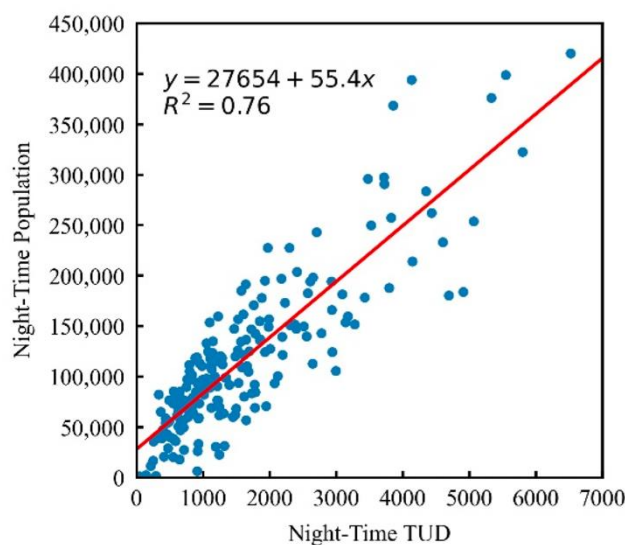
### 3.1.1 Need for Risk Identification, Prediction and Forecasting

Several challenges such as risk identification, state estimation (e.g., by data fusion), prediction and forecasting are tackled within UC1.

#### Risk identification during events

While the number of incidents involving crowding at events is substantial, *identification* of these risks is a complex affair. This for one is due to the lack of ubiquitous data on crowding and crowding risk, for situations where incidents have actually occurred obstructing the derivation of a ground truth. Second, the sources of risk are manifold and diverse in nature. Mass gatherings pose several significant safety risks, including high crowd density leading to crushes and fatalities, particularly in confined spaces during concerts, sports events, or religious gatherings. Stampedes, often triggered by sudden rushes towards exits or entrances, can similarly result in injuries or deaths. Minor incidents can escalate, causing disproportionate crowd reactions and cascading risks such as infrastructure collapse. Terrorism and violence, including bombings, shootings, or acts of activism, can significantly endanger crowds, with the potential for crushes or stampedes in the ensuing chaos. Extreme weather conditions, like heatwaves or thunderstorms, pose risks of heat exhaustion, dehydration, or electrocution, with crowds moving abruptly to seek shelter possibly causing crushes. Structural issues with bridges, balconies, or temporary stands can lead to collapses under the stress of overcrowded conditions. Fire and explosion hazards increase with the use of pyrotechnics or flammable materials in crowded events. Additionally, the spread of communicable diseases is a heightened risk due to the challenge of enforcing social distancing and other preventive measures in densely populated settings [19].

Previous risk identification approaches have mostly focused on density (number of visitors in an area) as the key factor determining crowding risk. For instance, [20] propose a risk identification framework for city level crowding using data on average densities, employing Tencent user density (TUD) data generated by the location of smartphone users who use Tencent applications use. They show the extent in which these sampled data are representative for the number of people in an area; see Figure 3-2.



**Figure 3-2 – Example from Yang et al (2022) showing linear regression fit of the TUD data and the residential population of each subdistrict in the area within the Outer-ring Road of Shanghai, China.**

They use factors including place characteristics, management measures, and crowd characteristics using a standard “Risk = Probability x Impact” approach; see Figure 3-3. The relevant part of the picture shows how different factors (density, tendencies, dangerousness of the location, and the traffic situation) determine the probability of an incident; the crowding determines the severity, while the reputation and type of the location may increase or reduce the risk further. The factors together determine the overall “crowd-gathering risk”. For more details, we refer to the paper [20].

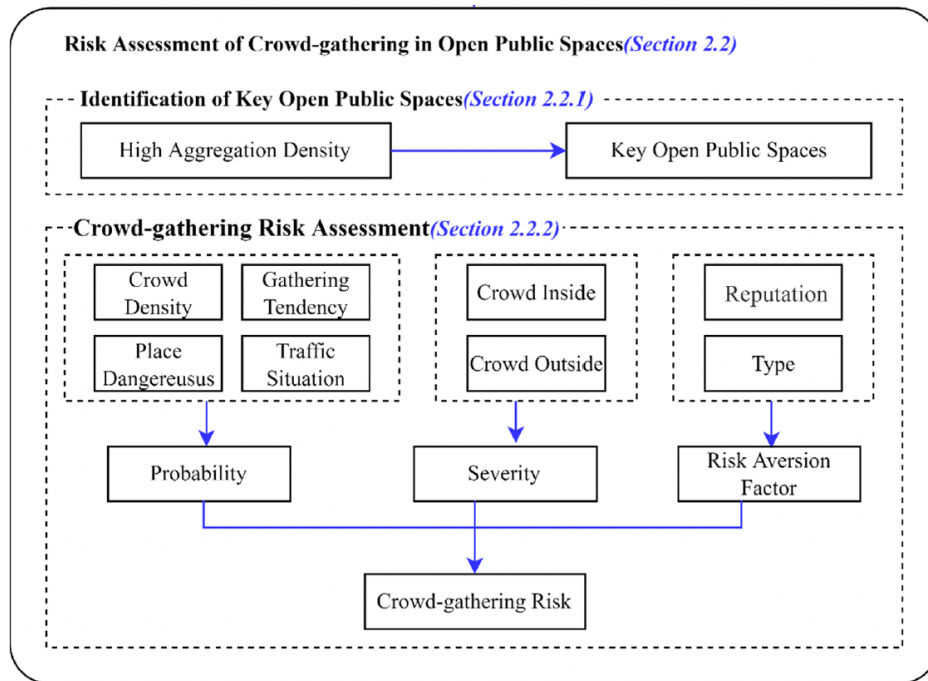


Figure 3-3 – Determining crowd risk (adapted from [20]).

[19] provides a roadmap for the future of crowd safety practise and highlights some of the key challenges. They emphasize that “One of the biggest challenges in crowd safety is the risk of crowd-related incidents.”. They furthermore highlight that the integration of new technologies is one of the key knowledge gaps, next to a comprehensive understanding of the complexity of crowding. This is particularly the case when it comes to the deployment of Decision Support Systems. This is one of the challenges that we aim to address in UC1.

### Short-term predictions and longer-term forecasts

Most of the literature on predicting crowd dynamics focuses on physical modelling and computer simulation models of crowds. For instance, [19] discuss the use of computer simulation in crowd risk practise. They discuss how computer simulation models are increasingly vital tools for enhancing crowd safety in mass gatherings, allowing planners and designers to test and refine crowd management strategies within a virtual setting. These models facilitate predictive analysis, helping to foresee crowd behaviours under various scenarios to pre-empt potential risks and identify solutions to prevent accidents, such as bottlenecks or overcrowding. They also enable the testing of emergency response plans, allowing for the identification and rectification of plan weaknesses in advance. Moreover, simulations are instrumental in optimizing crowd management strategies, including the strategic placement of entry/exit points, barriers, and security personnel, ensuring effective crowd control. Additionally, they serve as a valuable training tool for personnel, offering hands-on experience



in managing crowds and responding to emergencies in a risk-free environment. They conclude that, despite their benefits, the deployment of computer simulation models for crowd safety faces several challenges, primarily the need for accurate data regarding crowd characteristics, which is essential for creating realistic simulations. The complexity of modelling unpredictable crowd behaviours adds another layer of difficulty, potentially limiting the models' predictive accuracy and the effectiveness of suggested safety measures. The success of these models also hinges on the quality of data input, the assumptions about crowd behaviour, and the technological limitations. [19] conclude that addressing these challenges requires ongoing monitoring, improvement by developers, and effective communication between academics and commercial developers, alongside detailed guidelines for data handling and model interpretation to maximize the utility and reliability of simulation models in improving crowd safety.

The role of AI in crowd management is by and large limited to computer vision applications. Several applications focus on the prediction of pedestrian trajectories (e.g., [21]), for instance for AV applications. One of the few examples of using Machine Learning for predicting crowd dynamics is [22]. In their study focusing on pedestrian traffic data from Melbourne, Australia, the Koopman operator methodology is highlighted through the Extended Dynamic Mode Decomposition (EDMD) framework. This approach involves approximating the Koopman operator within a function basis, derived using time delay embedding and the Diffusion Map algorithm. The model successfully identifies the unique temporal patterns across 11 traffic sensors, rendering the complex dynamics into a linear framework for easier analysis. Through this linear perspective, the model is decomposed into spectral components, aiding in the interpretation of the traffic system's dynamics and improving the scientific understanding of both the system itself and the model's functioning. The spectral components revealed in this study are particularly insightful, as they correlate with the geometry of the underlying state space and demonstrate the model's stability across various forecasting horizons. This illustration of the Koopman operator framework not only demonstrates its potential for making explainable and accurate data-driven predictions in complex traffic systems but also suggests its applicability to other traffic environments, offering a significant advancement in the field of urban traffic analysis and management. The authors conclude that there are clearly exogenous factors missing that can improve prediction, including weather or public events. Also, their approach accounts for the stochasticity in the traffic and includes uncertainty information in the model predictions, but the impact of interventions in the predictions and forecasts.

The work performed in UC1 is carried out and monitored under T5.2.

### ***3.1.2 The Hague UC Background and Developments so far***

In Scheveningen, a district of The Hague bordering the sea, different authorities work together in teams to ensure the safety. These teams are together (e.g., The Hague and the law enforcement authorities) responsible for the safety of visitors, inhabitants, shops, beach club owners in the vicinity.

Managing large crowds does not start on the day itself but has a 10-day planning horizon. Until last year there were hardly any digital tools available to support this planning process, and not every team had access to the same information. The focus of this planning process was to determine the amount of personnel that is needed on a specific day. Because of the structural shortage of staff, these predictions are vital. The same personnel could also be deployed at a festival site or in a busy city centre. This planning process is already a joint process in which both Law enforcement and city authorities are involved.

On the day itself (8-hour horizon), the available personnel need to be at these locations that are most critical. The operational team is responsible for that. Most of the information necessary to make decisions was until recently entirely obtained by people on the ground (walking and cycling through the area). They kept the control centre updated of the observed condition. Again, there were hardly any digital support systems available. Efficiently allocating crowd management personnel or receiving an evidence based early indications that more personnel is needed at specific points of interest or overall, in Scheveningen, would assist in the mitigation of potential risky situations or even contribute to towards their prevention.

TUD and ARG started in 2022 a process to develop a first set of digital tools:

- ARG developed in the first half of 2022 a common operational picture of the crowdedness situation in different areas in Scheveningen in an online digital twin dashboard called "Digitwin". Digitwin is a cloud native web application which can be integrated with other data platform and ML ops environment via API. A first step was bringing together public and semi-public data sources and visualizing them in such a way that they were easy to comprehend. During busy summer days, mainly the real-time viewer was used. Since this is a web-based tool, it could be accessed at different locations, and everybody had access to the same visualized information. If teams wanted to evaluate certain days / events/ incidents, they have the opportunity to search and visualize the relevant historical data.
- In the second half of 2022, TUD developed a mid-term prediction model (see [23] and Figure 3-5 for some predictions) to help the tactical / planning team to get a prediction of the expected crowdedness and of risks based on that crowdedness. Also, a short description of the model is included in the next section. These predictions were expected to advance the planning process. An impression of the Digitwin with the predictions is shown in [23] and Figure 3-6. After a multi-variate regression analyses, advanced machine learning methods were developed and cross-compared. As the XGBoost framework resulted in the most accurate forecasts, TUD continued with detailing its design (feature selection, regression tree depth, number estimators) of the one-day and multi-day ahead predictions. Overall, it turned out that the predictions were sufficiently accurate. In illustration, the picture below shows the MAE values for the one-step ahead forecasts presented in [23].

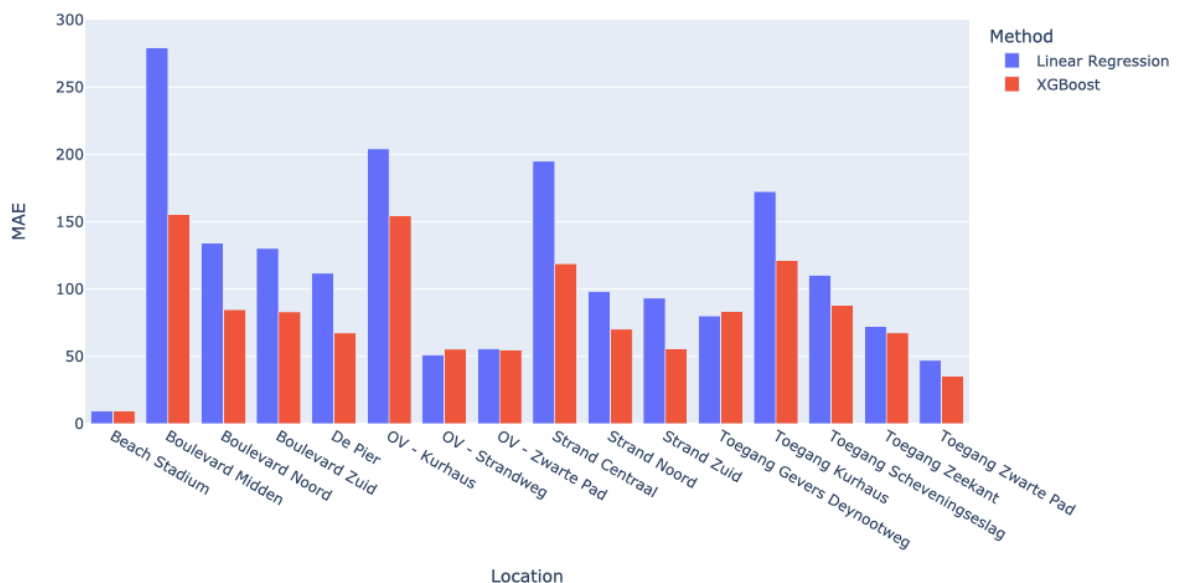


Figure 3-4 Comparison between forecasting accuracies for Linear Regression and XGBoost.

For some specific locations, however, additional input large scale data is likely needed to further improve the prediction quality: while the considered features (previous visits, temperature forecast, wind forecast) were suitable for most locations, for some additional information (e.g., event calendar) is likely required to improve prediction accuracy. That said, we can conclude that the proposed technology and methods are very promising and provide important support for planning events where crowding may potentially become an issue. However, further development of the forecasting models is necessary. In the current models limited historical data is used and just one type of data source is used. Increasing the scale of historical data ranges and different data sources is necessary to develop better forecasting models. This is where the Emeralds come into place.

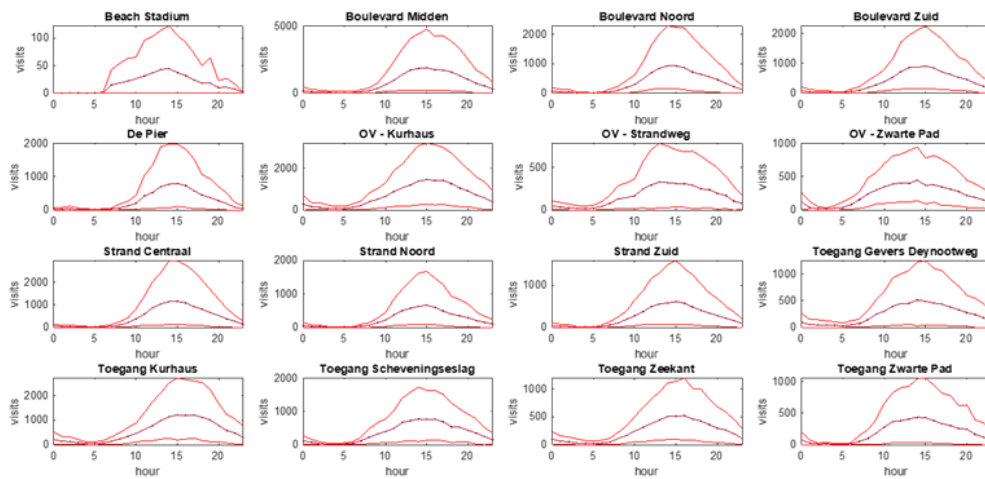


Figure 3-5 – Hourly forecasts of visits in the Scheveningen area.

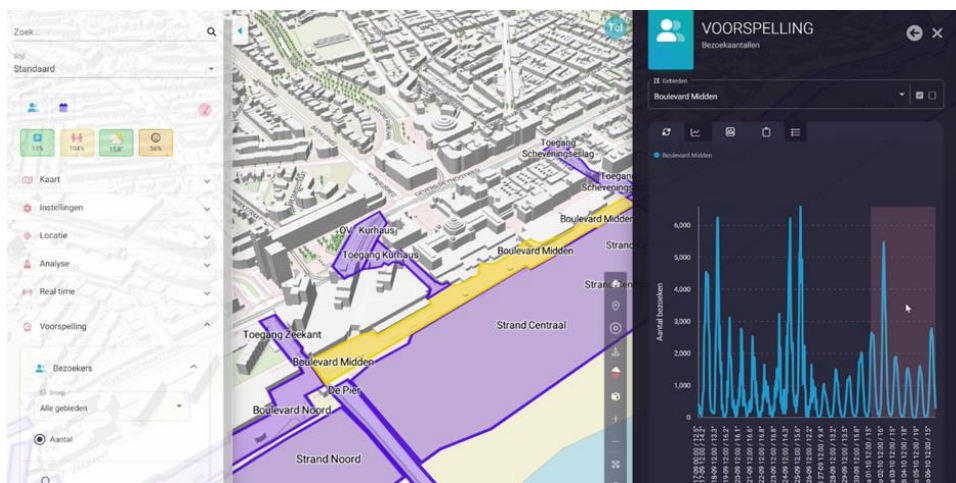


Figure 3-6 – The ARG DigiTwin showing the central beach area of Scheveningen (on the left) and the number of visitors and its predictions (on the right).

Note that the Digitwin is not integrated in Law enforcement or city systems. If something goes wrong, stakes are high. It is an additional tool that runs parallel to regular city and Law enforcement systems.



Emeralds will not be included in planning and operational Law enforcement or city systems. For demonstration purposes, emeralds will be tested in the operational environment of the Digitwin platform through integration points that are described in D2.2.

The underlying modelling approach for assessing the risks related with crowdedness is shown at Figure 3-7. On the left-hand side different factors that are relevant to estimate or predict crowdedness in an area are listed. The right-hand side of the 'bowtie' picture includes factors that determine risk levels. Indicatively, crowdedness is just one element that determines risk. The purpose of crowds densifying at a certain location, their moods, (changes) in weather conditions, the crowd safety infrastructure, etc., or personnel on site also have an influence on risk levels. They can both reduce / mitigate risk but can – in certain circumstances - also induce risk (by for example making people more aggressive). For instance, roads can be closed by traffic authorities to divert the traffic to other parking facilities to prevent traffic and ensure an ongoing flow of traffic. Or safety authorities can send to specific areas on the beach where the crowdedness is expected to reach critical levels, manage the crowd to other area and so prevents overcrowding.

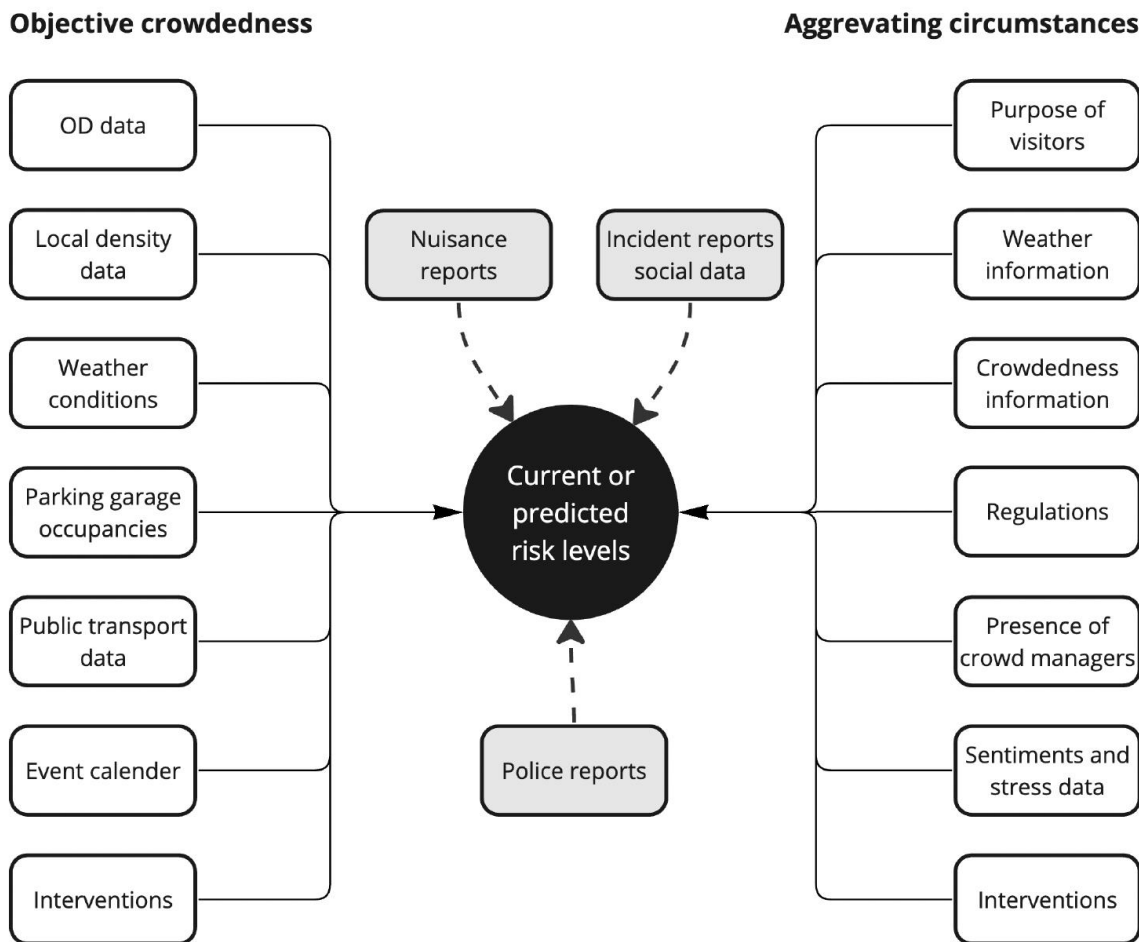


Figure 3-7 – Bow tie model for the crowd management system

### 3.1.3 Extreme Scale Data Challenges and Need for Emeralds

In the context of using diverse data sources for the risk identification, prediction, and forecasting of pedestrian crowds in cities, the use case faces extreme data challenges that stem from the heterogeneity, volume, and real-time processing requirements of the data. Integrating data from location-based services, smart cameras, counting mechanisms, social media, parking garages, and public transportation systems presents a complex scenario due to varying data formats, quality, and throughputs. The sheer volume and velocity of data generated necessitate advanced AI technology for efficient processing, analysis, and interpretation. Additionally, ensuring data privacy, security, and compliance with regulations adds another layer of complexity. These challenges demand robust data management strategies, sophisticated AI algorithms capable of handling high-dimensional data, and seamless integration techniques to accurately assess risks, predict crowd movements, and forecast potential safety concerns in urban environments. In light of this, EMERALDS toolset components such as mobility data fusion and crowd risk forecasting present an opportunity to improve data processes and design of models and pipelines for the understanding of complex urban mobility problems.

