



Intelligent Urban Traffic Management via Semantic Interoperability Across Multiple Heterogeneous Mobility Data Sources

Mario Scrocca¹(✉) , Marco Grassi¹ , Marco Comerio¹ ,
Valentina Anita Carriero¹ , Tiago Delgado Dias² , Ana Vieira Da Silva² ,
and Irene Celino¹

¹ Cefriel – Politecnico di Milano, Viale Sarca 226, 20126 Milan, Italy
{mario.scrocca,marco.grassi,marco.comerio,valentina.carriero,
irene.celino}@cefriel.com

² A-to-Be Mobility Technology, S.A. Edificio Brisa, 2785-599
São Domingos de Rana, Portugal
{tiago.delgado.dias,ana.vieira.silva}@a-to-be.com

Abstract. The integrated exploitation of data sources in the mobility domain is key to providing added-value services to passengers, transport companies and authorities. Indeed, multiple stakeholders operate and maintain different kinds of data but several interoperability issues limit their effective usage. In this paper, we present an architecture enabled by Semantic Web technologies to overcome such issues and facilitate the development of an integrated solution for mobility stakeholders. The proposed solution is composed of different components that address challenges for enabling data interoperability, from the findability of data sources to their integrated consumption adopting standardised data formats. We report on its implementation and validation in four European cities to enable data-driven tools for the dynamic management of multimodal traffic. Finally, we discuss the feedback received by users testing the solution and the lessons learnt during its development.

Keywords: Mobility · Semantic interoperability · Data Integration

1 Introduction

Interoperability is one of the main challenges in enabling collaboration between travel and transport industry players. The interoperability of data and services is essential for creating an ecosystem of transport stakeholders, enabling the definition of new data-driven solutions. The benefits range from enhancing the operations of transport companies to providing integrated and seamless mobility services to users. The TANGENT project¹, co-funded by the European Commission under the Horizon 2020 Programme, is developing new complementary tools

¹ Enhanced Data Processing Techniques for Dynamic Management of Multimodal Traffic (TANGENT), <https://doi.org/10.3030/955273>.

for optimising traffic operations in a coordinated and dynamic way from a multi-modal perspective. In this paper, we discuss how knowledge graphs and semantic technologies can help in tackling interoperability in the mobility domain.

The development and testing of the data-driven solutions developed by TANGENT in the Athens, Lisbon, Greater Manchester, and Rennes Metropole case studies asked for a solution to different data interoperability challenges. The definition of a proper solution for data sharing and usage is not straightforward due to several issues to be addressed: datasets in different formats and/or using different data models, data services relying on different specifications and technologies, and metadata describing data sources according to different profiles. This paper describes the design and implementation of an integrated set of tools that employs Semantic Web technologies to address data interoperability issues for heterogeneous data sources from different stakeholders.

While the proposed solution can be applied to a generic domain, we describe how we implemented and tested it considering the specificities of multimodal transportation. Furthermore, we describe the reusable resources, like metadata specification and ontologies, that are made publicly available.

Within the integrated TANGENT solution, data interoperability is leveraged to support the deployment of an innovative dynamic traffic management platform. The platform provides intelligent services and cutting-edge user interfaces and is enabled in each city by the integrated consumption of several data sources. The data sources were retrieved and harmonised involving different stakeholders and systems, thus demonstrating the flexibility and scalability of the proposed solution. We received positive feedback from users during testing sessions that acknowledged the benefits of integrating data for different transportation modes within a single solution and expressed their interest in adopting the solution within their operations.

The remainder of the paper is organised as follows. Section 2 discusses the motivating challenges and the related work. Section 3 describes the proposed solution and its implementation considering the heterogeneous mobility data sources of the four cities involved in the validation. Section 4 discusses the evaluation and lessons learned. Finally, Sect. 5 draws the conclusions.

2 Challenges and Related Work

Data interoperability is a challenging objective to enable different stakeholders to communicate and exchange information effectively without losing meaning. Indeed, stakeholders adopt different (legacy) systems for data management and exchange that cannot be directly integrated or harmonised. To better understand the problem within the considered domain, we first investigated the existing products and services adopted by the case studies involved in TANGENT. The current landscape is characterised by: (i) open data portals managed at different levels (regional, national, European) that contain datasets and data services not well-documented (i.e., bad quality metadata), provided in not interoperable data formats (e.g., custom CSVs) and often not updated because not directly used

by the relevant stakeholders; (ii) solutions from third-party vendors (usually associated with the vendor of the sensors generating the data) that deal with specific sets of data (e.g., road traffic data), use custom data formats and do not provide easy access (e.g., API) to the raw data. Both of these options restrict the ability to find relevant data sources, especially considering non-public data, and to access pertinent information for comprehensive traffic management.

Five major challenges can be identified and should be addressed [5]: 1. **Locate** (*which data is available and where?*), 2. **Access** (*how to obtain the needed data?*), 3. **Harmonise** (*how to convert data according to the required data model?*), 4. **Integrate** (*how to ensure different data sources can be merged?*), 5. **Extract** (*how to consume harmonized and integrated data?*).