## 3.2 SMART Definition and UC Objectives

The Scheveningen area has appointed a crowd manager, who is responsible for the coordination between all stakeholders, the planning of the assets and defining measures and interventions. Each Tuesday, the crowd manager has a meeting with all stakeholders to discuss scenarios for the upcoming weeks, specifically the upcoming weekend. Decisions on measures and interventions are made based on event calendar, weather forecasts, holidays and historical personal experience. -

### Objective 1: Mid-term forecasting Scheveningen Beach

In order to improve the results of the decision-making process, the first objective of UC1 is to develop a data-driven forecast that predicts the expected crowdedness on Scheveningen in the upcoming 10 days.

Specific	<ul style="list-style-type: none"> <li>A) TUD developed a 6-day forecast of the expected visitor numbers in the different areas in Scheveningen on an hourly basis. The current model XGBoost already provided a reasonable prediction based on a small number of variables [26]. The event calendar was for example not included. That event information needs to be included, because regardless of weather conditions, the amount of people visiting a sports tournament or pop concert can be relatively high.</li> <li>B) Apart from including more variables, also the model structure itself needs to be improved or new model structures need to be explored.</li> <li>C) The current model contains a 6-day prediction horizon. The tactical team in Scheveningen prefers a 10-day prediction horizon. Extend the prediction time horizon of the model.</li> <li>D) Besides generating the predictions, the model should provide a clear estimate of the prediction uncertainty. Deviations of the crowdedness predictions need to be determined. Deviations are not necessarily caused by the algorithm, but are also inherited in training/testing data due to the inaccuracies of weather forecasts as well as noise in source data.</li> <li>E) Currently, the 6-day prediction model is calibrated on a monthly basis. Every time new predictions are needed; the calibrated model is applied with the most recent inputs. This way, we always have a calibrated model</li> </ul>
----------	---

	available. Preferably, we continuously retrain the model based on updated historical data.
Measurable	<p>Resono data are estimates of the amount of people in an area that will be used to assess the quality of the midterm prediction models. In these predictions both accuracy and reliability need to reach certain preset values.</p> <p>To assess the performance of the TUD model developed in 2022, we utilized several evaluation metrics, specifically Mean Absolute Error (MAE), Root Mean Square Error (RMSE), and Mean Absolute Percentage Error (MAPE). Each metric provides a different perspective on the model's error, and their combined use gives a more comprehensive evaluation of model performance. For instance, the MAPE, and to a lesser extent MAE, is less effective when dealing with target variables that take very small values since a small absolute error can turn into a large percentage error if the actual value is very small, leading to an overestimation of the error size. Similarly, MAE may not reflect the true error magnitude when the target variable values are low. For evaluating the multi-step ahead prediction, we computed the absolute error for each one-step ahead prediction and then averaged the errors.</p> <p><b>Accuracy:</b> Maximum error in number of people <u>per area</u> as well as maximum <u>average</u> error in number of people computed <u>over all areas</u>. Within the project team we determine the time period over which the errors are computed.</p> <p><b>Reliability:</b> Maximum error in number of people in an area may not exceed a certain preset value. This value will be determined within the project team.</p>
Action - oriented	<ol style="list-style-type: none"> <li>1) Enrich existing prediction model with additional data sources.</li> <li>2) Revisit model formulation, examine new algorithms and methods.</li> <li>3) Recalibrate prediction model.</li> <li>4) Determine number of variables suitable for providing mid-term forecasts.</li> <li>5) Adjust the time resolution of the model.</li> <li>6) Process input data with extreme scale filtering and analytics tools</li> </ol>
Relevant	<p>In order to plan the amount of personnel (from <u>The Hague and the Law enforcement</u>) needed in Scheveningen, we provide the tactical team responsible for personnel planning with an indication of the expected crowdedness in the different areas in Scheveningen.</p> <p>We assume that the developed model can easily be extended to other coastal areas in the Netherlands. Predictions can (likely in a simpler form) be shared with other stakeholders in Scheveningen, like Beach clubs and shop owners, public transport companies, inhabitants, etc. However, this can be politically sensitive. Gaining experience making midterm predictions in Scheveningen, we can try to extend these models to more complex situations like a busy inner-city after this project.</p>
Timebound	Assessing predictions can be done independent of time of year. The coastal team - whose will use predictions for operational and planning purposes - is active between April, 1 <sup>st</sup> and October, 31 <sup>st</sup> (summer season). Their input might be relevant for further model development.

## Objective 2: Short-term forecasting Scheveningen Beach

This objective focussed on the operational interventions during the day. Based on the midterm predictions the crowd manager has assets and interventions planned. However, the crowd manager can change and/or add interventions to reduce crowdedness and related risks during the day. For instance, redirecting traffic flows or change the (off)boarding location of the public transport stops.

In order to support in the decision-making process a short-term forecast will be developed, which will help the Crowd Manager with data-driven indicators.

Specific	<p>A) Apart from the midterm prediction model (objective 1) that can be used by the technical team for the planning of personnel, the operational team needs a short-term forecasting model to distribute the available personnel over these areas that require most attention based on expected crowdedness and related risks. This short-term model has an 8-hours prediction horizon in 15-minutes intervals. For the short-term prediction model more data sets are available. Besides historic crowdedness levels and weather predictions real-time information on crowds and weather can be used. Besides that, also real-time parking garage occupation, bicycle storage/parking occupation, traffic flows and FCD on roads towards Scheveningen, and transport schedules can be used. So far, TUD has not developed such short-term prediction model for Scheveningen.</p> <p>B) Apart from the predictions itself, the deviations of the crowdedness predictions need to be determined. Deviations are not necessarily caused by the algorithm, but also by inaccuracy in for example weather predictions that will be used as well as noise in the Resono data.</p> <p>C) Preferably, we continuously retrain the model based on updated historical data.</p>
Measurable	<p>Mobile app data are estimates of the amount of people in an area that will be used to assess the quality of the midterm prediction models. In these predictions both accuracy and reliability need to reach certain preset values.</p> <p><b>Accuracy:</b> Maximum error in number of people <u>per area</u> as well as maximum <u>average</u> error in number of people computed <u>over all areas</u>. Within the project team we determine the time period over which the errors are computed.</p> <p><b>Reliability:</b> Maximum error in number of people in an area may not exceed a certain preset value. This value will be determined within the project team.</p> <p>Note that error margins might be lower than in objective 1, because more precise input variables are available. For example, actual or short-term weather conditions instead of a 10-day weather prediction as input.</p>
Action - oriented	<ol style="list-style-type: none"> <li>1) Determine type of prediction model and develop it.</li> <li>2) Determine model input.</li> <li>3) Determine number of variables suitable for providing short-term forecasts.</li> <li>4) Adjust the time resolution of the model.</li> </ol>
Relevant	<p>The operational team needs a short-term forecasting model to distribute the available personnel over the different areas in Scheveningen; areas that require most attention based on expected crowdedness and related risks.</p>

	We assume that the developed model can easily be extended to other coastal areas in the Netherlands. Predictions can (likely in a simpler form) be shared with other stakeholders in Scheveningen, like Beach club and shop owners, public transport companies, inhabitants, etc. However, this can be politically sensitive. Gaining experience making short-term midterm predictions in Scheveningen, we can try to extend these models to more complex situations like a busy inner-city (see objective 4 and 5).
Timebound	Assessing predictions can be done independent of time of year. If responses of people in the field are important in further model development, it is important to test the model between April, 1 <sup>st</sup> and October, 31 <sup>st</sup> .

The use case description is based on the breakdown of general assumptions and approaches to scope the business and technical requirements.

- **Areas and points of interest:** The focus area is Scheveningen Beach in The Hague, due to the severity of crowdedness in high-risk incidents. Historical data and different event type (tourism and events) that take place at the location are available.
- **Forecasting time range:** Based on the business requirements of the use case two main forecasting time ranges are distinguished:
  - o *Short-term forecast.* Aggregating the next 8 hours based on a 15min (or less) time step. To support **short-term operational decision making**.
  - o *Mid-term forecast.* Aggregating the upcoming 10 days on an hourly time-step. To support **decision making on tactical level and planning purposes**.
- **Mobility data vs other data.** The initial forecast models are fed with mobility data (e.g., pedestrian count estimates, traffic flows and parking data), and weather predictions. During the project more and more other data sources will become available (e.g., social media, incident monitoring, and sound measurements) [24]. Therefore, a dynamic approach is required in which we assume not all data will be available from day one and we will further finetune the models based on onboarded data sources that come available in parallel with the EMERALDS project. Including more relevant data sources improves prediction quality as well as strengthens the robustness of AI/ML models.
- For most objectives, the approaches will be fully data driven. Knowledge about crowd behaviour and crowd dynamics, e.g., how people pass or cross one another, and with which speeds they move in areas given a certain crowd density, will not be used as input [25].

### **Dynamics of crowdedness and risks**

Crowdedness and risks can be managed by taking appropriate measures. Therefore, both crowdedness and risk are dynamic in nature which should be taken into account when making predictions. These dynamics are progressively considered throughout the developments taking place in WP3, WP4 and WP5. An initial planning has been drafted as follows: In year 1, these dynamics do not have to be taken into account when making predictions. In years 2 and 3, these dynamics become more and more important specifically when predicting risks. In turn, they will (if available) be included as input for predictions.

The quality of the predictions is not only determined by the quality of the developed EMERALDS tooling. Quality might differ between areas depending on the functionality of the areas. Therefore, the development of the EMERALDS tooling needs to be performed in close cooperation with TUD and

ARG. Since both organizations have been closely working together for a longer time with the city of The Hague, TUD and ARG have regular meetings with these cities to get a reflection on the results.

### **3.2.1 User story 1: Planning for safe and efficient operations during beach days**

- Scenario 1-0: current way of working before September 2022 – no midterm forecast was used. The crowd managing team used an infographic including simple classes based on weather forecasts and holidays. Experience of crowd managing team plays an important role.
- Scenario 1-1: 6-day mid-term forecast developed by TUD (XGBoost) in 2022 based on historic crowdedness and weather forecast.
- Scenario 1-2: 10-day mid-term forecast based on first round of data sets, e.g. historic crowdedness, weather forecast, event calendar and holiday calendar (as far as developed by EMERALDS and available in M18).
- Scenario 1-3: 10-day mid-term forecasting based on data sets mentioned under scenario 1-2 possibly enriched with additional data sets (as far as developed by EMERALDS and available in M36).

Note that in scenarios 1-1 to 1-3 the mid-term forecast were / will be used by the crowd managing team as an additional source of information (scenario 1-0) and not meant to fully replace it. The mid-term forecast models are not implemented in law enforcement or city systems.

### **3.2.2 User story 2: Crowd Management Planning of an event**

- Scenario 2-0: current way of working. Similar to Scenario 1-0. Since events occur less frequently, expert knowledge of the type of event in the type of visitors is more important.
- Scenario 2-1: 10-day mid-term forecast based on first round of data sets, e.g. historic crowdedness, weather forecast, event calendar and holiday calendar (as far as developed by EMERALDS and available in M36).

### **3.2.3 User story 3: Real-time monitoring and decision-making during an event**

- Scenario 3-0: no short-term predictions were used; estimates of crowdedness and potential risks were determined by personnel on site moving around Scheveningen (mainly by bike).
- Scenario 3-1: short-term predictions based on crowdedness data (Resono) (as far as developed by EMERALDS and available in M36).
- Scenario 3-2: short-term predictions based on multiple data sets, e.g., Resono data, parking garage data, public transport data, shared mobility data, and weather forecasts (as far as developed by EMERALDS in M36).

## **3.3 Computational and Data Resources**

---

### **3.3.1 Data Resources**

The table below lists the available data sources needed for the execution of UC1. As explained above this list is not exhaustive, since additional data resources will be onboarded throughout the first and second implementation cycle contemplating for dynamics in the crowdedness profiles studied. The



examined data sources relevant to the data analytics problem are indicated while specific data attributes that will be elicited in D5.2. A detailed description of the datasets is available in D1.4 Data Management Plan v.1, while corresponding metadata templates have been filled out. For modelling purposes, the datasets on crowdedness prediction will be split into training (observational data) and testing datasets during the development cycles in WP4 and WP3.

Name of the dataset	Description	Data Type	Related objective
<b>Event Calendar Scheveningen</b>	Date, times and additional information about events and holidays at Scheveningen Area	Time-based data	1, 2
<b>Parking data</b>	Real-time information on occupancy of public parking garages in Netherlands	Numeric	1, 2
<b>Public transport data</b>	Information on public transport schedules, timetables and real-time locations of buses, trains and trams.	Temporal-spatial traffic data	2
<b>Shared Mobility data</b>	Real-time location data of the parked shared mobility objects (scooters, bicycles) in the Netherlands.	Location of parked shared mobility objects, occupancy of designated parking areas	2
<b>Weather information</b>	Historical, real-time and future information on weather conditions.	Weather related data	1, 2
<b>Mobile App Counting</b>	Historical data on hourly visits on different beach areas in Scheveningen and The Hague Area	Numeric	1, 2
<b>Bicycle counting Netherlands</b>	Historical data on bicycle counting systems in the Netherlands.	Numeric	2
<b>Floating Car Data</b>	Travel time data based on floating car data collected from a smart phone app	Temporal-spatial traffic data	2
<b>Loop Detector Data</b>	Speed and flow data from double loop detectors in the network	Temporal-spatial traffic data	2

### 3.4 Assessment and Corresponding KPIs

UC1 concentrates on developing reliable predictions of various variables needed for crowd management. For the predictions – as indicated before – several evaluation metrics, specifically Mean Absolute Error (MAE), Root Mean Square Error (RMSE), and Mean Absolute Percentage Error (MAPE) should be used together. Each metric provides a different perspective on the model’s error, and their combined use gives a more comprehensive evaluation of model performance. For instance, the MAPE,

and to a lesser extent MAE, is less effective when dealing with target variables that take very small values since a small absolute error can turn into a large percentage error if the actual value is very small, leading to an overestimation of the error size. Similarly, MAE may not reflect the true error magnitude when the target variable values are low. For evaluating the multi-step ahead prediction, we computed the absolute error for each one-step ahead prediction and then averaged the errors.

The timing and reliable mid-term forecasts is most important for situations where an area or parking garage close to capacity. Based on that measures can be planned beforehand. For short-term predictions, next to timing your reliability and timing also availability and latency are important. The time window to take actions is often relatively short.

Therefore, next to MAE, RSME and MAPE, we will consider availability and latency as extra quality measures. Availability will be expressed in terms of the percentage of time predictions are available; latency describes the (extra) time needed from data collection to publication of the outcome of the Emeralds computation.

All these measures can be computed for the considered variables. For the estimations, the computations will be done using ground truth values, insofar available. For the predictions (say,  $k$ -minutes ahead), we will use the future observations when performing the assessment, e.g.:

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

Here,  $\hat{y}_i^k$  denotes the  $k$ -minute ahead prediction at time period  $i$ , and  $y_{i+k}$  denotes the observed value at time period  $i+k$ .

In general, we will work with 1-minute time periods (i.e.,  $y_{i+k}=1\text{min}$ ).

Success will be determined by considering the following aspects:

- The **quality** improvement of the estimation / prediction compared to the benchmark. This shows the (relative) improvement that is made by deploying EMERALDS data service.
- The **data and computation costs** (e.g., the EMERALDS data service can achieve similar quality without the need to collect / buy specific data).

The target values that determine the success are hence not straightforward to determine on beforehand. Moreover, for estimates and predictions different target values will be relevant. It is important to note that updates to the KPIs definition, measurement methods and impact will be delivered in M24 with D5.2, whilst further refinement and iterative design is carried out within the frame of T5.2.

For each KPI a norm or target value should be defined. This differs per situation and cannot be obtained from literature. These norms or target values / triggers should be discussed with the crowd management team and other boots on the ground.

The business KPIs follow directly from the scenario objectives. Below we list which qualitative business KPI's will be measured.

- Capacity planning of traffic managers. Based on the expected crowdedness staff is hired and planned to guide the crowds and traffic through the area. As currently the 'predictions' are based on personal experience this often leads to inefficient use (over-capacity) of staff. With a better midterm forecast a better staff planning can be made, which results in a more efficient use of personal and so lower costs.
- Police interventions. If during the day it gets too crowded and interventions are taken too late, the probability of incidents will increase, which leads to the necessary intervention of the police. If based on a better forecast the crowd manager can intervene earlier, the probability



of police intervention will decrease, and so policy capacity can be used on other places which results in higher general safety level.

- Economic impact. The beach has a high economic value due to the restaurant, hotels and other leisure activities. If certain hotspots or areas getting too crowded, visitors can decide to not go a leisure location which results to lower revenue. If the crowd is spread more evenly over time and over the whole area the capacity of the leisure hotspots is use more efficiently and so will lead to a higher economic impact of the whole area.

### 3.5 Mapping to EMERALDS Toolset

In this section, we describe the current approach of this use case in relation to the novel data services EMERALDS can provide. We present the pipelines designed to execute the scenarios in Figure 3-8. The above-mentioned datasets will be used to develop forecasting models which are used in the three user stories. The current pipeline design mostly serves as a reference guideline for the work in progress additional data flows and computational resource can be added throughout the project.

Extreme Scale Cloud/Fog Data processing:

- Hotspot analysis. Used for density/crowded plots of location based mobility information

Active & Federated Learning over Mobility Data (T4.2)

- Crowd density forecasting. Used for predicted crowdedness in the beach areas.
- Parking garage occupancy forecasting. Used for predicting the occupancy of the parking facilities near the beach.
- Active Learning & XAI for crowd/flow forecasting. Used for improved crowdedness prediction based on combined data sources.
- Active Learning (AL) for risk category classification. Used to predict the risk based on the scenario's which are used by the authorities.

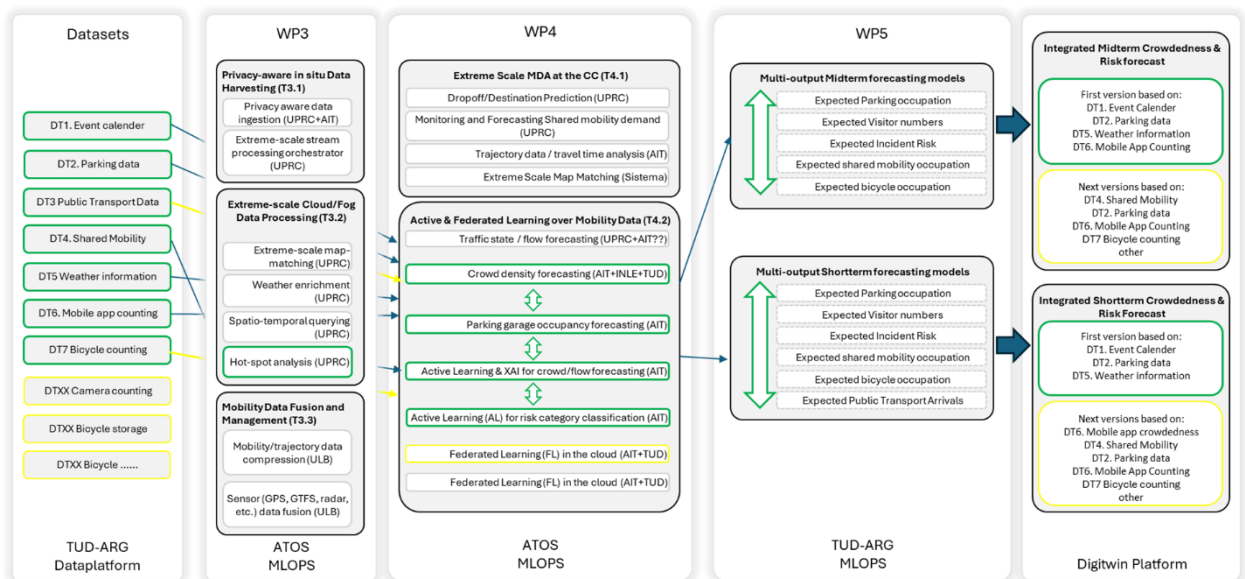


Figure 3-8 – Pipelines for use case 1.

### 3.6 Implementation Plan

---

Apart from the technical implementation of EMERALDS tooling, it is important to determine which information should be provided to whom, e.g., to avoid information overload and to provide information on a need-to-know basis only. Also, the way in which EMERALDS results are presented to users will be determined in close corporation among TUD, ARG, different departments for the city of The Hague and of the Law enforcement authorities. Developed EMERALDS' tooling will be incorporated (as described above) by ARG in the Scheveningen / The Hague Digitwin solution.

Apart from properly including EMERALDS' results into the Digitwin of Scheveningen, more mobility related aspects need to be taken into account in assessing their value and reported in D5.2 (M24), D5.3 (M33). For example, what is the societal impact of the improved predictions (and measures taken accordingly). Such as air quality, vehicle hours lost, overcrowding in other areas outside Scheveningen, safety, and legal aspects. These elements will be studied outside the course of this project.

Testing of EMERALDS' results for most objectives can be done throughout the year, since large (>1 year) historic data sets are available. However, if tactical and operational teams will be involved in the testing process and EMERALDS' results need to be compared to their everyday observations, this can only be done during April and October when Scheveningen beaches are most crowded.

- **Phase 1** UC Preparation, Alignment and Preliminary Data Collection **M1-M6**
- **Phase 2** Model development **M7-M15**
- **Phase 2** Model testing phase **M9-M15**
- **Phase 3** Implementation and first validation cycle **M12 – M24**
- **Phase 4** Preparation for second validation cycle **M24 -M30**
- **Phase 5** Implementation and second validation cycle **M30 – M33**

The assessment approach presented in Chapter 2 will be tailored to the requirements and implementation plan of the UC, within T5.1 in collaboration with T5.2, and results will be reported in two stages, documenting the findings of the first agile development sprint (M9-M24) on M24 (D5.2) and the refined developments of the second agile sprint (M25-M33) on M33 detailed in D5.3.

### 3.7 Key Stakeholders

---

The key stakeholders in UC1 are:

**Data providers:**

- NDW, provider of mobility data
- Resono, provider of mobile app counting
- RDW, provider of parking data
- CROW, provider of shared mobility data
- City of The Hague, provider of event agenda.
- Law Enforcement, provider of risk scenarios

**System Providers:**

- ARG, provider of data sharing platform and data visualisation platform Digitwin, TUD, developer of the system.
- 3rd parties (e.g., server & data hosting)

**System Users:**

- City of The Hague, responsible for safety on the beach.
- Law Enforcement, deciding on the need to intervene

**System Operators:**

- N/A

## 4 UC2: Multi-modal integrated Traffic Management

UC2 relates to T5.3. T5.3 aims to design a multi-modal traffic modelling and prediction module comprised of tightly coupled processing elements such as extreme scale data collection of big sensory data real-time and long-historical data. The implementation in T5.3 can be further broken down to the following subtasks:

- ST5.3.1 Self-identification of underperformance and self-tuning of traffic control algorithms, including multi-modal network decomposition
- ST5.3.2 Mobility AI-driven data fusion and predictive analytics methods for identification, assessment, and solvability of bottlenecks (the off-line state estimator and related functions)
- ST5.3.3 Traffic prediction model which learns from the training data set

### 4.1 Background – Context – Challenges

---

#### 4.1.1 *The value of data in multi-modal traffic management*

In recent years, the concept of active traffic management has evolved significantly, with a primary aim to enhance multi-modal traffic operations. This evolution aligns with prevailing policy objectives such as accessibility, safety, fairness, equity, sustainability, and liveability. By harnessing real-time and anticipated multi-modal traffic data, active traffic management dynamically intervenes to optimize traffic flow and ensure these policy objectives are met. This shift from prioritizing throughput and efficiency to embracing a broader set of policy goals has marked a significant transformation in traffic management approaches.

Particularly in The Netherlands, a frontrunner in network-wide multi-modal traffic management, there has been a focused effort to develop and implement multi-modal traffic management systems. The Dutch approach, known as "Multimodaal Netwerkkader" represents a systematic strategy to tailor multi-modal traffic management to meet area-specific objectives. This often involves prioritizing modes of transport such as bicycles, pedestrians, and public transportation over automobiles, reflecting a nuanced understanding of local traffic dynamics and policy goals.

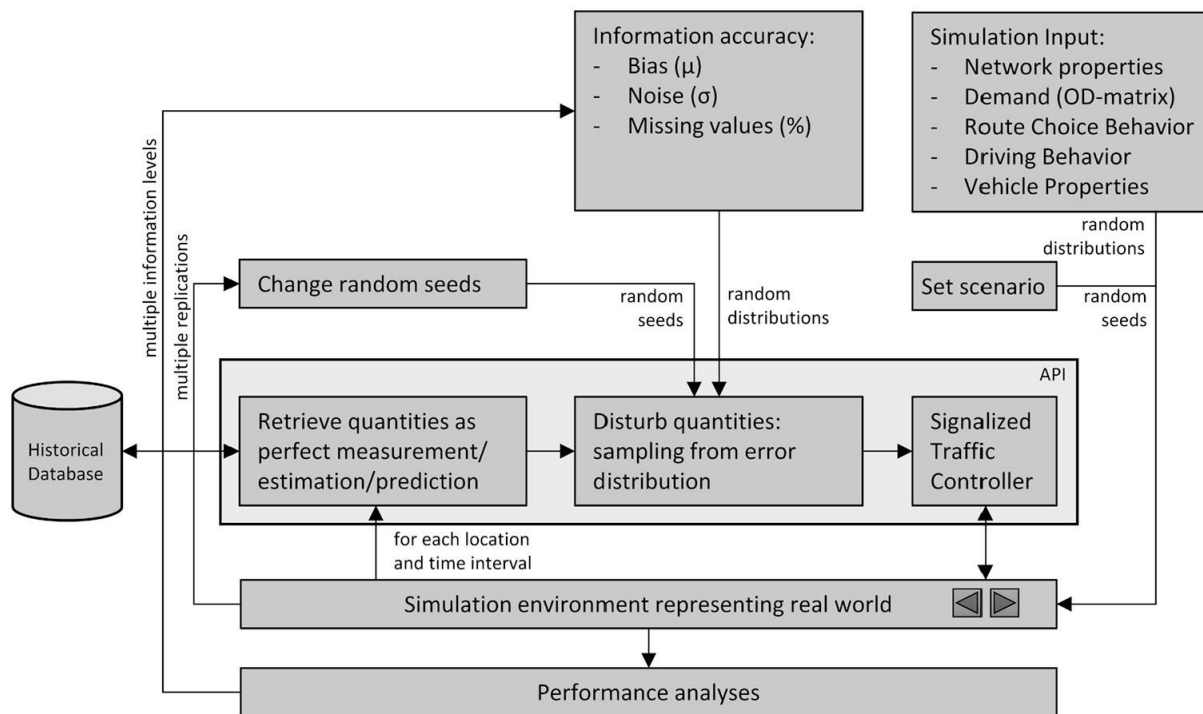
A critical aspect of this approach is the dual application of multi-modal traffic and transportation data. As discussed in section 4.2, use case 2 illustrates two main uses: the off-line use during the strategy design phase and real-time use during operational phases. These applications demand high-quality data—accurate, reliable, and timely—to effectively inform traffic management strategies. The impact of data quality on application performance is undeniable, yet, surprisingly, there has been limited research exploring the relationship between data quality and the effectiveness of multi-modal integrated traffic management applications.

The work of Klunder et al., particularly their studies in 2014 and 2017, has been instrumental in understanding how data quality affects traffic management measures such as ramp metering and dynamic routing. In 2014, Klunder and her colleagues explored the effectiveness of ramp metering, a traffic management strategy designed to mitigate congestion by controlling the flow of vehicles onto highways [27]. Their research underscored the pivotal role of data quality, demonstrating that accurate and timely traffic data are crucial for optimizing ramp metering algorithms to delay congestion onset and prevent capacity drops effectively.

Furthering this investigation, Klunder et al. [28] delved into the impact of data quality on dynamic routing applications, examining how Floating Car Data (FCD) penetration rates influence system performance. Their findings revealed that the improvement in total traffic delay through the use of in-vehicle dynamic routing systems is significantly dependent on the penetration rate and accuracy of FCD. For a 10% penetration rate, the delay improvement varied between 2.0% and 3.4%, translating

to yearly savings of approximately 15 million euros when delays are monetized using standard Value of Time (VOT) rates. This research not only highlighted the economic benefits of improving data accuracy but also illustrated the value of precise and reliable traffic data for enhancing traffic management strategies.

On the other hand, Poelman et al.'s work, particularly their 2020 and 2023 studies, provided critical insights into the quality of estimated and predicted input quantities for traffic control systems. In 2020, they proposed a general framework for sensitivity analysis to evaluate the impact of erroneous input quantities on the performance of signalized traffic controls [29], see Figure 4-1. This framework, applied to predictive control with varying adaptivity levels, demonstrated that accurate predictions significantly improve the performance of traffic controllers. Even in scenarios with incorrect data, controllers with higher adaptivity levels outperformed those with lower adaptivity, underscoring the importance of flexible and responsive control systems.



**Figure 4-1 – Experimental framework from Poelman et al (2020).**

Extending their research, Poelman et al. in 2023 investigated the impact of inaccurate predictions on model predictive control at traffic intersections [30]. They found that longer prediction horizons generally lead to better system performance, up to an optimal horizon length. High update frequencies were shown to mitigate the negative effects of prediction errors, enabling controllers to correct inaccuracies more swiftly. However, the studies also pointed to significant performance losses due to aggregation errors and biases in prediction models, suggesting substantial gains from more reliable predictions and the integration of detailed individual vehicle information into control strategies.

These studies collectively underscore the critical role of information quality in traffic management. Whether it concerns the estimation of current state variables or the accuracy of predictive models, the quality of information directly influences the effectiveness of traffic control strategies and algorithms. Despite the limited number of studies, the existing research highlights the profound impact of data quality on traffic management outcomes, advocating for more reliable predictions and the incorporation of detailed vehicle information in future applications.

#### 4.1.2 Challenges in data processing and analytics posed from extreme scale

In general, the quality of information critically depends on the characteristics – including abundance, scale, granularity, accuracy – of the data. In the past decades, the traffic and transportation management domain has transformed from a domain which is data poor, to a domain in which data is ubiquitous, due to developments in data collection, and increased connectivity. The magnitude, scale, and the heterogeneity of the data, however, result in new challenges in real-time and off-line data processing in multi-modal traffic management. The main challenges associated with managing and leveraging extreme data in the context of multi-modal traffic management include the following:

**Data Volume and Velocity:** The sheer volume of data generated by various sources can be overwhelming, requiring substantial storage and processing capabilities. Handling this data in real-time or near-real-time, considering its velocity, adds another layer of complexity. The infrastructure and algorithms must be scalable and efficient to process and analyse data streams effectively.

**Data Variety and Integration:** Extreme data comes in various formats and structures, from structured data like loop detector information to unstructured data from social media. Integrating these disparate data sources into a coherent system that can provide actionable insights for traffic management is a significant challenge. This requires sophisticated data integration tools and techniques to normalize, correlate, and aggregate different data types.

**Data Quality and Reliability:** Not all data sources provide accurate or reliable information. The quality of data can vary significantly, influenced by sensor errors, data transmission issues, or inaccuracies in user-generated content. Ensuring data quality and filtering out noise to make accurate traffic management decisions is a critical challenge.

**Privacy and Security:** With the increase in data sources, especially personal data from location-based services and social media, privacy and security concerns become paramount. Ensuring that traffic management systems comply with data protection regulations and safeguard user privacy while utilizing this data for traffic optimization is a complex challenge.

**Data Analysis and Interpretation:** The ability to extract meaningful insights from extreme data requires advanced analytical tools and expertise. Employing machine learning algorithms and artificial intelligence to predict traffic conditions, identify patterns, and recommend management strategies involves challenges in model development, training, and validation.

**Real-time Decision Making:** Translating data into real-time or near-real-time traffic management decisions requires not only fast data processing but also sophisticated decision-making frameworks. These systems must be capable of dynamically adjusting traffic signals, managing congestion, and recommending route changes to drivers instantaneously.

**Scalability and Flexibility:** Traffic management systems need to be scalable to accommodate growing data volumes and flexible enough to integrate new data sources as they become available. Developing systems that can evolve with the changing landscape of data sources and traffic management technologies is challenging.

**Policy and Regulatory Compliance:** Ensuring that traffic management practices derived from extreme data analysis comply with local and international regulations and policies is crucial. This includes adherence to traffic laws, environmental regulations, and urban planning guidelines, which can vary widely across different jurisdictions.

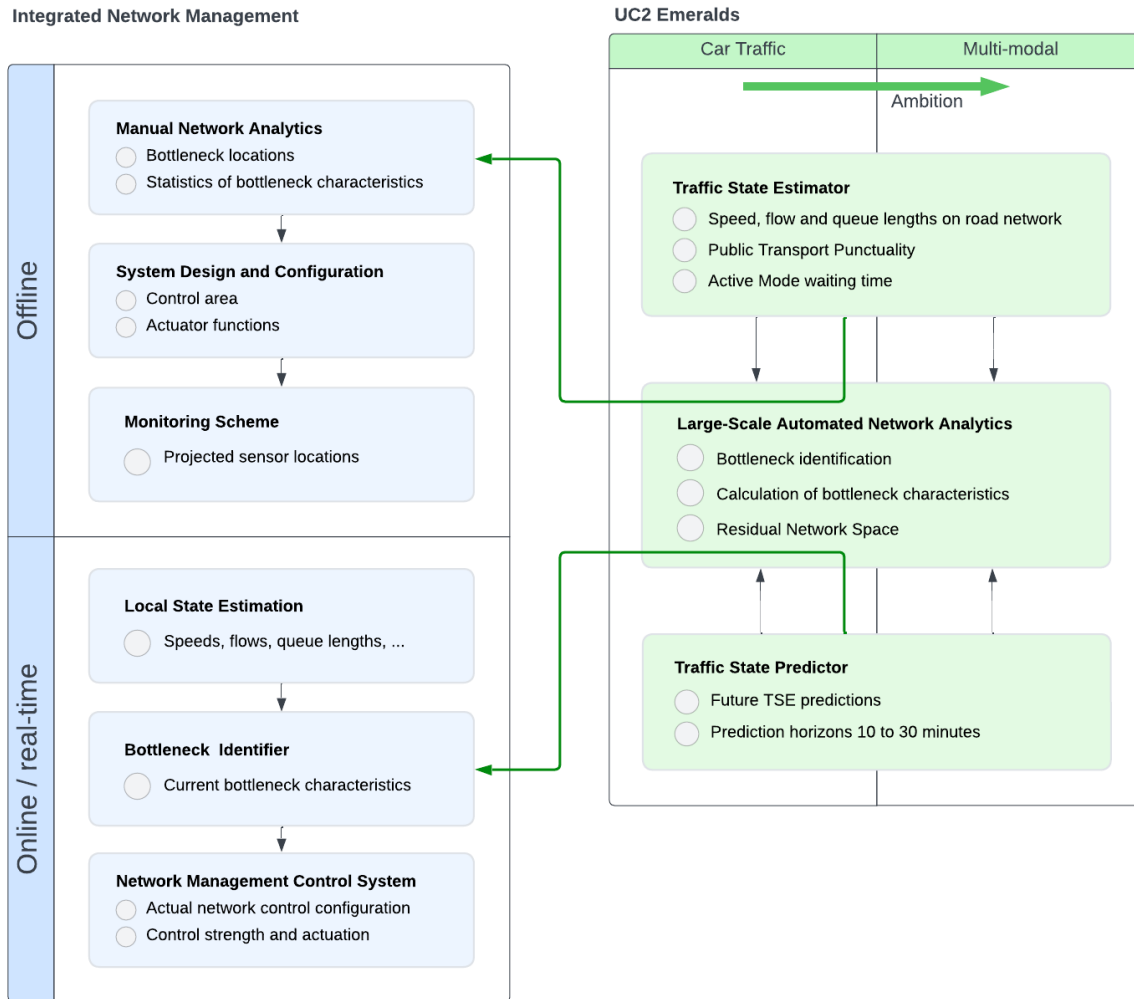
Addressing these challenges requires a multidisciplinary approach, combining expertise in data science, traffic engineering, information technology, and legal and ethical considerations. The successful management of extreme data for multi-modal traffic purposes can lead to significant improvements in traffic flow, safety, and sustainability, but it necessitates careful consideration of the complexities involved.

#### 4.1.3 Rotterdam use case in the context of EMERALDS

In Rotterdam, a pioneering approach to multi-modal traffic management has been initiated, with the Adaptive Flow Management Maastunnel serving as the starting point. This innovative strategy is part of a broader framework that seeks to harmonize traffic flow across various modes of transportation, aligning closely with the policy objectives of the city, in line with the “multi-modaal referentiekader” (Multi-Modal Network management framework). These objectives are meticulously integrated into Rotterdam's multi-modal network management framework, ensuring that every aspect of traffic control and flow supports the city's overarching goals for transport.

Also, for the Rotterdam multi-modal traffic management system, the process of integrated network management involves a variety of offline and online analytical techniques. **Offline**, manual network analytics are used to identify bottleneck locations and to study the statistical characteristics of these congested areas. This is focused on the analysis of the multi-modal network, where we focus on finding the bottlenecks (pinch points) where capacity is scarce, and at areas where capacity is available. These network analytics aim to support the design of the INM system [18]. The system design and configuration are meticulously planned, encompassing control areas and actuator functions. A detailed monitoring scheme also plays a crucial role, outlining where sensor locations should be projected for optimal data collection. In the **online**, real-time arena, local state estimation is performed to ascertain speeds, flows, queue lengths, etc., while a bottleneck identifier focuses on the current characteristics of traffic congestion. The Network Management Control System then utilizes this data to configure and adjust control measures in real time, ensuring a responsive and adaptive traffic management system.

Figure 4-2 furthermore shows the role of the EMERALDS in UC2. A Traffic State Estimator measures speed, flow, and queue lengths on the road network, punctuality of public transport, and active mode waiting time. This data feeds into the Large-Scale Automated Network Analytics, which identifies bottlenecks and calculates residual network space. Furthermore, the Traffic State Predictor is an essential tool that offers future Traffic State Estimations (TSE) predictions with prediction horizons ranging from 10 to 30 minutes. This predictive capability allows for proactive traffic management, adapting strategies before congestion becomes unmanageable.



**Figure 4-2 – Offline and online tool traffic management components**

Figure 4-3 introduces the fundamental process of network analytics, which is pivotal in understanding and managing traffic dynamics, and supports the design of multi-modal integrated traffic management strategies. The figure shows how the process begins with the policy that dictates the desired functional use of the network, and a reference framework that offers clear definitions of what constitutes a bottleneck. Network analytics then uses this framework to examine data across various scales, such as intersections, corridors, routes, and subnetworks, and across different modes of transportation, including cars, bicycles, public transport, and pedestrians.

The goals of network analytics are multi-faceted, focusing on policy monitoring, identifying the causes of bottlenecks (seeds), and examining the relationships between these bottlenecks over time and space. Additionally, the analysis considers the properties of these seeds, like severity and frequency, and seeks residual space within the network or other modalities to propose solutions. Real-time monitoring of bottlenecks is critical to the responsive and adaptive nature of the system.

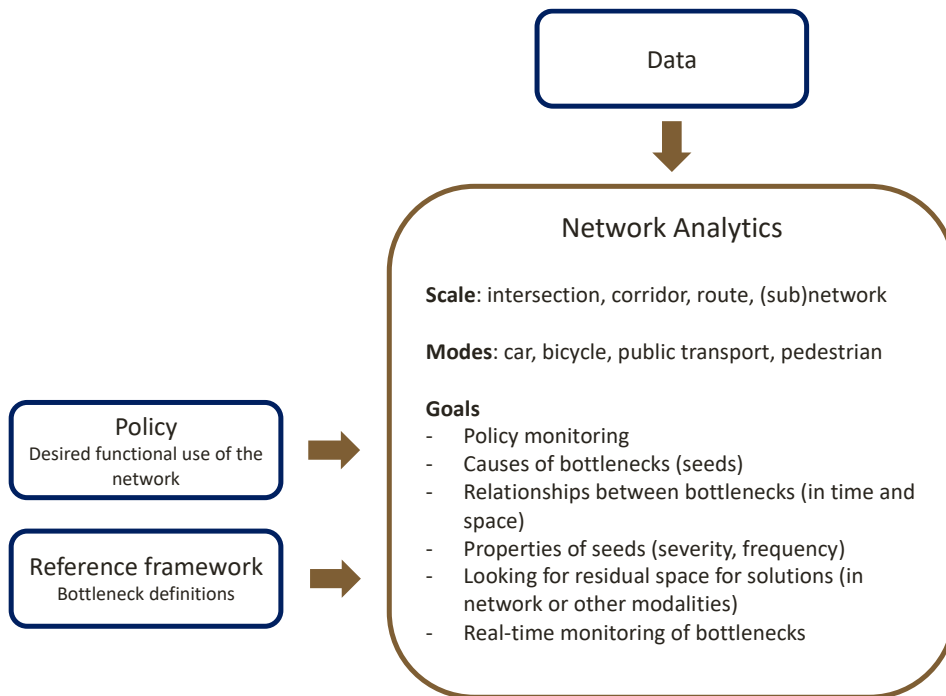


Figure 4-3 – Network analytics overview

In the remainder of this chapter, we will discuss how (some of) these challenges are addressed in UC2. Establishing several user stories, we will show how the EMERALDS will be pivotal additions to the existing multi-modal traffic management applications, both off-line and real-time.

## 4.2 SMART Definition and UC Objectives

For UC2, we have identified **six** objectives, distinguishing estimation, analytics and prediction for car traffic and multi-modal traffic respectively. The objectives are listed in the tables below.

### Objective 1: Traffic State Estimator for Car Traffic 1

Specific	<p>The car traffic state estimator (TSE) entails the estimation of speeds, flows, travel times, and densities of car traffic on motorways and urban roads. This requires a road network as a basis to project the estimates on. To achieve this, the following challenges have to be dealt with:</p> <ul style="list-style-type: none"> <li>• Construct a road network with upstream and downstream relations between links</li> <li>• Project measurements to the network (map-matching)</li> <li>• Calculate estimates for road segments that have no measurement. The geometry of the network (e.g., prohibited turns) and flow-consistency has to be maintained.</li> <li>• Performance on large-scale networks</li> </ul>
Measurable	<p>The TSE has to fulfil the following KPI's:</p> <ul style="list-style-type: none"> <li>• Scalability: the TSE has to be able to process all measurements for the full network of the Rotterdam metropolitan area</li> </ul>



	<ul style="list-style-type: none"> <li>• Completeness: the TSE has to provide estimates for each link in the network</li> <li>• Validity (quality): Since the measurements are the only ground truth, we run the TSE with reduced measurements and determine the quality of the estimate. No direct KPI can attached to this, since the TSE on its own is not valuable. Therefore, an impact assessment of the “network analytics” objective 3 is performed. What are the differences of their outcomes when measurements are omitted?</li> </ul>
Action - oriented	<p>The following actions are taken to establish the TSE:</p> <ul style="list-style-type: none"> <li>• Building the network</li> <li>• Map match measurement locations to the network. We include all loop detector sensors and floating car data available in the Netherlands.</li> <li>• Determine rules for filling in blanks based on traffic flow theory. This is an iterative process where expertise from big mobility data science and traffic engineering are required. In each iteration we aim to improve the KPIs</li> </ul>
Relevant	The TSE is crucial to allow network analytics and is therefore input for objective 3 /ST 5.3.2.
Timebound	As one of the building blocks for UC2, we project the system to be operational by M19 for testing.

### Objective 2: Multi-Modal Traffic State Estimator

Specific	<p>The multi-modal state estimator (MMTSE) entails determining the required multi-modal traffic and transportation variables that are needed for the design and operations of integrated multi-modal traffic management [31]. These include, in addition to objective 1:</p> <ul style="list-style-type: none"> <li>• Speeds, flows, travel times, delays, and stops for bicycle traffic</li> <li>• Location, speeds, travel time and delays, and stops of public transport vehicles</li> </ul> <p>By adding modes, new challenges occur:</p> <ol style="list-style-type: none"> <li>A) Networks for each mode need to be available and the interfaces between modes are crucial</li> <li>B) Data for public transport and bikes have a different origin than car-data. These need to be understood to make the correct estimates.</li> </ol>
Measurable	Identical to objective 1: Number of modes studied, data granularity and temporal scope.
Action - oriented	<p>For the development of the MMTSE we are taking an incremental approach. In each increment we add a data source. These steps are considered within each iteration:</p> <ul style="list-style-type: none"> <li>• Exploratory Data Analysis: what are the spatial-temporal characteristics, including coverage/sample size, frequency and latency</li> </ul>

	<ul style="list-style-type: none"> <li>Adjusting the MMTSE to allow the data to be added</li> <li>Map-match the new data source</li> <li>Estimate missing values in time and space, in particular estimate queue length on urban corridors.</li> </ul> <p>Possible datasets for which iterations will be made are:</p> <ul style="list-style-type: none"> <li>Talking bike data (partial bike trajectories)</li> <li>Bike counts by loop detectors (NDW)</li> <li>Public Transport location data (NDOV)</li> </ul>
Relevant	The MMTSE is crucial to allow multi-modal network analytics and is therefore input for objective 4.
Timebound	As one of the building blocks for UC2, we project the system to be operational by M19 for testing.