Each of these challenges is associated with several issues and identifying a single solution is impossible since a single interoperability problem cannot be formulated. Indeed, data interoperability scenarios are widely heterogeneous and pose various requirements [24] that can be possibly faced only by considering a set of tools appropriately configured. To select such tools and define an integrated solution [4], we reviewed state-of-the-art data interoperability solutions based on Semantic Web technologies and their application for the mobility domain.

2.1 Locate and Access

The first challenge is the findability and discoverability of data. Data cannot be re-used and (made) interoperable if they cannot be found. For this reason, data catalogues/portals are implemented to describe data sources through a set of metadata. The challenge is associated with the need for a proper, structured and machine-readable description of data sources that could also support interoperability across different data catalogues. Once data sources are located, the second challenge is related to data accessibility. Data catalogues adopt different strategies for data access mainly associated with the architectural choices for the hosting and storage of static and dynamic data sources. The challenge is to enable uniform access to heterogeneous data sources for end users.

The *locate* and *access* challenges are also being addressed by the European Commission through National Access Points (NAP) for mobility data. Each Member State should operate a NAP to enable the sharing of mobility data by transport stakeholders as mandated by dedicated Delegated Regulations [7–9] supplementing the Intelligent Transport Systems (ITS) Directive 2010/40/EU [21]. The concept of NAP leverages the one of Data Catalogue, i.e., a digital platform to facilitate the sharing of data sources and their findability by other stakeholders. However, several mobility data platforms exist but are not interoperable. Even in the case of NAPs, each Member State adopted different approaches for their implementation, thus creating interoperability issues at the European level. For this reason, the NAPCORE² project is currently working on coordinating and harmonizing such platforms around Europe. One important objective is supporting the findability of data contained in each mobility data

² <https://napcore.eu/>.

platform [25] and defining mobilityDCAT-AP³, a uniform metadata specification to access the data sources. The adoption of structured metadata descriptors according to well-known vocabularies, e.g., the Data Catalog Vocabulary (DCAT) [1] and the corresponding DCAT Application profile (DCAT-AP) for data portals in Europe [29], is fundamental to facilitate search within one or multiple data catalogues. Moreover, proper data governance must be defined to regulate the usage of the catalogue between the different involved stakeholders. Finally, data catalogues should support the harmonisation of technological access to data sources.

For these reasons, we identified the two core components of the proposed solution as a shared *Data Catalogue* to enable the findability of data sources and a uniform *Data API* for accessibility.

2.2 Harmonise, Integrate, and Extract

The remaining three data interoperability challenges (harmonise, integrate, and extract) are related to the processing of (meta)data to enable their integration and exploitation according to common semantics. A flexible solution is required to address heterogeneous requirements in terms of: (a) **schema and data transformation**: information manipulation to obtain syntactic (structural) and semantic interoperability of (meta)data; (b) **integration** with existing information systems as data sources (i.e., components generating or storing the data) and/or data sinks (i.e., components consuming or archiving the data).

Different approaches can be exploited and implemented, spanning from ad-hoc solutions targeting a specific scenario to more general and scalable solutions supporting multiple stakeholders and data representations. The semantic any-to-one mapping approach based on [30] and validated in [27] reduces the number of mappings, i.e., translations from one representation to another, that are needed to implement interoperability by different stakeholders. Such an approach is based on the identification of a reference model for the domain of interest. Each stakeholder is responsible for defining mappings from their own data representation to the reference model (*lifting*) and vice versa (*lowering*). In this paper, we discuss how we adopted this approach.

Considering the mobility domain, different reference models are proposed based on existing standards. Chouette [12] and the SNAP solution [23] rely on a reference model based on Transmodel⁴ for the conversion of Public Transport (PT) data in different formats. The `transit_model` tool⁵ adopts the Navitia Transit Feed Specification (NTFS)⁶ to manage, convert and enrich transit data from/to different formats. Moreover, considering traffic data, the Datex II⁷ specification is often used as a reference model to convert custom data formats and share harmonised data [14].

³ <https://w3id.org/mobilitydcat-ap>.

⁴ <https://www.transmodel-cen.eu/>.

⁵ https://github.com/hove-io/transit_model.

⁶ <https://github.com/hove-io/ntfs-specification>.

⁷ <https://www.datex2.eu/>.

Different approaches for lifting and lowering can be suitable considering a specific scenario. Moreover, the harmonisation, integration and extraction process may require the definition of custom pre- and post-processing, considering different interoperability issues. Therefore, composing and configuring different components should be possible considering the specific requirements for integrating certain data sources.

Different approaches for lifting and lowering can be suitable based on a specific scenario. Moreover, the harmonisation, integration, and extraction process may require the definition of custom pre- and post-processing, taking into account different interoperability issues. Therefore, it should be possible to compose and configure different components while considering the specific requirements for integrating certain data sources.

Different semantic-based ETL (Extract, Transform and Load) tools have been proposed to define composable procedures with Semantic Web technologies [13]. Technologies for declarative knowledge graph construction [28] can effectively support lifting transformations, while a standardised lowering solution to convert RDF to any format using a generic declarative language is currently missing [26]. Moreover, other components, such as message filtering or routing, are usually required within a transformation pipeline. Enterprise Integration Patterns [16] offer a relevant categorisation of the components and techniques for system integration.