### Objective 3: Large-scale Network Analytics for Car Traffic

Specific	<p>The network analytics for car traffic comprises the labelling and rating of the traffic state (which can either be an estimate or a prediction). Therefore, the time-space of speeds, flows and queues is processed. The first step is to identify where and when congestion occurs or when indicators exceed policy norms.</p> <ul style="list-style-type: none"> <li>Highway congestion is determined by creating the fundamental diagram of traffic on each location and assessing when traffic is outside the free flow branch.</li> <li>Urban roads queues are labelled problematic when they block upstream intersection (for a traffic flow perspective). Also, when flows, speeds or queue length exceed policy norms (e.g., related to noise, air quality, liveability, crossability for active modes, hindrance for PT) urban roads are labelled problematic.</li> </ul> <p>The labels are created locally. The next step is to aggregate the labels in time-space to determine bottleneck characteristics.</p> <ul style="list-style-type: none"> <li>For highways clusters of congestion labels are made. Based on the size and properties of these clusters a distinction is made between stationary congestion and shockwaves. Per cluster the size in time-space and vehicle loss hours are retrieved.</li> <li>For urban networks labelled roads are clustered over upstream and downstream intersection. This provides 'jam trees' of labelled roads. Also, for these trees characteristics are calculated such as depth, loss hours, duration and possibly policy norm specific measures.</li> </ul> <p>Ideally, we determine the cause for each of the bottleneck. That can either be infrastructure, interaction of (multi-modal) flows at intersections. The bottlenecks and their characteristics are then used as monitoring part of INM.</p> <p>Finally, the bottlenecks are determined for a longer period of time. We derive frequencies and statistics of characteristics (averages, outliers). This information helps designing and configuration of INM system.</p>
----------	---

Measurable	<p>Scalability: the analysis has to be performed on the complete network of the traffic state estimate or prediction</p> <p>Validity (precision and recall): manual labelling and clustering vs. automated. Then determine false/true positive/negative for several indicators. The precision (<math>TP/(TP+FP)</math>) and recall (<math>TP/(TP+FN)</math>) have to be at least 95%.</p>
Action - oriented	<p>The following steps are undertaken:</p> <ol style="list-style-type: none"> <li>1) Perform labelling</li> <li>2) Develop and Perform clustering</li> <li>3) Determine characteristics of bottlenecks</li> <li>4) Calculate statistics based on bottleneck characteristics.</li> </ol>
Relevant	<p>The tools developed under this objective provide crucial input to the INM system design process. It enables the identification of bottlenecks in the traffic network, the causes of these bottlenecks, but also shows where control space is available to resolve these bottlenecks. This 'space' can stem from other locations in the network, but also to other traffic modes (relation to objective 4).</p>
Timebound	<p>As one of the building blocks for UC2, we project the system to be operational by M19 for testing.</p>

#### Objective 4: Large-scale Multi-modal Network analytics

Specific	<p><b>Network-wide multi-modal data analytics.</b> This includes all tools that allow judging the performance of a multi-modal network, identify the main pinch points, computation and visualisation of the main performance indicators, etc., with the aim to improve the design of the multi-modal integrated traffic management approach.</p> <p>Specifically, the analytics will involve the performance of the active modes (focus on bikes) and PT (buses and trams).</p> <p>For <b>bicycle traffic</b>, traffic performance (or level of service) is related firstly to the number of stops on the route, the route delay, need to detour, and secondly to flow rates (crowdedness). These indicators need to be computed from (the combination of) the available data sources (incomplete FCD data for bikes, counting data insofar available, detection at intersections) [32]. Performance is to be determined both at the (sub)route level and at bottlenecks (intersections, crossings). In summary, the aim is to estimate and analyse:</p> <ul style="list-style-type: none"> <li>• Number of stops on a given set of routes (e.g., based on the preferred routes stemming from Multimodal Network Frameworks)</li> <li>• Route delay for a given set of routes</li> <li>• Average detour factor (determined for GPS of bikes where chosen sub routes are compared to shortest distance routes)</li> <li>• Delays at intersections (GPS data of bikes, possibly combined with intersection control data insofar historically available)</li> </ul>
----------	---

	<ul style="list-style-type: none"> <li>• Link bicycle flow rates</li> </ul> <p>For <b>public transport</b>, performance is determined by regularity and punctuality (depending on the service frequency). Punctuality relates that the service operations according to the planned schedule and can be determined by identifying delays at the bus stops, which in turn can be determined from the position data of the buses. Regularity is only relevant for high frequency services and is less relevant for UC2. Punctuality is influenced by delays incurred along the route of the vehicle, e.g., due to oversaturation at intersection control, and due to long boarding times at stops. In sum, we consider the following PT performance indicators:</p> <ul style="list-style-type: none"> <li>• Punctuality at stops</li> <li>• Boarding times at stops</li> <li>• Delays at intersections</li> <li>• Estimation of demand profiles (from boarding times)</li> </ul> <p>The indicators are computed and visualised (i.e. via the DigiTwin platform).</p>
Measurable	<p>For both the bicycle analytics and the PT analytics individual vehicle data is available. The bicycle data is historic and comprised 3 years of data. The data are potentially augmented with intersection control data (insofar historically available).</p> <p>For the PT operations, NDOV data are available. These comprise position data of all PT vehicles at a specific frequency based on GPS.</p> <p>For all indicators presented, accuracy and reliability are important:</p> <p><b>Accuracy:</b> the (relative) error in the value of the indicator against the real value as a function of time (each minute, each hour, per day); note that with absence of a ground truth, we need to resort to proxies to determine the accuracy (e.g., change in accuracy when leaving out data or adding noise; use of synthetic data from simulation).</p> <p><b>Reliability:</b> the percentage of time that error in the computed indicator is larger than p%.</p> <p>Next to the indicators, the multi-modal network analytics toolbox aims to identify the underlying causes from the bottlenecks determined by comparing the indicators to the Multimodal Network Framework performance reference maps.</p>
Action - oriented	<p>We see the following activities:</p> <ul style="list-style-type: none"> <li>• Specify and compute the different indicators described above from the enhanced data (from objective 2)</li> <li>• Determine pinch points / bottlenecks (e.g., by comparison with reference map from the multimodal Network Framework and combination with car bottleneck analysis (objective 3) [multi-modal bottleneck inspector])</li> <li>• Develop tools for detailed problem analysis / bottleneck causes, e.g., by looking for spatial correlations in the multi-modal network operations</li> <li>• Develop visualisation tools</li> </ul>

	<ul style="list-style-type: none"> <li>• Functional assessment</li> </ul>
Relevant	Accurate performance indicators for active mode and public transportation traffic operations are key to improve the multi-modal INM system for UC2; they are used to assess, improve and finetune the system.
Timebound	As one of the building blocks for UC2, we project the system to be operational by M19 for testing.

### Objective 5: Traffic State Prediction for Car Traffic

Specific	<p><b>The car traffic state predictor (TSP)</b> determines (short-term) estimated for the traffic conditions in the network. The conditions involve speeds, flows, route travel times, densities, queue lengths and link / intersection delays for the motorways and / or urban roads. The predictor uses the output of the TSE (objective 1) as the main input.</p> <p>The TSP will predict traffic states for different time horizons (e.g., 5, 15, 30, 60 minutes ahead). The TSP needs to provide network wide conditions and should be scalable to other, larger networks.</p> <p>One of the key challenges is that the predicted conditions are dependent on the traffic management interventions taken (i.e., if the INM system acts on the predicted occurrence of a bottleneck, the bottleneck may actually not occur).</p>
Measurable	<p>Similar to UC1, the quality of the predictions is assessed using accuracy and reliability measures. These are computed for all relevant variables (density, speed, flow, delay per link; queue lengths, delays at intersections) for the specified prediction times (5, 15, 30, 60 minutes ahead).</p> <p><b>Accuracy:</b> the (relative) error in the value of the variable for the <math>T</math> minute ahead prediction compared to the estimated value of this variable at <math>t+T</math> (e.g., predicted 15-minute ahead link speed <math>v^{15}(t)</math> is compared to the estimated link speed <math>v(t+15)</math>).</p> <p><b>Reliability:</b> the percentage of time that error in the computed variable is larger than <math>p\%</math>.</p> <p>The developed Emeralds will be compared against benchmark prediction approaches.</p>
Action - oriented	<p>For reaching the objectives, the following activities are foreseen:</p> <ul style="list-style-type: none"> <li>• Data preparation for analysis and training purposes</li> <li>• Statistical analysis of correlations in data / state estimates (direction of information, correlation of links in network for different prediction horizons; influence of external variables, etc.)</li> <li>• Development of Emeralds prediction algorithms / implementation of benchmark prediction approach (hybrid AI method, e.g., Graph-based NN)</li> <li>• Functional assessment</li> </ul>
Relevant	If successful, the predictions can be used both for predictive traffic management and for route information provision to the end users.

Timebound	The objective depends on reaching of objective 1; completion in M30 for testing.
-----------	--

### Objective 6: Multi-modal Traffic State Prediction

Specific	<p><b>Multi-modal prediction and forecasting.</b> Predictions (short-term, up to 12 hours) and forecasts (med-term or long-term) are needed to improve traffic management functions, including the provision of information to travellers.</p> <p>For multi-modal prediction, we distinguish functions for the active modes and for public transport. For the active mode, we consider the following objectives:</p> <ul style="list-style-type: none"> <li>- Route bike travel time prediction, including delays at intersections and crossings [33]</li> <li>- Bicycle traffic demand forecasting based on time of day, day of week, weather forecast, holidays and events</li> </ul> <p>For public transport (buses), we have the following objectives:</p> <ul style="list-style-type: none"> <li>- Route travel time prediction (real-time) for improving regularity and multi-modal traffic management</li> <li>- Passenger demand prediction at bus stops</li> </ul> <p>To which extent we achieve those goals will depend on the actual data availability. See objective 2.</p>
Measurable	<p>Similar to UC1, the quality of the predictions and forecast are assessed using accuracy and reliability measures. These are computed for all relevant variables (travel time for bikes and buses, bike demand, passenger demand at stops). For predictions, we look at the predicted travel times (e.g., the travel time that will be realised from the moment of departure); for the predictions of demand at a bus stop, we use the predicted value of the number of passengers waiting at the planned bus arrival time; the bicycle demand forecast will be done <math>d</math> days ahead).</p> <p><b>Accuracy:</b> similar to objective 5, the accuracy will be determined by comparing the predictions and the forecast to the (shifted) estimates (objective 2).</p> <p><b>Reliability:</b> the percentage of time that error in the computed variable is larger than <math>p\%</math>.</p> <p>The developed Emeralds will be compared against benchmark prediction and forecasting approaches.</p>
Action - oriented	<p>For reaching the objectives, the following activities are foreseen:</p> <ol style="list-style-type: none"> <li>1) Data preparation for analysis and training purposes</li> <li>2) Statistical analysis of correlations in data / state estimates</li> <li>3) Development of Emeralds prediction algorithms / implementation of benchmark prediction approach</li> </ol>

	4) Functional assessment of Emeralds against ground truth / benchmark applications
Relevant	<p>The real-time predictions of bike route travel times can be used for real-time information for cyclist, allowing them to choose their preferred route through the city thereby improving city bike-ability. The real-time predictions of bus route travel times can be used to improve the multi-modal integrated traffic management functioning, by making priorities for buses conditional on the predicted regularity rather than on the observed regularity. Real-time predictions of number of boarding passengers at stops can be used for real-time dwell time analysis, which in turn can further improve conditional priority interventions for the buses.</p> <p>The <i>forecasts</i> of demand are mostly used for planning purposes: they enable adaptations of signal plans in case high (or low) demands are expected for bikes, or delay forecasts of PT due to long dwell times.</p>
Timebound	The objective depends on reaching of objective 1; completion in M30 for testing.

For UC2, User stories describe processes, workflows and performance of traffic management systems related to the tasks of the traffic managers in the city of Rotterdam. Every user story has scenarios that each describe a different implementation of the user stories. We have identified the following user stories:

#### 4.2.1 *User-story 1: using off-line data for design of multi-modal measures*

This user-story deals with the task of policy makers of the design and implementation of multi-modal traffic management measures. Usually, this process follows four typical steps. First, policymakers determine the policy objectives of the Rotterdam transport system (step 1). These objectives are then translated into a Multi-Modal Network management Framework, which includes a reference map (step 2). An accurate picture of the existing multi-modal traffic conditions is then compared against the reference framework, which shows the locations of the key bottlenecks. This results in an assessment of the networks' performance (step 3). Based on the network performance, a roadmap for measures will be developed. Such measures include Traffic Management, Mobility Management, and Infrastructural measures (step 4). User-story 1 focuses on step 3 specifically. It primarily concerns determining properties and interactions of bottlenecks. This consists of:

- Automatic identification of bottleneck locations (UC objectives 1, 2, 3, 4)
- Calculating relevant statistics of the bottleneck performance (UC objectives 3, 4)
- Automatically finding interactions and causes of, for example, incidents (UC objectives 3, 4)

User-story 1 has two new scenarios:

- Scenario 1-0 (base scenario): current practice with ad-hoc data analyses of datasets
- Scenario 1-1: multi-modal local intersection analysis (assessment with reference framework) (UC objectives 1, 2, 3, 4)
- Scenario 1-2: public transport assessment at different scales (network-wide assessment of PT) (UC objectives 2, 4)

#### 4.2.2 *User-story 2: Integrated network management (automated traffic control)*

Integrated network management systems use accurate data on queues (length, number of vehicles in queues) to determine issues (which problem needs to be solved) and actions (where do we have space in the network for temporary storage of traffic to solve the problem) in real-time. An example application is the prevention of spillback of queues into the Maastunnel. Here, control actions limit the inflow into the tunnel and increase the outflow from the tunnel by changing the signal settings of various controllers in the network. Improving monitoring functions can enable Integrated Network Management to anticipate on future bottlenecks, prioritise bicycles (by adding the performance of bicycle network) and/or public transport (by using actual line punctuality). User-story 2 focuses on:

- Anticipating on future bottlenecks, by using predicted traffic states (UC objectives 3, 4, 5, 6)
- Finding (cost-effective) alternatives for expensive radar-based queue length estimators, by using data infusion based on loop detector data and aggregate floating car data (UC objective 1)

User-story 2 has three new scenarios:

- Scenario 2-0 (base scenario): current system with radar-based queue estimators
- Scenario 2-1: queue length estimation based on intersection data (vlog) and floating car data (cost reduction by replacing radar) (UC objective 1)
- Scenario 2-2: queue length prediction based on intersection data (vlog) and floating car data (prediction allows anticipation) (UC objective 5)
- Scenario 2-3: 5-10 minute ahead prediction of bottleneck activation on the main road network (preventing or postponing capacity drop) (UC objectives 1, 3, 5)

#### 4.2.3 *User-Story 3: network operations*

This user-story concerns the work of traffic management operators. In the traffic control centre, traffic management operators control the network by activating measures, e.g., they re-route traffic with dynamic message signs. Based on the current traffic performance, they also make decisions to activate traffic management scenarios. The current traffic state includes incidents, congestion, etc. User-story 3 contributes to:

- Making decisions based on the predicted performance (UC objectives 3, 4, 5, 6)
- Making the paradigm shift from uni-modal to multi-modal traffic control (UC objectives 3, 4, 5, 6)

User-story 3 has three scenarios:

- Scenario 3-0 (base scenario): current operation with live data feeds
- Scenario 3-1: 15-60 minute ahead prediction of the traffic state on the main road network (preventing or postponing capacity drop) (UC objective 5)
- Scenario 3-2: 15-60 minute ahead prediction of the traffic state on the urban road network (anticipation with measures (intersection controllers and variable message signs) on future traffic states) (UC objectives 5, 6)



## 4.3 Computational and Data Resources

In this section, the key resources that enable the processing, analysis, and visualization of data, as well as the computational capabilities necessary to run the algorithms and models are highlighted. These resources are instrumental in harnessing the power of data analytics and supporting decision-making processes.

We use the EMERALDS infrastructure as described in D2.1 for the implementation of UC2. In addition, offline analyses are performed by the TU Delft and Arane on their own infrastructure. For UC2 no additional computational resources are required.

Furthermore, the importance of data resources is emphasized, including both public and proprietary datasets, that serve as inputs for the use case. These datasets provide valuable information for understanding the multi-modal transportation system, traffic patterns, and various contextual factors that influence network-wide traffic management.

The availability and quality of computational and data resources directly impact the accuracy, effectiveness, and efficiency of the use case outcomes. Therefore, careful consideration and planning are given to ensure that the necessary resources are accessible, scalable, and capable of handling the computational and data processing requirements of the use case. Instilling interoperability with the Emeralds services and laying the groundwork for their experimentation, testing and validation.

### 4.3.1 Data Resources

Table 3 lists the available data sources needed for the execution of UC2. They are also linked to the relevant scenarios and UC objectives. The examined data sources relevant to the data analytics problem are indicated while specific data attributes will be elicited in D5.4. More detailed description of the datasets is available in D1.4 Data Management Plan.

**Table 3 – Description of available data sources and their connection to the scenarios and use case objectives.**

Code	Name of the dataset	Description	Scenario	UC objectives
<b>DT1</b>	Public transport data	Information on public transport schedules, timetables and real-time locations of buses, trains and trams.	S1-1, S1-2, S3-2	2, 4, 6
<b>DT2</b>	Bicycle counting Netherlands	Bicycle counts are available at several locations in the Rotterdam network. At these points, the number of bike passages per time unit are provided; see Figure 4-4.	S1-1, S3-2	2, 4, 6
<b>DT3</b>	Bridge openings	Information on bridge openings	S3-2	1, 2, 5, 6
<b>DT4</b>	Incidents	Location of incidents	S2-3, S3-1, S3-2	3, 4
<b>DT5</b>	Roadworks	Location of roadworks	S2-3, S3-1, S3-2	5, 6
<b>DT6</b>	Floating Car Data	Travel time data based on floating car data collected from a smart phone app	S1-1, S2-1, S2-2, S2-3, S3-1, S3-2	1, 5

<b>DT7</b>	Loop Detector Data	Speed and flow data from double loop detectors in the network	S1-1, S2-1, S2-2, S2-3, S3-1, S3-2	1, 5
<b>DT8</b>	Queue lengths from radar	Arane has access to the radar-measured queue length data from the City of Rotterdam. It contains queue length and buffer space data at 10 s intervals for selected roads around the Maastunnel	S2-1	1

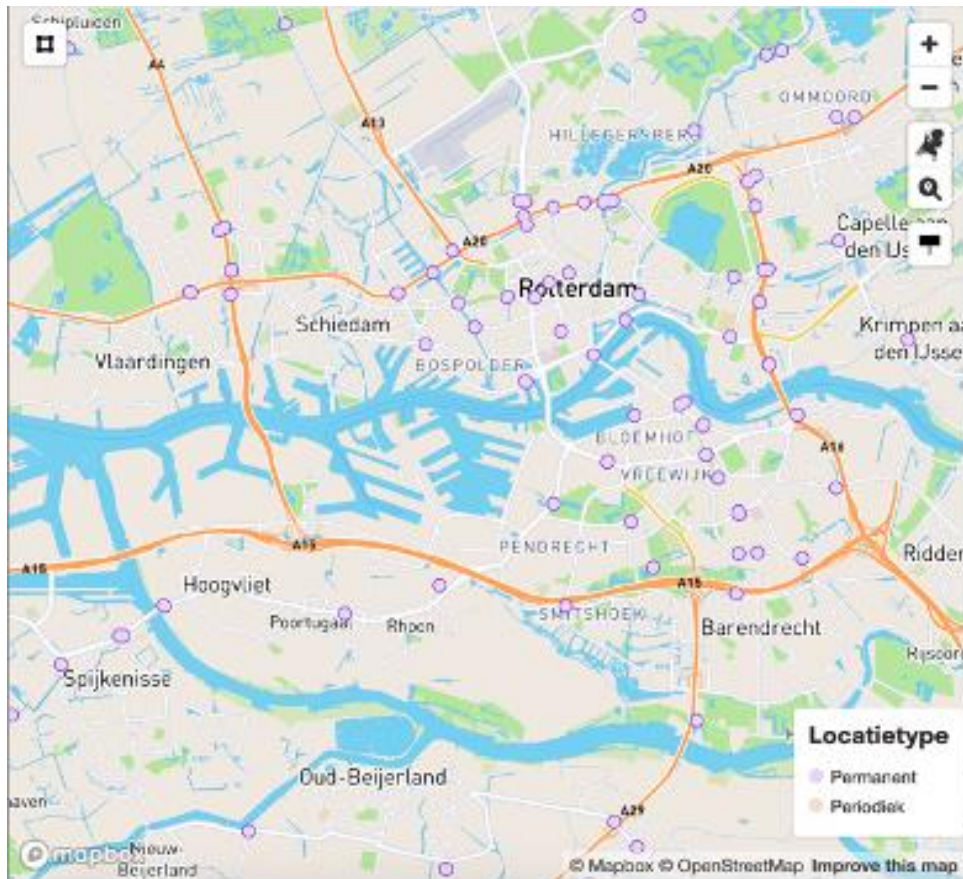


Figure 4-4 – Example of bicycle count data map.

#### 4.4 Assessment and corresponding KPIs

As introduced in chapter 2, we make a distinction between performance and business KPIs. UC2 is about estimation and prediction of various variables needed for Multi-Modal Integrated Network Management. For the estimates and predictions, the main performance KPIs used are as follows (see also tables in section 4.2):

1. RMSE (Root Mean Square Error) - Measures prediction accuracy, highlighting average error magnitude.
2. MAE (Mean Absolute Error) - Measures prediction accuracy, highlighting average error magnitude.
3. MAPE (Mean Absolute Percentage Error) - Measures prediction accuracy, highlighting average error magnitude.
4. Reliability - the percentage of time that error in the computed variable is larger than p%.

5. Availability - percentage of time the estimate / prediction is available.
6. Latency - the (extra) time needed from data collection to publication of the outcome of the Emeralds computation.
7. Scalability - the analysis has to be performed on the complete network of the traffic state estimate or prediction.
8. Completeness - the TSE has to provide estimates for each link in the network.
9. Validity (precision & recall) - Evaluates manual versus automated labelling and clustering accuracy, aiming for at least 95% precision (TP/(TP+FP)) and recall (TP/(TP+FN)).
10. Validity (quality) - Assesses if an estimate's quality is sufficient by comparing outcomes with and without certain measurements.

Table 4 below gives an overview of the different KPIs that are relevant for the various scenarios.

**Table 4 – KPIs for the different scenarios.**

KPI vs. Scenario	1-1	1-2	2-1	2-2	2-3	3-1	3-2
RMSE			X	X	X	X	X
MAE			X	X	X	X	X
MAPE			X	X	X	X	X
Reliability			X	X	X	X	X
Availability	X	X	X	X	X	X	X
Latency	X	X	X	X	X	X	X
Scalability	X	X	X	X	X	X	X
Completeness	X	X	X	X	X	X	X
Validity (precision & recall)			X	X	X	X	X
Validity (quality)	X	X	X	.X	X	.	.

The business KPI's follow directly from the scenario objectives. Below we list how the business KPI of each scenario is quantified:

- Scenario 1-1 (multi-modal local intersection analysis) and scenario 1-2 (public transport assessment at different scales): These scenarios extend the value chain of the City of Rotterdam's decision-making process by allowing better design of traffic management measures and network optimization. Based on interviews with the City of Rotterdam a valuation will be executed.
- Scenario 2-1 (queue length estimation): This scenario relates to the INM business case. By replacing radar with readily available data sources, a cost reduction will be achieved. This will be quantified based on the current operating costs. Eventually, this allows the city to implement INM on their complete road network more easily. We quantify the KPI by assessing the updated cost-benefit ratio of INM.

- Scenario 2-2 (queue length prediction) and Scenario 2-3 (prediction of bottleneck activation on the main road network): These scenarios aim to improve the efficiency of INM by incorporating predictions. Higher anticipation levels to future bottlenecks can be achieved. We quantify this in terms of monetized travel time savings of anticipated INM on a toy network.
- Scenario 3-1 (prediction of the main road traffic state) and scenario 3-2 (prediction of the traffic state on the urban road network): The scenarios improve the process of (manual) traffic control by operators. We assess this by interviewing traffic control operators about expected improved operations. We then quantify this by estimating the monetization potential of expected travel time savings.

## 4.5 Mapping to EMERALDS Toolset

In this section, we describe the current approach to multi-modal integrated traffic management in Rotterdam in relation to the novel data services EMERALDS can provide. We present the pipelines designed to execute the scenarios. We designed three pipelines covering all the scenarios, and describing the inputs, outputs, and intermediate processes. In the pipelines, the data requirements for the input and the relevant EMERALDS are also depicted. Figures 4.5, 4.6, and 4.7 present the three pipelines respectively. Table 4.3 presents the EMERALDS and their coding used in the pipelines.

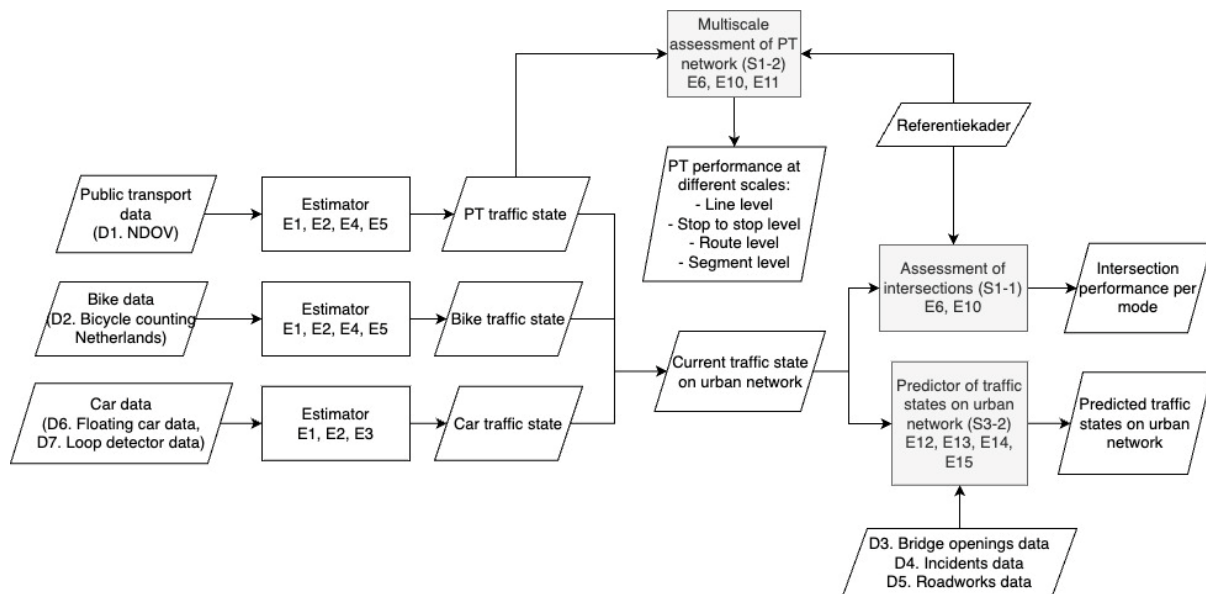


Figure 4-5 – Pipeline 1, covering Scenarios 1-1, 1-2, and 3-2. “E” indicates the EMERALD codes found in Table 5.

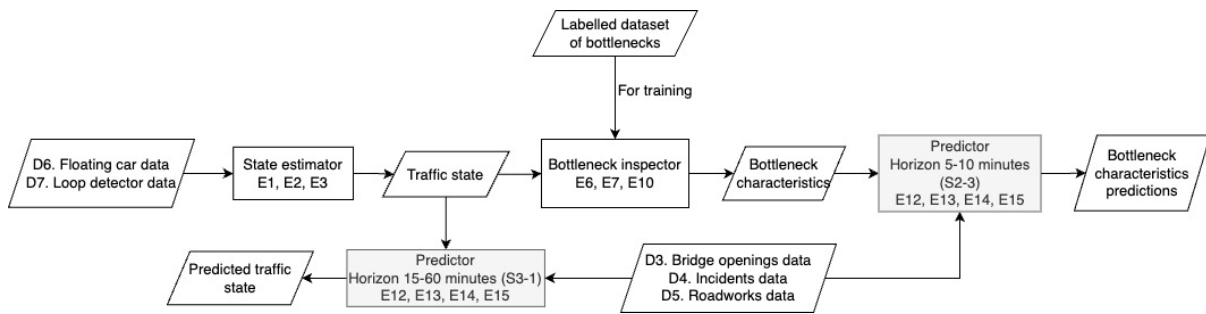


Figure 4-6 – Pipeline 2, covering Scenarios 2-3, and 3-1. “E” indicates the EMERALD codes found in Table 5.

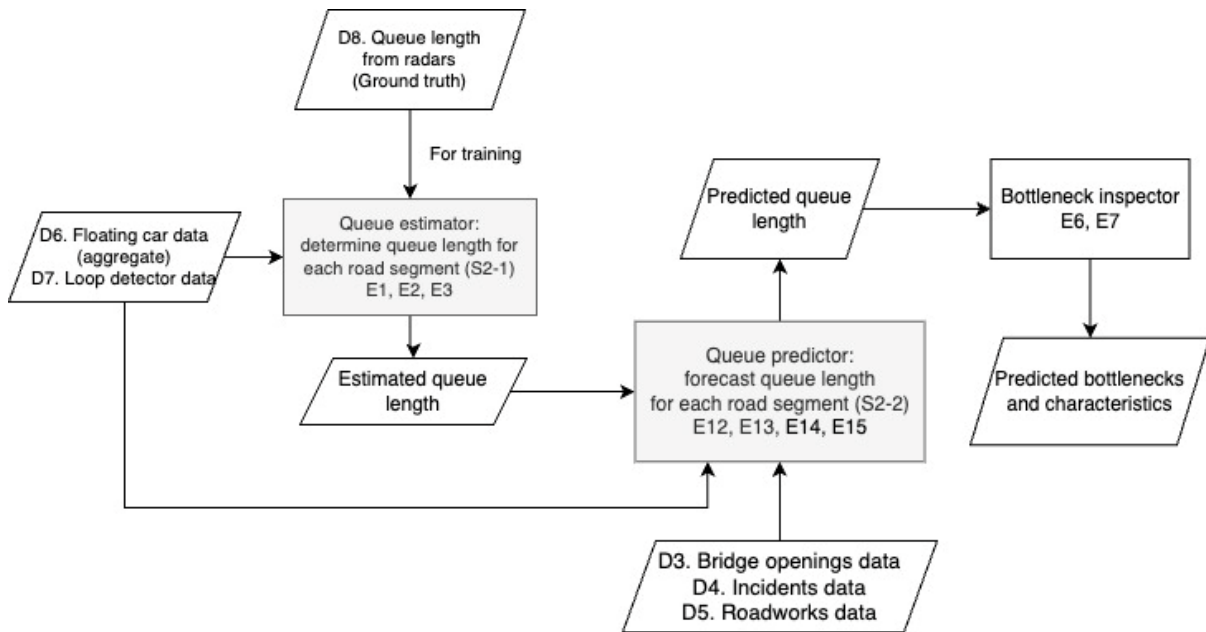


Figure 4-7 – Pipeline 3, covering Scenarios 2-1, and 2-2. “E” indicates the EMERALD codes found in Table 5.

Table 5 – EMERALDS and their codes referred to in the Pipelines.

Code	EMERALDS component
E1	Extreme-scale stream processing orchestrator (UPRC) (T3.1)
E2	Sensor (GPS, GTFS, radar, etc.) data fusion (ULB) (T3.3)
E3	Traffic State Estimation (T4.1)
E4	Extreme-scale map-matching (UPRC) (T3.2)
E5	Multi-modal Traffic State Estimation (T4.1)
E6	Spatio-temporal querying (UPRC) (T3.2), Hot-spot analysis (UPRC) (T3.2)
E7	(Multi-modal) Bottleneck Inspector (T4.1)
E8	Connected vehicles and Infrastructure data analytics (T4.1)
E9	Incident prediction, Trip habits and route forecasting by indicating origin and destination (T4.1)
E10	Multi-modal network-wide traffic analytics (T4.1)
E11	Trajectory Data / Travel Time Analysis (AIT) (T4.1)
E12	Weather enrichment (UPRC) (T3.2)
E13	Multi-modal Traffic State Prediction (T3.3)

<b>E14</b>	<i>Explainable AI for flow forecasting (T4.1)</i>
<b>E15</b>	Traffic State / Flow Forecasting (UPRC) (T4.2)

The three pipelines have different data requirements. In Table 6, we map the pipelines to the data requirements.

**Table 6 – Mapping of the pipelines to the data and specific data requirements.**

<b>Pipeline</b>	<b>Data requirements</b>	<b>Data specification</b>
<b>Pipeline 1</b>	DT1. Public transport data	Stop to stop and trajectory data for the whole city of Rotterdam. Time period: 1 year
	DT2. Bicycle counting Netherlands	We make a selection of 4 intersections for S1-2 based on bicycle count availability. Time period: 1 year
	DT3. Bridge openings data	Whole city network. Time period: 1 year
	DT4. Incidents data	Whole city network. Time period: 1 year
	DT5. Roadworks data	Whole city network. Time period: 1 year
	DT6. Floating Car data	Around selected intersections for S1-2. Different spatial aggregations. Time period: 1 year
	DT7. Loop Detector data	For selected intersections. Period: 1 year.
<b>Pipeline 2</b>	DT4. Incidents data	Period: 1 year
	DT5. Roadworks data	Period: 1 year
	DT6. Floating Car data	A subnetwork of the ringroad of Rotterdam will be selected containing several. Period: 1 year
	DT7. Loop Detector data	The selected subnetwork. Period: 1 year
<b>Pipeline 3</b>	DT6. Floating Car data	A selection of 5 road segments. Period: 1 month
	DT7. Loop Detector data	A selection of 5 road segments. Period: 1 month
	DT8. Queue length from radars	Road segments have different traffic flow characteristics. We select 5 typical road segments for further analysis. Period: 1 month

## 4.6 Implementation Plan

The results from EMERALDS will not be implemented directly in the INM system of the city of Rotterdam (AFM) in the runtime of the EMERALDS project. Before possible implementation new tools will have to be thoroughly tested in practice, benchmarked and validated. This will likely be done in parallel with regular processes. The technical and operational teams need to gain confidence in the outcomes of the developed tooling. Apart from the technical implementation of EMERALDS tooling, it is important to determine which information should be provided to whom, e.g., to avoid information

overload and to provide information on a need-to-know basis only. Also, the way in which EMERALDS results are presented to users will be determined in close corporation between TUD, Arane, ARG, different departments for the city of Rotterdam.

The assessment approach presented in Chapter 2 will be tailored to the requirements and implementation plan of the UC, within T5.1 in collaboration with T5.3, and further detailed in D5.4 (M24) and final results following the second development-assessment cycle in D5.5 (M33)

- **Phase 1** Use Case Preparation, Modelling and Training phase **M7-M15**
- **Phase 2** Model testing phase **M9-M15**
- **Phase 3** Implementation and first validation cycle **M12 – M24**
- **Phase 4** Preparation (Adjustments, necessary modifications) for second validation cycle **M24 -M30**
- **Phase 5** Implementation and second validation cycle **M30 – M33**

## 4.7 Key Stakeholders

---

The key stakeholders in UC2 are:

### System Providers:

- NDW (National Access Point multi-modal mobility data)
- Municipality of Rotterdam (road authority responsible for traffic operations, responsible for all traffic management hardware and the traffic control center)
- Technolution (system and service provider for MobilMaestro system)
- Fileradar (providing algorithm to process Radar data)

### System Users:

- Municipality of Rotterdam (road authority responsible for safe, smooth and sustainable traffic flow operations)
- Travellers
- emergency services (police, ambulances, fire brigade)
- PT companies (bus- and tram companies)

### System Operators:

- Municipality of Rotterdam (road authority responsible for running the system, traffic control, etc.)
- Technolution (provides platform for multi-modal traffic management)
- Arane (responsible for control system development, control strategy development in line with policy objectives of the city)

The stakeholders are partaking in the project either as partner (Arane), or as associated partner (Rotterdam), or are indirectly involved (engagement via Arane).

## 5 UC3: Public Transport Trip Characteristics Inference and Traffic Flow Data Analytics

Riga, being the strategic and economic centre of Latvia, deals with an increased number of commuters from the outskirts and neighbouring municipalities to the city as people have relocated to areas with lower population density. This outflow has a direct impact on the transportation system and creates the need for AI-driven decision-making tools for mobility planning. To optimise the transport network and its efficiency, make Public Transport (PT) more attractive to the user, adapt to changes and forecast future demand, real-time data analytics of travel behaviour is essential. Detailed reporting on the implementation, validation and assessment of the developments carried out within UC3 will be presented in D5.6 (first cycle) and D5.7 (second cycle and refinements)

One of the key objectives of UC3 is to develop advanced algorithms and techniques for inferring trip characteristics from public transport data. This involves analysing various data sources, such as timetable information, ticketing and GPS data to extract meaningful information about trip duration, route choices, transfer patterns, and passenger preferences. Leveraging machine learning and data mining techniques, UC3 aims to uncover hidden patterns and correlations within the data to infer valuable trip characteristics.

### 5.1 Background – Context – Challenges

"Rīgas satiksme", Ltd. was founded in 2003 as a municipal corporation to provide accessible, safe, and sustainable public transport services. Today, the company operates 6 trams, 22 trolleybuses and 51 bus routes (see Figure 5-1).

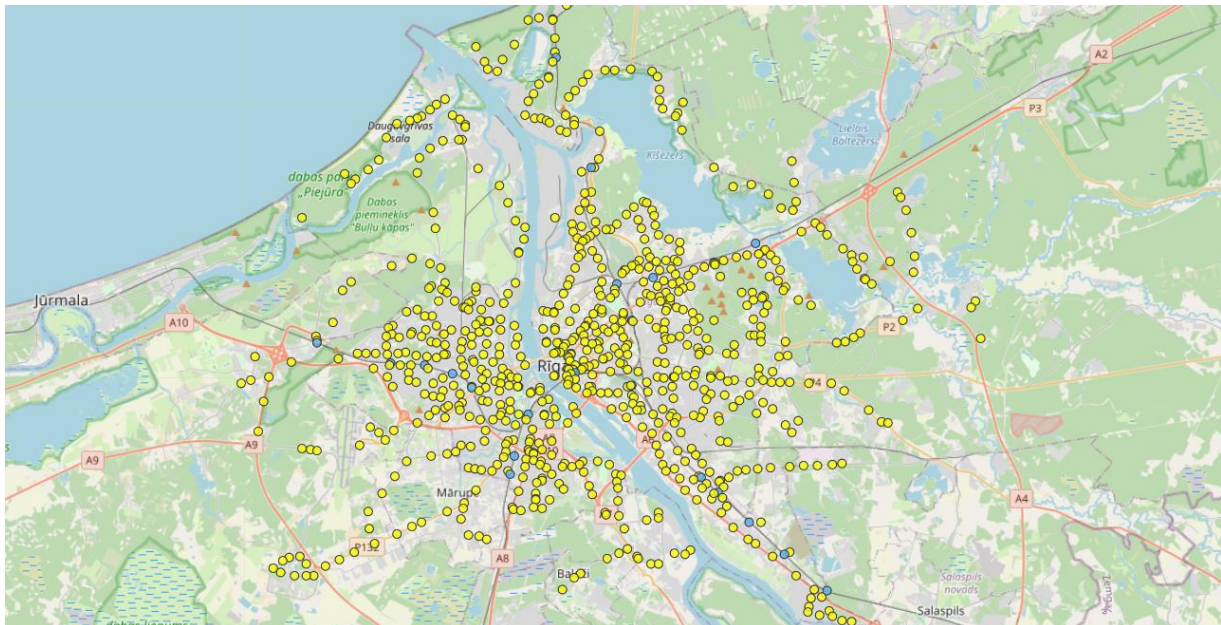


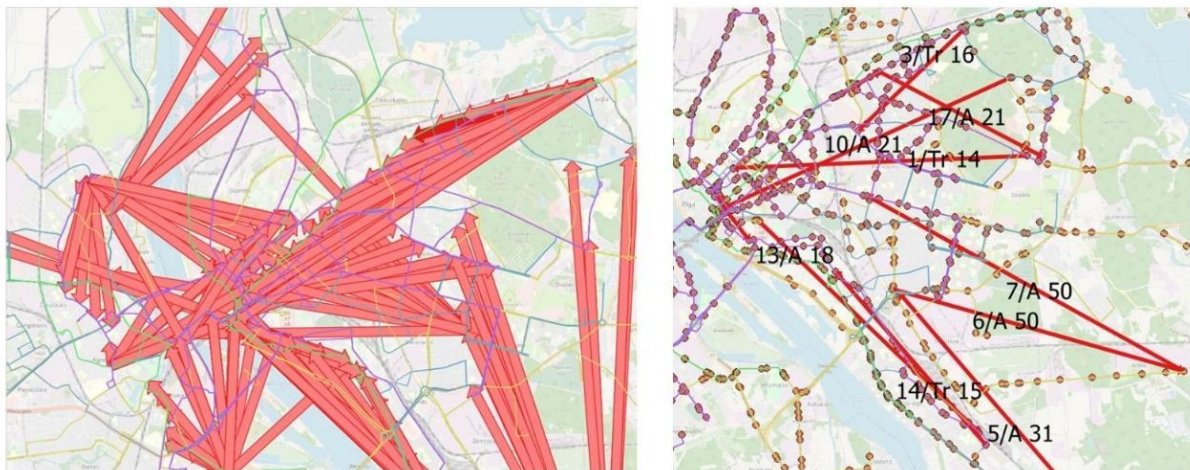
Figure 5-1 – The map of public transport stops in the city of Riga [34].

In 2023, the vehicles of the Riga public transport company (trams, trolleybuses, buses) covered around 35 million kilometres per year and transported around 114 million passengers [35]. Although the number of passengers has increased by 23% compared to 2022, the numbers still lag behind the popularity of PT in 2017. Based on periodic surveys on the movement habits of the citizens, conducted



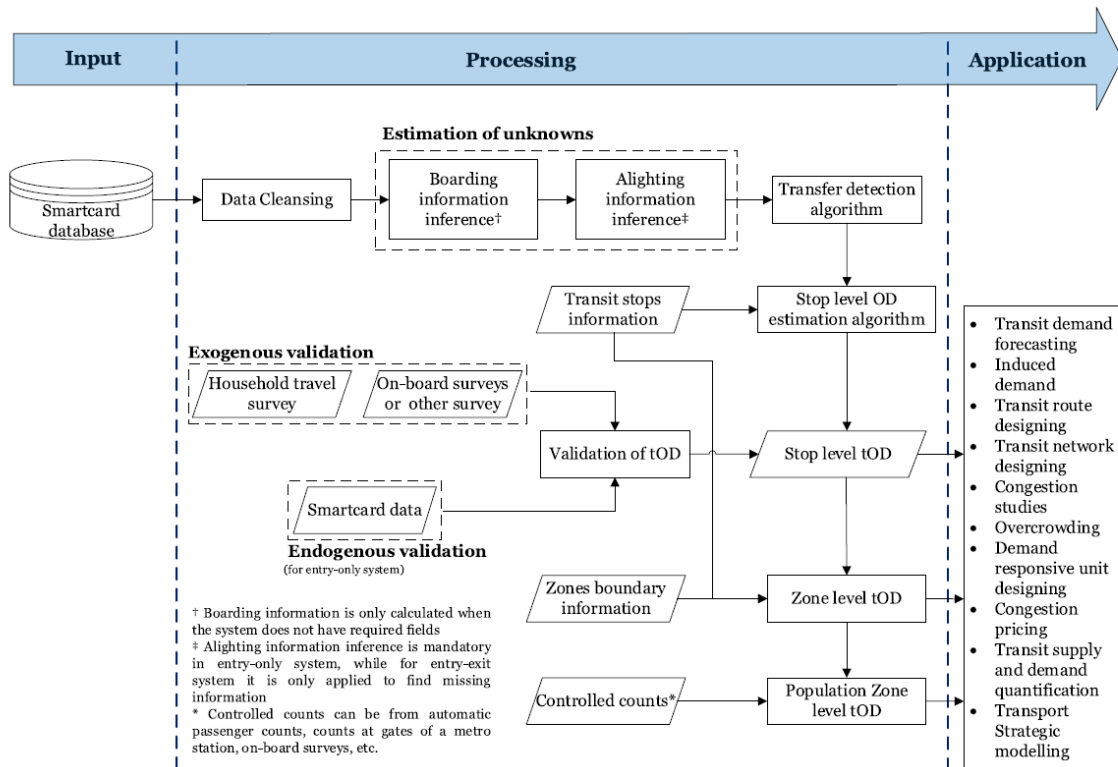
by the Central Statistical Bureau, in 2017 there was a significant increase in the use of PT as the primary form of movement. This was followed by the COVID-19 pandemic, which had opposite effects on the PT sector, decreasing demand. If in 2017, 19.1% of the population in Riga used PT as the primary means of transportation in the city and its outskirts (<100km), then according to the data of 2021 the percentage decreased to 11.5. These external factors that cannot be influenced and factors such as an efficient and accessible public transport system, sustainable mobility prioritisation and the reduction of the use of private cars in the city are some of the reasons why Riga needs new, smarter solutions for network management.

Until now, Grupa93 Ltd. (G93) and "Rīgas satiksme" Ltd. have collaborated on various projects and hackathons with the aim of jointly developing tailored urban mobility data analytics solutions and tools. The data provided by the PT operator have been used in the development of the methodology for reducing CO<sub>2</sub> emissions in case of new mobility hub implementation, planning the new route integration of passenger ferry traffic in the city of Riga, etc. However, the data mentioned below in this section have been used mainly experimentally in hackathons organised by the municipality. For example, e-ticket validation data, GPS and GTFS data were combined and the first draft of the algorithm for trip chaining (origin–destination) was developed with the vision in mind to produce detailed origin-destination data needed for more advanced network, fleet and urban space planning. (see Figure 5-2).



**Figure 5-2 – Visualisation of the origin-destination analysis conducted so far by G93: PT trip flows (left), Trajectory data with entry and exit spots (right).**

Aggregating a complete and versatile dataset on people's travel habits is challenging, requiring the validation of the entry, exit and transfer points of the trip (to determine the travel route), and the participation of multiple stakeholders, in the dataset development.



**Figure 5-3 – Typical steps of OD estimation problem [36].**

The framework of typical OD estimation is presented in Figure 5-3 where processing can be divided in following four parts:

- data cleansing;
- transfer detection;
- estimation of entry and exit locations;
- estimation of stop-to-stop transit and validation of results acquired [36].

For PT in Riga there is an entry-only system which means that passengers are required to validate their tickets only when entering the vehicle. Based on literature review done of multiple studies, the trip-chaining algorithm is the most used when there is an entry-only system. As mentioned before, a preliminary algorithm for joining the datasets and geolocation validations provided by Riga Traffic company has been developed and tested by G93. However, it is essential to improve the algorithm and real-time data aggregation and analytics methods to evaluate longer periods (so far only specific day travel behaviour has been analysed).

UC3 is working on not only providing the justification for prioritising PT in Riga but also **generating the data that is necessary for re-organising the PT network** in the city. It is already known that there should be more attention to PT as one of the main modes of transportation, but till today it is a politically sensitive question, given that mainly transport planning is around private cars. One of the goals of this UC is to show that the number of passengers that use PT is much higher than expected and to provide specific evidence on spots where the reorganisation and investments in the network are needed and will benefit the most.

Riga UC will contribute to defining the functionality of the tools EMERALDS provides and later on testing them. The work in the project will be used to create sustainable mobility and urban planning analytics services. The prototype algorithm for joining Riga public transport data sets, estimating entry points and trip chaining for exit points will be implemented for better efficiency building on the results of WP3-WP4.

In the scope of this UC, a traffic data analysis toolset to analyse and forecast passenger travel behaviour and optimise the PT network will be developed. Currently, it is possible to geo-locate 66% of all ticket validations/entry points and 37% of all exit points. EMERALDS aims to improve the algorithm to geo-locate >90% of the validations and improve its accuracy. The trip destination is determined by geo-locating the next entry point and searching for the possible exit points for the previous leg in a near radius. Thus, by performing joins with real-time or calculated data on private car traffic, the total cost of mobility per network segment could be assessed and used for planning. Combined with vehicle movement data, the algorithm will be able to calculate time losses via schedule and free-flow conditions. For trips with no subsequent trips, AI could be used to infer additional trip characteristics (most popular destinations, probability distribution, etc.). The toolset will be attractive to municipalities and PT providers to assess passenger travel habits and to evaluate the multimodality potential.

One of the most challenging parts of trip-chaining is the validation of the obtained data, which depends on the input data type and structure. The algorithm can be validated by using either endogenous data – based on the same datasets used for algorithm development or exogenous – datasets from an external source [36]. Examples include:

- Household travel survey - validation is done based on destination distribution of riders, trip patterns and the total number of transfers;
- Trip recording – in volunteers’ data, the significant proportion is of students, which can also possibly introduce biases in the results;
- Entrance and exit counts of passengers at the stops;
- Local travel demand surveys - the following questions in different part of week might be asked: (1) origin–destination of the trip leg on the surveyed line; (2) the connecting line before the surveyed one, if any; (3) the connecting line after the surveyed one, if any; (4) origin–destination of the full trip coded at the stop level; (5) socio-demographic of the passenger [36,37].

The following sections provide a more detailed vision of how UC will come together with EMERALDS provided by technical partners, objectives, KPIs and overall implementation plan.

## 5.2 SMART Definition and Use Case Objectives

Based on the planning documents of the municipality of Riga, PT has been identified as one of the priorities in the mobility hierarchy. However, in practice, actions show a completely different reality. For several years, data on e-ticket validations have been publicly available, and the PT company uses them to compile and publish statistics on the total number of transported passengers and the most demanded routes. Yet, further data analysis, such as on specific segments, delays and their causes, has not been carried out so far. Considering what is mentioned before, the following subsection presents a detailed explanation of the objectives set by UC3 and their connection with tools provided by EMERALDS.

**Objective 1:** Improve trip-chaining algorithm for determination of entry stops

Specific	<p>Accurate identification of entry stops is fundamental for understanding passenger movement patterns. It involves analysing trip-chaining behaviour, where passengers begin their journeys and how they connect to subsequent legs of their trips. Where, for example, machine learning algorithms play a vital role in enhancing the specificity of the trip-chaining algorithm.</p> <p>By focusing on the specificity of entry stops, the objective addresses a key element in public transport optimisation, aiming to provide a more granular and</p>
----------	---



	accurate understanding of passenger behaviours during the boarding phase of their journeys.
Measurable	<p>Until now, 285 thousand journeys or 66% of entry stops have been determined for one day in 2019. It was not possible to reach a higher value, because the trajectory data was missing in certain routes due to a server malfunction. Therefore, the aim is to increase the algorithm's precision in determining at least 90% of entry stops, as measured through comparative analysis with the current algorithm.</p> <p>To provide greater reliability of the results obtained in the project, it is expected to organise data validation campaign.</p>
Action - oriented	<p>For reaching the improvements, the following activities are foreseen:</p> <ul style="list-style-type: none"> <li>• Data provision, incl. already achieved results;</li> <li>• A description is prepared for the previously performed activities so that it is possible to duplicate and improve the algorithm;</li> <li>• Development of Probabilistic approach for trip chaining and Trajectory/Route Forecasting and Origin/Destination Estimation;</li> <li>• Validation of obtained results or ground truth to determine and improve the accuracy;</li> <li>• Development of Privacy aware Crowding data ingestion to make sure that there are no possibilities of data de-anonymisation.</li> </ul>
Relevant	The need for improvements is relevant both as a separate activity and as part of the OD matrix. Reaching a higher number of entry stop determination can help to identify the most crowded stops in different times of the day, which, however, can help to optimise the frequency of transport, knowing the required capacity. Furthermore, entry stops as part of the OD matrix can help to identify the most popular routes for PT network optimisation.
Timebound	M15 - first intermediate results; M18 - first results of improved algorithm by using 2019 data; then improved with other years data by month 30.

**Objective 2:** Improve trip-chaining algorithm for estimation of exit stops

Specific	The objective focuses on refining the trip-chaining algorithm specifically for the accurate estimation of exit stops in PT routes. Exit stops are critical for understanding where passengers conclude their journeys, providing insights into travel patterns and facilitating better route planning and resource allocation.
Measurable	<p>Until now 160 thousand journeys or only 38% of exit stops have been estimated for one day in 2019 via trip chaining algorithm, clearly stating that there is a need for improvements of the algorithm. The overall goal is to improve and use the algorithm also for other days. Therefore, the aim is to increase the algorithm's precision in estimating at least 75% of exit stops, as measured through comparative analysis with the current algorithm.</p> <p>To provide greater reliability of the results obtained in the project, it is expected to organise a data validation campaign.</p>
Action - oriented	For reaching the improvements, the following activities are foreseen:

	<ul style="list-style-type: none"> <li>• Development of Probabilistic approach for trip chaining and Trajectory/Route Forecasting and Origin/Destination Estimation;</li> <li>• Testing Euclidian and Network distance approach to trip chaining;</li> <li>• Progressing from use of one day data to several consequent days for greater prediction accuracy;</li> <li>• Validation of obtained results or ground truth to determine and improve the accuracy.</li> <li>• Development of Privacy aware Crowding data ingestion to make sure that there is no possibilities of data de-anonymisation.</li> </ul>
Relevant	Accurate exit stop estimations can lead to better-informed decision-making for route planning, resource allocation, and service optimisation. It directly impacts the quality of public transport services and addresses the needs of both passengers and transportation authorities. Furthermore, exit stops as part of the OD matrix can help to identify the most popular routes for PT network optimisation. The objective is especially relevant now that PT in Riga has lost its demand in recent years, where increasing the efficiency of the network and creating connections according to demand is an essential step.
Timebound	M15 - first intermediate results; M18 - first results of improved algorithm by using 2019 data; then improved with additional years of data by month 30.

**Objective 3:** Data fusion, compression and enrichment

Specific	The specific of this objective relies on merging diverse datasets related to PT, data compression for efficient long-term storage, and enrichment of the datasets with additional details such as weather, time of day, week, season etc. to support in-depth segment analysis.
Measurable	Success is measurable by the degree of integration achieved among different datasets, the extent of data compression, and the depth of enrichment, indicating the amount of additional valuable information incorporated specifically for segment analysis.
Action - oriented	<p>The following activities are foreseen:</p> <ul style="list-style-type: none"> <li>• Provision of the data;</li> <li>• Development of Weather enrichment, Extreme-Scale Map Matching, Mobility/trajectory data compression and Sensor GPS and GTFS data fusion;</li> <li>• Training and testing.</li> </ul>
Relevant	The objective aligns with the broader goal of improving data analytics for public transport. Given the large volume of GPS or vehicle trajectory data, it has not been possible to access all the data that is being produced, where data compression comes relevant. Data fusion is relevant to connect different data sets, as the data schema is different and also changes independently. However, by creating a consolidated and enriched dataset, the analysis becomes more insightful, supporting informed decision-making in public transport planning and management
Timebound	M15 - preliminary work on data fusion has been done, as it is a part of the improvement of the trip-chaining algorithm. M18 - it is expected to get

	intermediate results of all EMERALDS mentioned on this objective; Final results by M30.
--	---

**Objective 4: Public transport and street network segment analytics**

Specific	This objective focuses on a detailed analysis of street network and public transport route segments to identify potential areas for optimization, enhance connectivity, and improve overall transportation efficiency.
Measurable	As there is no prior work done in case of this objective, UC3 will keep in mind the <i>Pareto principle</i> , which states that for many outcomes, roughly 80% of consequences come from 20% of causes. In this case, meaning that most of the rides with public transport most likely happen in a small share of the whole network.
Action - oriented	To perform segment analysis, the following activities are foreseen: <ul style="list-style-type: none"> <li>• Finalised table of entry and exit stops, with missing data filled by Probabilistic approach for trip chaining and Trajectory/Route Forecasting and Origin/Destination Estimation;</li> <li>• Completed Trajectory data/travel time analysis for the whole public transport network in form of time losses (passenger hours);</li> <li>• Public transport route and street network scaling for which Map-matching will be developed;</li> <li>• Detailed specifications for segment analysis, encompassing parameters displayed, views, and data filters. s.</li> </ul>
Relevant	The purpose is to pinpoint prospective zones for optimisation, amplify connectivity, and elevate the overall efficiency of transportation systems. By studying the complexity of street and transport network interactions, the analysis aims to uncover opportunities to streamline routes, enhance accessibility, and ultimately contribute to an advanced and more effective public transportation infrastructure. The relevance lies in the potential to identify areas of improvement, ensuring that the transport system is adjusted to meet the evolving needs of the community and promote sustainable urban mobility.
Timebound	Month 15 – first intermediate results of trip-chaining algorithm in form of OD matrix table and analysis of travel time after which work on data categorisation will start; M18 - preliminary results of segment analysis delivered in CARTO platform, continuing work until M30. Moreover, development of Map-matching.

**Objective 5: Visualization and analysis of priority areas**

Specific	This objective is specific in its focus on visualising and more in-depth analysis of specific aspects of PT operations. It aims to delve into the details of passenger travel patterns and time spent, specifically their origin-destination dynamics, while concurrently conducting a comprehensive segment analysis of varying scales within the PT routes and street network. The overall aim is to present analysis of priority areas of the PT network to identify the locations that are the most significant for further inspection.
----------	---

	In the scope of this objective with analysis of priority areas it is understood that data will be presented in the way of showing the most problematic locations/points or even routes.
Measurable	The objective's outcomes can be measured through the development and application of visualisations that illustrate passenger travel patterns, analysis of priority areas to identify spots of concentrated passenger activity, and detailed segment analyses that provide insights into the efficiency and effectiveness of different scales of public transport routes and street networks.
Action - oriented	To perform this objective the following steps needs to be done: <ul style="list-style-type: none"> <li>• <i>1st implementation cycle</i>: Data analytics visualisation in CARTO platform (2019 data) with basic functionality.</li> <li>• <i>2nd implementation cycle</i>: The final data analytics visualisation in the CARTO platform is performed with expanded functionality.</li> </ul>
Relevant	The objective is highly relevant as it directly contributes to optimising PT systems. By visualising travel patterns and conducting analysis of priority areas, the city and PT companies can identify high-traffic areas, understand passenger behaviour, and make informed decisions to enhance the efficiency and responsiveness of the public transport network.
Timebound	M18- Data analytics visualisation of the 2019 data is performed. The main work on this objective is expected to begin of the 2nd implementation cycle around M30 as there is need for other objectives to be done, including real-time prediction of crowding and delays.

#### Objective 6: Crowd density forecasting

Specific	This objective aims to predict crowdedness levels for trips conducted on public transport. It focuses on creating a model/algorithm that can analyse various factors influencing crowding, such as passenger volume, weather and other events such as time of day, week, season etc.
Measurable	<p>There are multiple ways to measure this objective's success. Models' accuracy and metrics will be further investigated in D5.6 as the work on this objective is planned to begin on the 2nd implementation cycle and further discussion with technical partners are needed.</p> <p><b>Accuracy Metrics:</b> Measure the accuracy of the predictions by comparing the predicted crowding levels with the actual observed values.</p> <p><b>Prediction Error Rates:</b> Involves quantifying the difference between predicted and actual values for crowding levels.</p> <p><b>Validation and Testing:</b> Using historical data and real-time observations. This process involves dividing the data into training and testing sets, training the models on the training data, and evaluating their performance on the testing data.</p>
Action - oriented	The following steps needs to be done: <ul style="list-style-type: none"> <li>• Determination of the scope and timeframe that will be covered;</li> <li>• Gathering of the data that is necessary, such as historical data of passengers' trips and weather conditions;</li> <li>• Development of Weather enrichment and Crowd density forecasting;</li> </ul>

	<ul style="list-style-type: none"> <li>• Training, testing, and validating the model.</li> </ul>
Relevant	By providing real-time information on crowding levels, passengers can make informed decisions about their travel routes and schedules, while transport authorities can optimise operations to mitigate congestion and improve service reliability.
Timebound	The main work is planned to start on the 2nd implementation cycle, but as there is need for the input data preparations of other EMERALDS used for achieving goals of UC3, some steps might already be performed by the end of 1st cycle.

For UC3 user stories address the diverse needs of both PT users and Riga Traffic Company/City of Riga, ranging from optimising routes and prioritising PT in urban infrastructure to efficiently managing the PT fleet. For each of the user stories there are scenarios that describe a different implementation. The following user stories have been identified:

### 5.2.1 User story 1: Optimisation of existing PT network

This user story deals with the task of transport planners to analyse passenger behaviour, to understand individuals' travel patterns and preferences. In Riga, the public transport network hasn't undergone a significant reform for over a decade now, however, there are ongoing discussions about the necessity of it. Despite having data on e-ticket validation, it's currently used for general statistics on the demand for different routes. Analysing origin-destination data at different times of the day and understanding how passengers transfer between public transport nodes will provide crucial information for transport planners for further analysis and network optimisation.

User story 1 has 3 new scenarios:

- Scenario 1-0 (base): public transport route use and time statistics
- Scenario 1-1: stop blocks analysis/specific stop clusters (less than 50m proximity) (UC objective 1, 2, 3, 4, 5)
- Scenario 1-2: defined public transport zones analysis (UC objective 1, 2, 3, 4, 5)
- Scenario 1-3: planning mobility hubs (busiest interchange locations, interchange time) (UC objective 4, 5)

### 5.2.2 User story 2: Prioritising public transport in the mobility hierarchy

This user story deals with the task of mainly city planners but also transport planners with designing sustainable and accessible transport, more specifically, PT network. In all the intermediate and long-term strategies and policies that the city of Riga has implemented, PT has been identified as one of the highest priorities in mobility hierarchy. However, these objectives and aims don't match the reality, private cars still being the centre of infrastructure planning. Indicating that current possibilities are used up and there is a need for street network reorganisation. To do so, one of the steps is data analysis on where delays are the highest and most passengers flows occur. As an example, the city plans to create a separate lane for public transport or to close the direction of a certain street for private transport, to prioritise the public transport, data on specific sections are needed. The number of passengers transported, and passenger hours lost in traffic would argue for the rationality and socio-economic return of such activities.

User story 2 has 3 scenarios:

- Scenario 2-0 (base): current system based on intuition and space (lane) availability.
- Scenario 2-1: segment travel time/delay analysis based on vehicle (UC objective 3, 4, 5)
- Scenario 2-2: segment travel time/delay analysis based on number of passengers (UC objective 1, 2, 3, 4, 5)



- Scenario 2-3: segment travel time/delay analysis based on lost time value (UC objective 1, 2, 3, 4, 5)

### 5.2.3 User story 3: Fleet planning and management

This user-story deals with the task of PT providers to efficiently organise and oversee the public transport fleet based on dynamic factors, including real-time data on passenger demand. So far, the decision on fleet planning has been made based on the data of average passenger numbers for specific line and time. However, there are other external factors such as weather, and periods of time as season, weekday and even hour that affects it.

User story 3 has 1 new scenario:

- Scenario 3-0: based on average numbers per line and specific times of day
- Scenario 3-1: detailed crowding prediction for specific lines and locations and times of day (UC objective 5, 6)

## 5.3 Computational and Data Resources

This section details the datasets and their significance in achieving the established objectives. It proceeds to explore the computational requirements and workflow of UC3. Subsequently, it illustrates the convergence of these UC3 requirements, aligning with the principles outlined in D2.1 of the EMERALDS Reference Architecture.

### 5.3.1 Data Resources

Table 7 provides general information about the datasets that are going to be shared with the consortium and used within the project. Based on the availability, additional datasets might be added after M18. For detailed information of the datasets, mentioned in this section, including data schemas, description of metadata etc. see the Data Management Plan (D1.4 v.1.). The examined data sources relevant to the data analytics problem are indicated while specific data attributes will be elicited in D5.6.

Table 7 – UC3 (Riga) Data Resources.

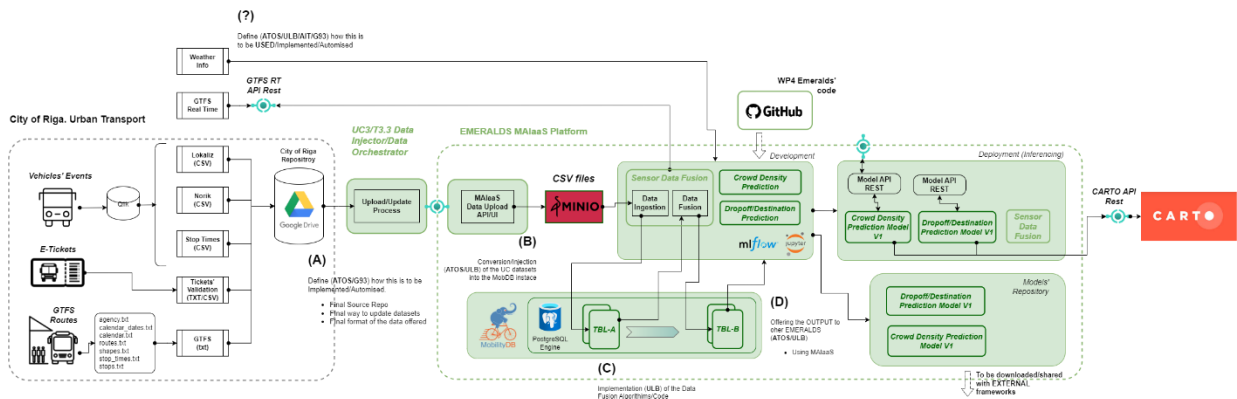
Name of the dataset	Accessibility level	Description	Data format	Relevant objective
DT13 E-ticket validation data	Open access	The dataset contains compiled e-ticket validation data that has not undergone data processing and/or analysis. The data is an important link for UC3 as combining with GTFS data it provides the results of the origin-destination matrix.	txt.	1, 2, 3 as an input data where the outputs of data processing become as an input for objective 4, 5 and 6
DT14 Riga's public transport route list data, GTFS	Open access	Public transport routes list data in different periods for which no data analysis has been performed. The data set is one of the intermediate	txt.	1, 2, 3 as an input data where the outputs of data processing become as an

		data sets to join the validation data set with the vehicle messages data set (GPS trajectories) to determine check-in location and populate validation data with the route and stop information for further data analysis and visualisation.		input data for objective 4, 5 and 6
DT15 Vehicle events data, GPS	In this stage of the project data is shared only with the consortium and therefore only used in the scope of the EMERALDS project.	Compiled public transport GPS data. The data set is used to determine check in location/stop for a respective validation/check-in.	csv., qvd.	1, 2, 3 as an input data where the outputs of data processing become as an input data for objective 4, 5 and 6.

### 5.3.2 Computational Requirements

The computational requirements for the UC3 description will be elaborated in D5.6, along with the technical specification of the EMERALDS toolset integration and testing in collaboration with the technical partners from WP2, WP3 and WP4. EMERALDS. An initial set of requirements includes:

- **Data Collection:** The Riga Traffic Company collects data on vehicle trajectory, e-ticket validation, route schedule and other relevant metrics from their transportation system. Data is received through various sensors and ticketing mechanisms. Data is made available to the consortium during the implementation cycles, including the possibility to integrate further data resources and availability.
- **Data Fusion and Processing:** G93, as an external partner, processes the collected data using their expertise and tools. Processing involves tasks such as data congestion analysis, compression techniques, and applying analytics to derive meaningful insights. Following by utilisation of novel data fusion and management techniques to combine and enrich the data.
- **G93 Contribution to the Process:** G93 adds value by leveraging their capabilities in data analysis and providing valuable insights into traffic patterns, passenger behaviour, and potential areas for improvement.
- **Development of EMERALDS:** Technical partners together with G93 experts, develop tools/algorithms based on UC3 objectives, which will later be delivered into ATOS platform.
- **EMERALDS deployment in ATOS:** Figure 5-5 presents data gathering and integration process of UC3, however it is still under the development and the final version will be provided within D5.6.



**Figure 5-4: – Data gathering and ingestion process as a blueprint for the UC scenario data pipelines that will be reported in D5.6 (M24), D5.7(M33)**

- **CARTO Platform Integration:** Processed results are fed into the CARTO platform through APIs developed within the project. CARTO serves as a visualisation and analysis tool, allowing for a detailed examination of specific routes or segments.
- **Future Service Offering:** The ultimate goal is to evolve the current project into a service that can be extended to parties involved in PT system management, such as the city of Riga and Riga Traffic Company.
- **Creating Value:** The workflow creates value by providing actionable insights into traffic management, leading to improved efficiency and reduced delays. Utilising advanced analytics and visualisation on the CARTO platform, decision-makers can make informed choices for optimising specific routes. The service-oriented approach allows for scalability and the potential to offer valuable solutions to stakeholders.

## 5.4 Assessment and Corresponding KPIs

As established in Chapter 2, a differentiation exists between performance metrics and business KPIs. UC3 revolves around estimation, comprehensive analytics, and prediction of multiple variables essential for effective PT management. The main software performance KPIs used are described below. The KPI1 and KPI2 are further subdivided into sub-indicators. These encompass the designated target value for the share of validation, for which entry and exit points are determined, as well as the accuracy of the results. The specific target value for data accuracy will be defined during the second phase of the project implementation after deciding how to define ground truth.

KPI	Existing value	Target value	Measurement process	Impact on domain
<b>1. Entry stops determined by vehicle position at a certain time</b>				
<b>1 A. Share of unique ID e-ticket validation data for which origins are determined/estimated</b>	66%	>90%	The data will be compared with preliminary validation results	It will contribute to public transport trip characteristics inference and traffic flow data analytics
<b>1 B. Data accuracy</b>	N/A	To be determined in the	To validate the accuracy of data, it is planned to engage students and other	It will provide greater reliability in the obtained results and

		second implementation cycle	volunteers, to spend their day using PT and record rides or counting passengers entering and exiting the vehicle. The objective is to assess whether the algorithm accurately predicts and calculates the correct entry stop.	further data analysis and interpretation.
--	--	-----------------------------	---	---

Existing values for KPI1 are calculated as a number of validations (check-ins) for a representative business day with check-in locations determined through the algorithm that joins validation data with vehicle trajectory data and scheduling (GTFS data) as well as probabilistic estimation of missing data divided by the total number of validations (check-ins). For example, on September 25, 2019, the total number of validations was 417828, whereas the number of validations for which check-in locations had been determined was 274183, making the baseline value of this KPI 66%. These were results obtained before starting participation in the EMERALDS project, where, together with technical partners and applying new methods, it is expected that the new value will reach at least 90%.

After the determination of the entry stops for unique IDs e-ticket validation, it becomes feasible to estimate the exit stops or destinations using a trip-chaining algorithm, which KPI2 is about.

KPI	Existing value	Target value	Measurement process	Impact on domain
<b>2. Exit stops estimated via trip chaining algorithm</b>				
<b>2 A. Share of unique ID e-ticket validation data for which destinations are determined/estimated</b>	37%	>75%	The data will be compared with preliminary validation results	It will contribute to public transport trip characteristics inference and traffic flow data analytics
<b>2 B. Data accuracy</b>	N/A	To be determined in the second implementation cycle	To validate the accuracy of data, it is planned to engage students and other volunteers, to spend their day using PT and record rides or counting passengers entering and exiting the vehicle. The objective is to assess whether the algorithm accurately predicts and calculates the correct entry stop.	It will provide greater reliability in the obtained results and further data analysis and interpretation.

The existing value for KP2 is calculated as a number of validations (check-ins) for a representative business day with exit locations estimated through trip chaining algorithm as well as probabilistic estimation of missing data divided by the total number of validations (check-ins). For example, on September 25, 2019, the total number of validations was 417828, whereas the number of validations

for which check-in locations had been determined was 154187, making the baseline value of this KPI 37%. So basically, for the same person represented by a unique ticket ID, the algorithm, which is based on the literature, chains the trips and it is assumed where the exit stop is located.

The sub-indicator A of the first two KPIs does not reflect the accuracy of the obtained results, but the number that needs to be filled in, given that there is missing data, mostly because of malfunction of the servers. UC3, with the help of tools developed within the EMERALDS project, aims to fill the gap of missing data, to reach set target values. To ensure the accuracy of the outcomes, the plan is to engage students and other volunteers in recording their journeys throughout the day or passenger counting at the PT stops. However, the specifics of this and the target values will be established at the onset of the project's second implementation phase. If new concepts or techniques for representing data accuracy emerge in the coming months, UC3 remains open to adopting other approaches.

Once entry and exit stops on specific routes have been determined/estimated, it will be possible to proceed to assess the number of people entering and exiting at particular stops. Additionally, it will be possible to calculate the number of people travelling between these stops. Another aspect is to use delay information from public transport GPS data to justify the reasons for making public transport improvements a priority.

KPI1 and KPI2 also derive from the goals within the work of WP4, where targets for per cent point improvement over the current algorithm's entry/exit stop estimation percentage rates are identified (see D4.1.). Furthermore, data accuracy where the validation campaign is intended will help to achieve the work of technical partners on Prediction accuracy/Performance and overall project KPI of performing processing and analytics in situ, as its main goal is to be efficient enough to apply to low-resource environments.

KPI	Existing value	Target value	Measurement process	Impact on domain
<b>3. Share of public transport network (segments) analysed</b>	N/A	27%	The data will be extracted to measure the performance of the existing PT network in territories with the highest demand and hot-spot areas.	Share of public transport network analysis will help to address the need for more comprehensive and contextual information to equip public transport planners with improved decision-making solutions and will add valuable insights into the performance of the existing transport network and the call for improvement in mobility services and planning.

As was already mentioned at the beginning of this chapter, the primary objective is to understand the process behind the analytics. After which it is possible to move forward on data analysis. During the first cycle of project implementation, the work is expected to be more focused on the analysis of public transport segments. So far, no existing results have been presented for KPI3. These are the results that UC3 aims to develop within the project and the tools that partners are going to provide. The target value for KPI3 is relatively low, given that this is a new segment. Keeping in mind the Pareto principle, which states that for many outcomes, roughly 80% of consequences come from 20% of causes. In this case, meaning that most of the rides with public transport probably happen in a small share of the whole network.

KPI	Existing value	Target value	Measurement process	Impact on domain
<b>4. Share of street network analysed</b>	N/A	20%	The territory will be determined where the majority or 80% of the passenger flow is concentrated, which will then be the segments of interest for further analysis.	Share of street network analysis will contribute to plan the public transport network as a part and not as a separate unit of the overall transport network on a smaller scale. Leading to various socio-economic benefit calculations, including time and money saved.

Achieving the target value set for KPI4 presents a challenge for UC3, as no prior work has been undertaken in this regard. To analyse the share of the street network, it is necessary to transfer data from public transport into street network segments. This links together UC3 and UC2, considering that in Rotterdam public transport data has already been integrated with street network data. This implies a promising opportunity for knowledge exchange between Riga and Rotterdam.

KPI3 and KPI4 directly relates to the work of WP3-WP4 as the level in new in-depth knowledge about public and street network analysis obtained also depends on that. Here Trajectory Data Analysis can be mentioned and its application of the tool to the PT network in Riga.

In addition to technical (software) KPIs, there are business KPIs relevant to the specific user stories. The KPIs as well as user stories will be refined during the project, below is the present vision of the UC3:

- Increased accessibility to primary destinations served by direct routes (Scenario 1-1) and across defined public transport zones (Scenario 1-2).
- Enhanced service efficiency and customer satisfaction by reducing total passenger waiting time at transfer points, measured in person-hours (Scenario 1-3).
- Increased reliability of trips completed within the scheduled time (vehicle hours: time lost compared to schedule and free flow) (Scenario 2-1).
- Passengers' feedback on satisfaction with service quality and experience (person hours: time lost compared to schedule and free flow) (Scenario 2-2).
- Demonstration of cost-effectiveness and sustainability by decreasing time loss compared to schedule and free flow (EUR) (Scenario 2-3).
- Optimisation of public transport resources to ensure accessibility and comfort, thereby raising service standards and enhancing customer satisfaction (Scenario 3-3).

## 5.5 Mapping to EMERALDS Toolset

Table 8 displays the outcomes of the mapping exercise performed with technical partners to align EMERALDS tools with UC3 objectives through a series of meetings. It is important to note that this does not rule out the potential consideration of other tools in further implementation phases based on their potential for added value. Objective 5 is not part of this exercise as UC plans to use CARTO platform for data presentation and analysis.

Table 8 – Mapping of emeralds components to UC3 objectives

EMERALDS component	Objective	Outcome
<b>Privacy-aware in situ Data Harvesting (T3.1)</b>		
Privacy aware data ingestion	1, 2	Risk screening to ensure that there is no likelihood of data de-anonymisation, meaning that people's privacy while analysing the data won't be endangered.
<b>Extreme-scale Cloud/Fog Data Processing (T3.2)</b>		
Weather enrichment	3, 6	Enriched data for segment analysis and crowding density predictions.
Extreme-scale map-matching	3, 4	Map-matched public transport routes and street networks for further GPS data enrichment and segment analysis.
<b>Mobility Data Fusion and Management (T3.3)</b>		
Mobility/trajectory data compression	3	Compressed GPS, GTFS data for further accessibility and effective storage.
Sensor (GPS, GTFS, radar, etc.) data fusion	3	Merged different datasets for more comprehensive and in-depth PT network analytics.
<b>Extreme Scale MDA at the CC (T4.1)</b>		
Dropoff / Destination Prediction	1, 2, 5	Higher accuracy provided for already performed trip-chaining.
Probabilistic approach for trip chaining	1, 2, 5	Supplemented results of trip-chaining where more exit stops will possibly be determined.
Trajectory data / travel time analysis	4, 5	Acquired results of travel time analysis by presenting delays, where together with OD matrix it will be possible to calculate passenger hours lost and monetise them.
<b>Active &amp; Federated Learning over Mobility Data (T4.2)</b>		
Crowd density forecasting	5, 6	Forecasted how many passengers are going to be at the PT stops with enriched PT data.
<b>Mobility AI-as-a-Service (T4.3)</b>		
MLOps Platform	4,5	Real-time training and model inference
<b>Deployment of Services and Dashboards (T2.2)</b>		
CARTO Visual Analytics	All	Visual analytics and interactive dashboard for retrieving insights.

To achieve goals set on Objectives 1 and 2 which align with KPI1 and KPI2, all datasets provided by UC3 are used as input data. To do so, EMERALDS T4.1 Probabilistic approach for trip chaining and Trajectory/Route Forecasting and Origin/Destination Estimation will be developed with an OD matrix table on individual validation (ride) level as output data. GTFS and GPS datasets are used to develop functionalities provided by T4.1. Trajectory data/travel time analysis, with PT delays as outcome. After

which passenger analytics per stop/segment will be performed, using CARTO platform with customizable dashboards for visual presentation of data analytics.

In parallel with the work on Objectives 1 and 2, work on Objective 3 will be started since GTFS and GPS data fusion are already required in this step. The datasets (primarily GPS data) will be enriched with street network and weather information, which together with OD matrix and PT delays will provide an input for Objectives 4 and 6 (see Figure 5-5).

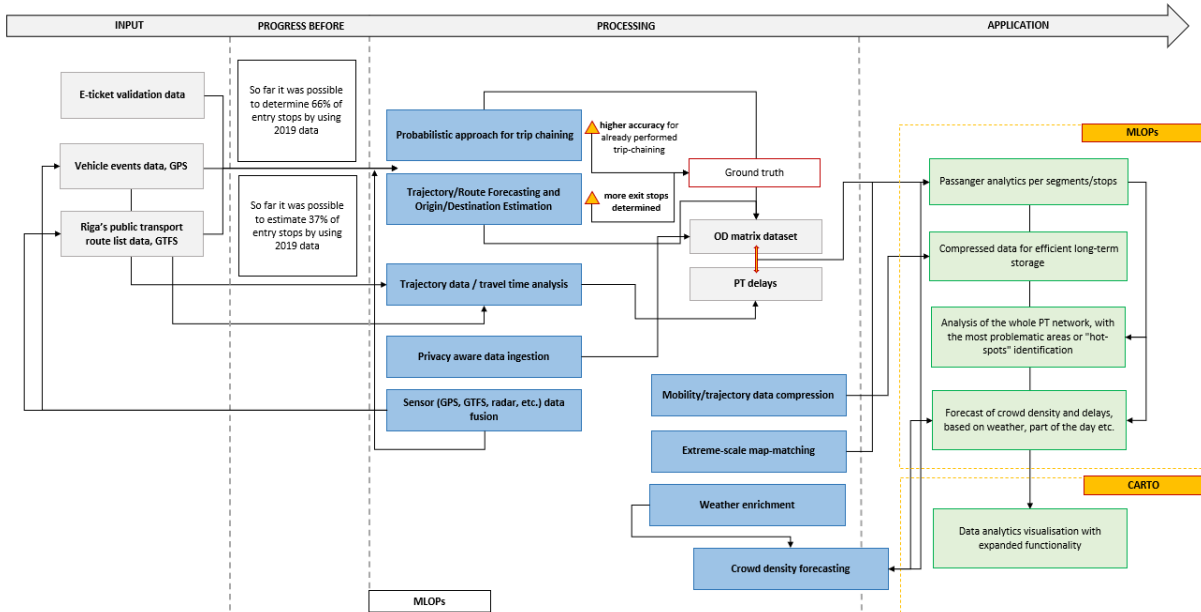


Figure 5-5 – UC3 pipeline.

Addressing extreme scale in the context of public transport data involves dealing with large volumes of data, diverse data types, sensitivity to real-time updates, rapid data velocity, intricate data complexities, and heterogeneity arising from the diverse sources.

- **Volume:** Public transport systems generate vast amounts of data, including information on routes, schedules, and more. Handling this immense volume efficiently requires robust computing infrastructure.
- **Variety:** Data in public transport comes in various formats, such as structured schedules, unstructured text data, and real-time sensor data. Managing this diverse range of data types is crucial for a comprehensive analysis.
- **Sensitivity:** Public transport data is time-sensitive, and decisions based on outdated information may lead to inefficiencies or disruptions.
- **Velocity:** The speed at which public transport data is generated and updated is high. Systems must be capable of processing data streams in real-time to provide actionable insights promptly.
- **Complexity:** Public transport data is complex, involving multiple variables like traffic conditions and system performance.
- **Heterogeneity:** Data in the public transport domain originates from various sources, such as GPS devices and ticketing systems. Integrating and interpreting this heterogeneous data is a significant challenge.

UC3 is critical in showcasing the computing continuum tech stack because it demonstrates the need for a comprehensive solution that can handle extreme-scale data effectively. The computing continuum tech stack ensures a seamless flow of data processing from edge to cloud, enabling real-time analytics, machine learning applications, and decision-making.



In summary, the importance of the Riga UC lies in its ability to stress the computing continuum tech stack by presenting a scenario where extreme-scale data characteristics are addressed cohesively to derive valuable insights and enhance the efficiency of public transport systems.

## 5.6 Implementation Plan

The implementation of UC3 can be divided into 3 main phases – Preparation, 1st, and 2nd Implementation cycle. Further subdivided into more specific actions, milestones and sequence of objectives mentioned in this document. As the 2nd cycle is more or less depending on the results of actions planned in the first two years, UC3 has prepared the overall timeframe of the Implementation throughout the project and a more detailed plan until the end of M24 (see Figure 5-6).

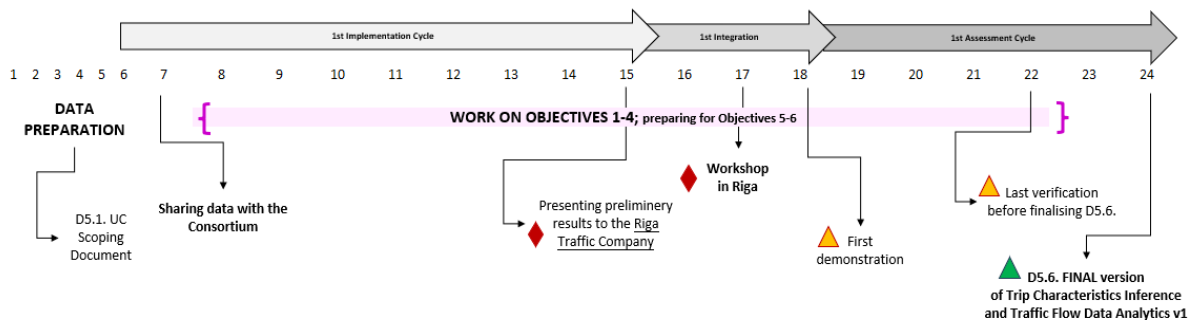


Figure 5-6 – UC3 milestones during the 1st implementation cycle.

The first 6 months were allocated for the preparation of the datasets and the outline of the first version of the UC’s objectives. To do so, several meetings with the Riga Traffic Company were organised, being both data providers and the potential user of the tools developed within the project. During the meetings experts from the IT department and Project Management of the Riga Traffic Company were introduced to the general goal of the project and what is the role of Riga being one of the UC. Further, the first draft of objectives and their mapping with EMERALDS offered within the project was done.

It was already mentioned that UC3 has done some previous work on developing the trip-chaining algorithm, therefore when sharing the datasets with the consortium description of the steps was prepared. This allowed technical partners working with the Trip Probabilistic approach to replicate work done previously and further improve the results aiming to reach goals set within this UC. Kick-off meetings with technical partners were an important milestone as they allowed not only to get a deeper understanding of the vision of how tools in the case of Riga could work but also to revisit objectives and KPIs.

It was decided that for initial work 2019 data will be used since it is more complete. However, in the later stage, it is planned to test tools and algorithms developed with newer data for which the structure differs. In that way lowering the risk of constantly changing data as the Riga Traffic Company currently is going through changes in their systems. Entering the second year of the project UC3 plans to continue close work with technical partners on Objectives 1-4 so that in M15 intermediate results can be presented to the Riga Traffic Company. This step is important as there is still a risk of PT companies’ low interest with the potential to provide proof of the concept and more active participation. According to the UC3 plan by M15, the first version of software tools is released giving a 9-month period for the UCs to test, measure and validate the tool’s capabilities to address their requirements, documenting this procedure in D5.6. That includes the development of specification of ground truth for accuracy of the results of Objectives 1 and 2, based on validation approaches mentioned at the beginning of this section.

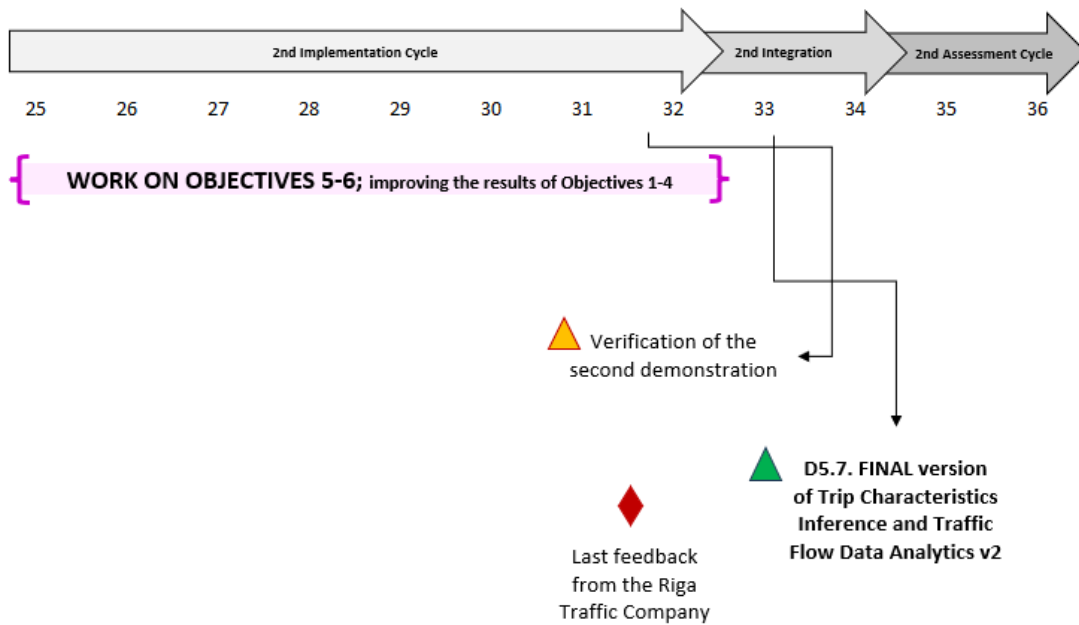


Figure 5-7 – Draft of UC3 milestones during the 2nd implementation cycle.

M17 is expected to be one of the key milestones as it is planned to organise a working group in Riga. The aim is to get feedback on work done so far, seek potential improvements, develop a demonstration plan and elaborate on user stories. After which it will be possible to set a more detailed programme for the 2nd Implementation cycle. **The assessment approach presented will be tailored to the requirements and implementation plan of the UC and further detailed in D5.6 Trip Characteristics Inference and Traffic Flow Data Analytics v1 (M24) and D5.7 Trip Characteristics Inference and Traffic Flow Data Analytics v2 (M33).**

## 5.7 Key Stakeholders

The primary stakeholder for this UC implementation is Riga Traffic Company, which provides data and vital business insights, adding significant value to the development and application of tools also in other similar contexts. Grupa93 has been at the forefront, facilitating discussions and exchanging information regarding project activities, engaging with transportation planners and IT experts. However, it's worth noting that there are other stakeholders whose roles can be identified based on their contributions to data provision, system operation, and utilization. **The key stakeholders in UC3 are:**

### Data Providers:

- Rīgas satiksme” Ltd.
- Latvian open data portal.

### System Users:

- “Rīgas satiksme” Ltd.;
- Municipality of Riga;
- Urban and transport planners and others
- Travellers.

### System Operators:

- Rīgas satiksme” Ltd.;
- G93.

## 6 Conclusions and Next Steps

---

The EMERALDS project is focused on developing innovative solutions for urban mobility challenges through the use of Mobility Analytics as a Service (MAaaS) tools. The deliverable D5.1 plays a crucial role in coordinating the various use cases (T5.2-T5.5) within the project. The main objective is to ensure a comprehensive assessment of the methods and to identify the added value of the MAaaS toolset. Task 5.1 takes a proactive approach, aiming to maximize the potential for assessing the developed innovations and coordinating the use cases towards that goal.

Chapter 2 of the project document discusses the assessment methods used within the EMERALDS project. The assessment process consists of several stages, including technical assessment, legal assessment, functional assessment, user-acceptance assessment, impact assessment, and socio-economic assessment. Each stage has specific assessment questions and corresponding Key Performance Indicators (KPIs). While all stages are considered, the primary focus is on the functional assessment, which evaluates the performance of the MAaaS tools in meeting the specific objectives set for each use case.

Chapter 3 focuses on Use Case 1, which deals with risk assessment, prediction, and forecasting in crowded events. The objective is to design a toolkit that utilizes extreme data analytics and AI methods to assess the risk of injuries or damage in crowded events. Relevant data sources, including social data, sensor data, police reports, weather data, and event data, are analysed to predict and forecast crowd situations and identify risk factors. The toolkit includes tools for prediction and forecasting, as well as assessing the current, predicted, and forecasted risk levels. A pilot of the toolkit will be conducted in The Hague, with stakeholders providing feedback on its performance and usability. The final step involves assessing additional data sources, enhancing performance, and analysing the potential for broader deployment and business cases.

Chapter 4 explores Use Case 2, which focuses on multi-modal traffic network management in Rotterdam. The goal is to improve the functionality of the current system by utilizing novel data services and tools. Traffic state estimators, network analytics tools, and traffic state predictors are developed to enhance the performance of the Integrated Network Management (INM) system. The focus is on improving traffic management operations, such as automated configuration of actuator control goals and incorporating multi-modal bottlenecks in the system. The use of predicted exceedances of policy norms and the ability to anticipate them with the use of actuators are also important aspects.

Finally, Chapter 5 discusses Use Case 3, which is centred around trip characteristics inference and traffic flow data analytics in Riga. The objective is to optimize the public transport network and provide data for its reorganization. A traffic data analysis toolset is being developed to analyse and forecast passenger travel behaviour, making public transport more efficient and attractive to users. The toolset utilizes real-time data analytics and algorithms to estimate entry points, trip chaining, and exit points. The integration of this data with real-time or calculated data on private car traffic allows for a comprehensive assessment of the total cost of mobility per network segment.

Moving forward, the project will continue to refine and implement the use cases (T5.2, T5.3, and T5.4), validating the MAaaS tools within the specific contexts, under the coordination of T5.1. Consolidation of the implementation plan will be facilitated through the creation of focus groups for each use case and the organisation of a series of technical workshops between use case stakeholders and technical partners (from WP2, WP3 and WP4). A concrete validation procedure for the EMERALDS overall technological offering comprises a point of discussion among partners, albeit the assessment approach introduced in Chapter 2. Stakeholder feedback will be crucial in improving the performance and usability of the tools. An initial version of the EMERALDS toolset will be released, followed by testing and refinements based on stakeholder input. The project aims to provide a revised version of

the toolset and conduct additional validation cycles to ensure its effectiveness and applicability in real-world scenarios.

The use cases of the EMERALDS technology and toolset are meticulously structured to ensure thorough development and effective deployment. In this direction, a hierarchical approach has been established, beginning from stakeholder engagement, which is crucial for understanding the needs, challenges, and expectations of the stakeholders. Upon building trust and establishing clear communication channels, work within WP5 orchestrated the determination of the grade of alignment between the stakeholders' requirements and priorities with the functional specifications of emeralds services. Engaging with stakeholders early in the WP5 activities from M1 assisted in identifying the key areas of focus as presented throughout the UC objectives and overall requirements, whilst laying down a solid foundation for the subsequent phases.

Following stakeholder engagement, the process advances to capturing requirements, data collection, and preparation, as well as the definition of UC objectives. At this point, detailed requirements are gathered from stakeholders, relevant data is collected, and preparations for analysis are made. This step is critical for developing a comprehensive understanding of the problem domain and ensuring that the necessary data is available for effective solution development.

Subsequent to defining objectives and data, the approach involves mapping these requirements to the specific tools and capabilities within the EMERALDS toolset. This step optimizes the use of the technology by ensuring that the appropriate tools are utilized for the respective tasks. Following this, use case technical focus groups are formed, comprising experts who address the specific needs of each use case from WP2, WP3, WP4 and WP6. These groups provide in-depth technical insights and interact with the aim to justify that the solutions are technically sound and effectively address the identified challenges. In parallel, any ethical considerations or implications associated with compliance with the ethics requirements outlined in D1.2 codeveloped with the Independent Ethics Advisors (IEA) are monitored. If any corrective measures are raised upon concerns from the IEAs these are directly communicated to the use case development and implementation is modified to enforce the compliant processes.

The next step involves the refinement of objectives, user stories, and scenarios. In this step, objectives are refined based on feedback and further analysis, and detailed user stories and scenarios are developed to guide the development and deployment of the use case. Subsequently, use case pipelines that correspond to the scenarios are designed, ensuring the seamless integration of the EMERALDS tools into the existing use case environment (M16-M22). The final step is the validation with use case data, where the developed solutions are validated with real use case data to ensure their effectiveness, reliability, and alignment with the defined objectives. These results will be demonstrated and reported in D5.2, D5.4 & D5.6 (M24) while the second iteration of this approach will produce the final use case results to be reported in D5.3, D5.5 and D5.7 (M33).

Finally, conducting a commercial value assessment of the innovations concurrently with use case testing and validations as well as leading networking and dissemination and communication activities is undertaken through the active link between WP5 and WP6. WP5 results will be utilized to reach a comprehensive understanding of the potential market impact and viability of the EMERALDS toolset and constituent modules. This assessment encompasses the evaluation of various factors, including market demand, competitive landscape, pricing strategies, and potential revenue streams. Moreover, aligning commercial considerations with testing and validation activities will enable the project to iteratively refine its offerings to effectively address market needs and enhance commercial success. Ultimately, early identification of market opportunities and challenges will be strengthened, enabling strategic planning for commercialization.

## 7 Annex – Ethics Checklist

---



### Ethics Checklist and Questionnaire

THIS FORM NEEDS TO BE FILLED-IN BY THE DELIVERABLE LEADER BEFORE, DURING AND AFTER THE WORK LEADING TO THE RELEVANT DELIVERABLE. DELIVERABLE LEADERS ARE ENCOURAGED TO DISCUSS EACH ACTIVE QUESTIONNAIRE WITH THE ETHICS COMMITTEE. EACH OPEN QUESTIONNAIRE SHOULD EITHER BE STORED IN THE PROJECT MS TEAMS DIRECTORY OR A LINK MADE AVAILABLE TO THE RELEVANT SHARED DOC. COMPLETED FORMS NEED TO BE SUBMITTED AS PART OF THE DELIVERABLE Q/A PROCESS. IN THE EVENT OF COMMENTS AND/OR QUESTIONS BY THE ETHICS COMMITTEE, THE DELIVERABLE LEADER HAS TO PROVIDE RELEVANT RESPONSES AND/OR CLARIFICATIONS IN A TIMELY MANNER

#### **A. PERSONAL DATA**

1. Are **personal data** going to be processed for the completion of this deliverable?

**NO**

- If “yes”, do they refer only to individuals connected to project partners or to third parties as well?

2. Are “**special categories of personal data**” going to be processed for this deliverable? (whereby these include personal data revealing racial or ethnic origin, political opinions, religious or philosophical beliefs, and trade union membership, as well as, genetic data, biometric data, data concerning health or data concerning a natural person's sex life or sexual orientation)

**NO**

3. Has the **consent** of the individuals concerned been acquired prior to the processing of their personal data?

**N/A**

- If “yes”, is it based on the Project’s Informed Consent Form, either on the provided Template or on other attached herein Template?
- If “no” is it based on a different legal basis?

4. In the event of processing of personal data, is the processing:

**N/A**

- obviously “**Fair and lawful**”, meaning executed in a fair manner and following consent of the individuals concerned or based on another - acknowledged as adequate and proportionate as per above - legal basis?
  - Performed for a **specific (project-related) cause** only?
  - Executed on the basis of the principle of **proportionality and data minimisation** (meaning that only data that are **necessary** for the processing purposes are being processed and such deductive reasoning is documented)?
  - Based on **high-quality, updated and precise personal data**?
5. Are there any provisions for a storage limitation period of the personal data-in case of storage- after which they must be erased?

N/A

6. Are all **other lawful requirements** for the processing of the data (for example, **notification of the competent Data Protection Authority(s)** or undergoing a **DPIA procedure** and consulting with the competent DPA, if and where applicable) adhered to and on what legislative basis are such notifications justified as necessary or dismissed as unnecessary?

N/A

7. Have individuals been **made aware of their rights** on the processing of the personal data as per the GDPR and the relevant and executive national legislation (particularly the rights to access, rectify and delete the personal data and their right to lodge a complaint with the relevant Competent Authority) and if yes, by what demonstrable means (e.g. the informed consent form as per above or as per other Templates, attached herein?)

N/A

8. Even if anonymized or pseudonymized or aggregated data are referred to, does the dataset contain **location data** that could potentially (even via the combined use of other datasets) be **traced back to individuals**? If yes, what specific measures are taken to ensure this data (i) is anonymized or pseudonymized and (ii) cannot be used to track individuals without their consent? If no, what is the scientific methodology used to collect and gather said data?

N/A

9. In the context of risk assessment, prediction and forecasting, as foreseen in the scope of the EMERALDS project, during traffic, population movement monitoring or weather events, **is there any risk** that personal data could be inadvertently revealed in the event of an **emergency or unusual event**, because of the dataset usage, either on its own or combined with other openly available datasets, triggering identification or unwanted disclosure of PII? What measures are in place to protect - still identifiable if the dataset allows such extraction - **personal data** in these circumstances?

NO

10. For the use case of Trip Characteristics Inference as per the EMERALDS project scope, are there specific measures to ensure that **inferences made about trip characteristics** cannot be linked back to **specific individuals** or reveal **sensitive information** about their **habits** or **routines** i.e. by identifying specific individuals' absence or presence routines whether in the home or in a professional environment or in other premises?

N/A

11. Are there any potentially Personal Data in the datasets, **disclosable by combination with other datasets**, either open data or proprietary (e.g., E-tickets validation data)? If yes, how are Personal Data adequately anonymized or pseudonymized or how other datasets that by combination may result in unwanted or illegal disclosures or identification before any processing takes place?

NO

## **B. DATA SECURITY**

1. Have proportionate security measures been undertaken for protection of the data, taking into account project requirements and the nature of the data?
- If yes, brief description of such measures (including physical-world measures, if any)
  - If yes, is there a data breach notification policy in place within your organization (including an Incident Response Plan to such a breach)?

N/A

2. Given the **large-scale nature** of some datasets, are there specific measures in place to protect **included personal data** at scale at the data source or in the possession of data processors?

N/A

3. Regardless of personal data, in the case of Multi-modal integrated traffic management as defined under the EMERALDS scope, are there specific measures in place to ensure **the availability and integrity of data spanning multiple modes of transport** from being disclosed in other manners than the ones intended and covered under an open data scheme?

N/A

4. Are there specific measures in place to secure **sensitive infrastructure data**, if present?

## **C. DATA TRANSFERS**

1. Are personal data transfers beyond project partners going to take place for this deliverable?

**NO**

- If “yes”, do these include transfers to third (non-EU) countries and if what policies apply?

2. Are personal data transfers to public authorities going to take place for this deliverable?

**NO**

3. Do any state authorities have direct or indirect access to personal data processed for this deliverable?

**NO**

3. Taking into account that the Project Coordinator is the “controller” of the processing and that all other project partners involved in this deliverable are “processors” within the same contexts, are there any other personal data processing roles further attributed to any third parties for this deliverable? And if any, are they conformed to the GDPR provisions?

**NO**

4. Given the geographical diversity of the datasets, are there measures in place to ensure compliance with specific personal data protection regulations **in different jurisdictions** i.e. at the place of the data source establishment as well as at the place of the establishment of a Data processor?

**N/A**

5. Are there additional protocols for data transfers involving **sensitive infrastructure data**, if present?

**NO**

#### **D. ETHICS AND RELATED ISSUES**

1. Are personal data of children going to be processed for this deliverable (ie. “underage” signified e-tickets)?

**NO**

2. Is **profiling** of identifiable individuals in any way enabled or facilitated for this deliverable?



NO

3. Are **automated decisions** for identifiable individuals made or enabled on the basis this deliverable?

NO

4. Have partners for this deliverable taken into consideration system architectures of **privacy by design** and/or **privacy by default**, as appropriate?

YES

5. Have partners for this deliverable taken into consideration gender equality policies or is there an explicit reasoning that dismisses such risk as unsubstantiated or such need as irrelevant as per the methodology of work and production of the deliverable?

**YES, if this question refers to the composition of the team. 4 out of 8 members working D5.1 are female. Three have a non-Dutch background.**

6. Have partners for this deliverable taken into consideration means of protecting the confidentiality of the dataset if it is not signified as open data?

N/A

7. Are there additional considerations around the collection and processing of **location data** and data **that could potentially be used to infer patterns about individuals' movements**?

NO

8. Have partners identified any **additional ethical issues** related to the processing of sensitive infrastructure data?

NO

9. Are shared economy (ie. "Uber" transfer services or "Lime" Scooters or other solution) or other shared mobility infrastructures used by the data sources? If yes, are there measures in place to ensure that the processing of **shared mobility data** respects privacy rights?

NO

10. In the context of Traffic Flow Data Analytics, are there specific considerations to ensure that the **analysis of traffic flow data** does not infringe on privacy rights or reveal sensitive information about individuals' movements or routines?

N/A

11. Is the Project taking into account the need for an all people-inclusive policy in the future within its overall goals and not only the “tech-savvy” (i.e. elderly people not familiar with some tech devices) and does it entail possible proposals for that?

YES

## 8 References

---

- [1] S. P. Hoogendoorn, R. L. Landman, J. Van Kooten, M. Schreuder, R. Adams, "Design and Implementation of Integrated Network Management Methodology in a Regional Network," *Transportation Research Record*, vol. 2489, issue 1, pp. 20-28, 2015.
- [2] F. Schneider, D. Ton, L. B. Zomer, W. Daamen, D. Duives, S. Hoogendoorn-Lanser, S. Hoogendoorn, "Trip chain complexity: a comparison among latent classes of daily mobility patterns," *Transportation*, vol. 48, pp. 953-975, April 2021.
- [3] M. D. Kamruzzaman, J. Hine, T. Yigitcanlar, "Investigating the link between carbon dioxide emissions and transport-related social exclusion in rural Northern Ireland," *International journal of environmental science and technology*, vol. 12, issue 11, pp. 3463-3478, Feb 2015.
- [4] M. Kamruzzaman, K. Deilami, T. Yigitcanlar, "Investigating the urban heat island effect of transit oriented development in Brisbane," *Journal of transport geography*, vol. 66, pp 116-124, 2018.
- [5] T. Yigitcanlar, *Technology and the city: Systems, applications and implications*, New York: Routledge, 2016.
- [6] F. Golbabaei, T. Yigitcanlar, J. M. Bunker, "The role of shared autonomous vehicle systems in delivering smart urban mobility: A systematic review of the literature," *International Journal of Sustainable Transportation*, vol. 15, issue 2, pp. 1-18, July 2020.
- [7] A. Paz, P. Maheshwari, P. Kachroo, S. Ahmad, "Estimation of performance indices for the planning of sustainable transportation systems," *Advances in Fuzzy Systems*, vol. 2, issue 12, pp. 1-13, Jan 2013.
- [8] D. Impedovo, G. Pirlo, "Artificial intelligence applications to smart city and smart enterprise," *Applied Sciences*, vol. 10, issue 8, 944, April 2020.
- [9] D. Luckey, H. Fritz, D. Legatiuk, K. Dragos, K. Smarsly, "Artificial intelligence techniques for smart city applications", in *the 18th International Conf. on Computing in Civil and Building Engineering*, San Paolo, Brazil, 2020, pp. 3-15.
- [10] L. Butler, T. Yigitcanlar, A. Paz, "Smart urban mobility innovations: A comprehensive review and evaluation," *IEEE Access*, vol. 8, pp. 196034-196049, Jan 2020.
- [11] C. Chakhtoura, D. Pojani, "Indicator-based evaluation of sustainable transport plans: A framework for Paris and other large cities," *Transport Policy*, vol. 50, pp. 15-28, Jan 2016.
- [12] B. C. Richardson, "Sustainable transport: analysis frameworks," *Journal of Transport Geography*, vol. 13, issue 1, pp. 29-39, March 2005.
- [13] Z. Toth-Szabo, A. Varhelyi, "Indicator framework for measuring sustainability of transport in the city," *Procedia-Social and Behavioral Sciences*, vol. 48, pp. 2035-2047, 2012.

- [14] R. Papa, C. Gargiulo, L. Russo, "The evolution of smart mobility strategies and behaviors to build the smart city" in *the 5th IEEE International Conf. on Models and Technologies for Intelligent Transportation Systems*, Napoli, Italy, 2017, pp. 409-414.
- [15] E. J. Tomaszewska, A. Florea, "Urban smart mobility in the scientific literature—bibliometric analysis," *Engineering Management in Production and Services*, vol. 10, issue 2, pp. 41-5, 2018.
- [16] T. Yigitcanlar, M. Kamruzzaman, "Smart cities and mobility: Does the smartness of Australian cities lead to sustainable commuting patterns?," *Urban Technologies*, vol. 26, issue 2, pp. 21-46, April 2019.
- [17] M. Z. Naser, A. H. Alavi, "Error metrics and performance fitness indicators for artificial intelligence and machine learning in engineering and sciences," *Architecture, Structures and Construction*, vol. 3, issue 4, pp. 499-517, 2023.
- [18] S. P. Hoogendoorn, J. Van Kooten, R. Adams, "Lessons learned from field operational test of integrated network management in Amsterdam," *Transportation Research Record*, vol. 2554, issue 1, pp. 111-119, Jan 2016.
- [19] M. Haghani, M. Coughlan, B. Crabb, et al., "A roadmap for the future of crowd safety research and practice: Introducing the Swiss Cheese Model of Crowd Safety and the imperative of a Vision Zero target," *Safety science*, 168:106292, Dec 2023.
- [20] Y. Yang, J. Yu, C. Wang, J. Wen, "Risk assessment of crowd-gathering in urban open public spaces supported by spatio-temporal big data," *Sustainability*, vol. 14, issue 10, pp. 6175, May 2022.
- [21] M. A. Hassan, M. U. G. Khan, R. Iqbal, et al., "Predicting humans future motion trajectories in video streams using generative adversarial network," *Multimedia Tools and Applications*, vol. 83, issue 3, pp. 15289-15311, Sept 2021.
- [22] D. Lehmberg, F. Dietrich, G. Köster, "Modeling Melburnians—Using the Koopman operator to gain insight into crowd dynamics," *Transportation Research Part C: Emerging Technologies*, vol. 133, issue 4, pp. 103437, Dec 2021.
- [23] P. Krishnakumari, S. Hoogendoorn-Lanser, J. Steenbakkens, S. Hoogendoorn, "Crowd Safety Manager: Towards Data-Driven Active Decision Support for Planning and Control of Crowd Events," arXiv preprint arXiv:2308.00076, July 2023.
- [24] V. X. Gong, W. Daamen, A. Bozzon, S. P. Hoogendoorn, "Counting people in the crowd using social media images for crowd management in city events," *Transportation*, vol. 48, pp 3085-3119, Jan 2021.
- [25] D. C. Duives, W. Daamen, S. P. Hoogendoorn, "Monitoring the number of pedestrians in an area: the applicability of counting systems for density state estimation," *Advanced Transportation*, vol. 2018, pp. 14, April 2018.
- [26] T. P. van Oijen, W. Daamen, S. P. Hoogendoorn, "Estimation of a recursive link-based logit model and link flows in a sensor equipped network," *Transportation Research Part B: Methodological*, vol. 140, pp. 262-281, 2020.
- [27] G. A. Klunder, H. Taale, L. Kester, S. Hoogendoorn, "The effect of inaccurate traffic data for ramp metering: Comparing loop detectors and cameras using information utility," *IFAC Proceedings Volumes*, vol. 47, issue 3, pp. 11318-11325, 2014.
- [28] G. A. Klunder, H. Taale, L. Kester, S. Hoogendoorn, "Improvement of network performance by in-vehicle routing using floating car data," *Advanced Transportation*, pp. 1-16, Dec 2017.
- [29] M. C. Poelman, A. Hegyi, A. Verbraeck, J. W. C. van Lint, "Sensitivity Analysis to Define Guidelines for Predictive Control Design," *Transportation Research Record*, vol. 2674, issue 6, pp. 385-398, 2020.



- [30] M. C. Poelman, A. Hegyi, , A. Verbraeck, J. W. C. van Lint, "Structure-free model-based predictive signal control: A sensitivity analysis on a corridor with spillback," *Transportation research part C: emerging technologies*, vol. 153, pp. 104174, Sept 2023.
- [31] G. Reggiani, A. Dabiri, W. Daamen, S. Hoogendoorn, "Clustering-based methodology for estimating bicycle accumulation levels on signalized links: a case study from the Netherlands," in *the IEEE Intelligent Transportation Systems Conf.*, Auckland, New Zealand, 2019, pp. 1788-1793.
- [32] G. Reggiani, A. Dabiri, W. Daamen, S. P. Hoogendoorn, "Exploring the Potential of Neural Networks for Bicycle Travel Time Estimation," in *the conf. on Traffic and Granular Flow*, Pamplona, Spain, 2019, pp. 487-493.
- [33] A. Vial, G. Hendeby, W. Daamen, B. van Arem, S. Hoogendoorn, "Framework for Network-Constrained Tracking of Cyclists and Pedestrians," *IEEE Transactions on Intelligent Transportation Systems*, vol. 24, issue 3, pp. 3282-3296, 2022.
- [34] The map of public transport stops in the city of Riga [ONLINE]. Available: <https://saraksti.rigassatiksme.lv/index.html#riga/map,page>
- [35] "Riga traffic" in vehicles 2023 [ONLINE]: Available: <https://www.rigassatiksme.lv/lv/aktu%C4%81la%20inform%C4%81cija/2023-gada-par-23-pieaudzis-rigas-satiksmes-parvadato-pasazieru-skaitis/>
- [36] E. Hussain, A. Bhaskar, E. Chung, "Transit OD matrix estimation using smartcard data: Recent developments and future research challenges," *Transportation Research Part C: Emerging Technologies*, vol. 125, art. number 103044, April 2021.
- [37] O. Egu, P. Bonnel, "How comparable are origin-destination matrices estimated from automatic fare collection, origin-destination surveys and household travel survey? An empirical investigation in Lyon," *Transportation Research Part A: Policy and Practice*, vol. 138, pp. 267-282, Aug 2020.