In conclusion, two additional components are identified to support the solution: a *Reference Conceptual Model* defining common semantics and the composition and configuration of *Semantic Harmonisation and Fusion Pipelines*.

3 TANGENT Solution for Dynamic and Intelligent Multimodal Traffic Management

This section describes the TANGENT solution for dynamic and intelligent multimodal traffic management. We discuss each macro-component of the proposed architecture and then its integration within the overall TANGENT solution [19] as shown in Fig. 1. The main original contributions of this work are the implementation and integration of different technologies to propose an holistic architecture for data interoperability based on Semantic Web technologies and its customisation for the multimodal traffic management domain. In the following we describe them, highlighting their value in solving the discussed challenges and their impact on the business scenarios.

3.1 Data Catalogue

This component is a catalogue of digital assets available online and accessible by users via a web browser. Different digital assets can be characterized by specifying a metadata descriptor according to a common metadata profile. The data catalogue may also harvest metadata descriptions from existing data portals (e.g., NAPs) [3]. Programmatic access to the list of assets published and

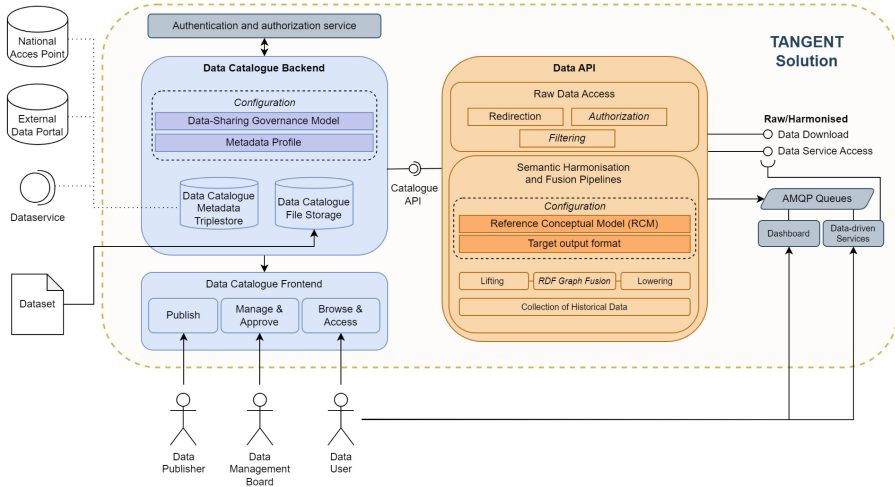


Fig. 1. Overview of the proposed solution for data interoperability

their metadata is implemented via a dedicated API (*Catalogue API*). Metadata serialized in RDF enables advanced functionalities based on querying and/or automated processing by agents of such metadata (e.g., in federation scenarios). Finally, the catalogue enforces processes for the governance of digital assets.

The *Data Catalogue* provides a single location where all the collected data sources for each city are described and can be explored. Its development is based on the Knowledge Catalog and Governance (KCONG) framework⁸, developed by Cefriel, which is customised considering the TANGENT *Data Sharing Governance Model* [2,6] and the TangentDCAT-AP metadata specification as metadata profile.

The *Data-Sharing Governance model* supports the strategic and operational management of data-sharing. The model focuses on the tasks/processes related to providing access to data sources needed by the various technical components. It identifies key processes (data publication, data quality, data access, data storage, data usage) to be addressed by stakeholders with different roles and following specific rules. The identified roles are: (i) *Data publisher*, a person responsible for publishing and describing a data source within the catalogue; (ii) *Data Management Board (TMB)*, a group of people responsible for the management and control of a (set of) data source(s) within the catalogue; (iii) *Data user*: a person accessing and using a data source available in the catalogue. As an example of the defined rules (fully described in [6]), the Lisbon case study leader acts as a data publisher and can create and modify only the metadata descriptions of data sources related to the Lisbon case study.

The definition of the TangentDCAT-AP metadata specification considered best practices, particularly the reuse of well-known vocabularies for metadata.

⁸ <https://kcong.cefril.com/>.

For this reason, TangentDCAT-AP is defined as an extension of DCAT-AP [29], considering the requirements for mobility data platforms elicited by the NAP-CORE project [25] and specific requirements for the description of data sources elicited within the TANGENT project [6]. The final release of TangentDCAT-AP is compatible with the first official release of mobilityDCAT-AP⁹ by the NAPCORE project. The mobilityDCAT-AP specification extends DCAT-AP by focusing on requirements for data sources in the mobility domain. It will be recommended as the reference metadata profile for National Access Points¹⁰ and adopted for the European Mobility Data Space¹¹.

TangentDCAT-AP is documented at <https://knowledge.c-innovationhub.com/tangent/tandcatap>. Following the best practices [25], the additional properties defined by TangentDCAT-AP have been also published online together with the defined controlled vocabularies. The published vocabularies are hosted on GitHub¹² and served through content negotiation. The vocabularies define possible statuses assigned to a data source, data requirements to categorise the content of data sources and their types.

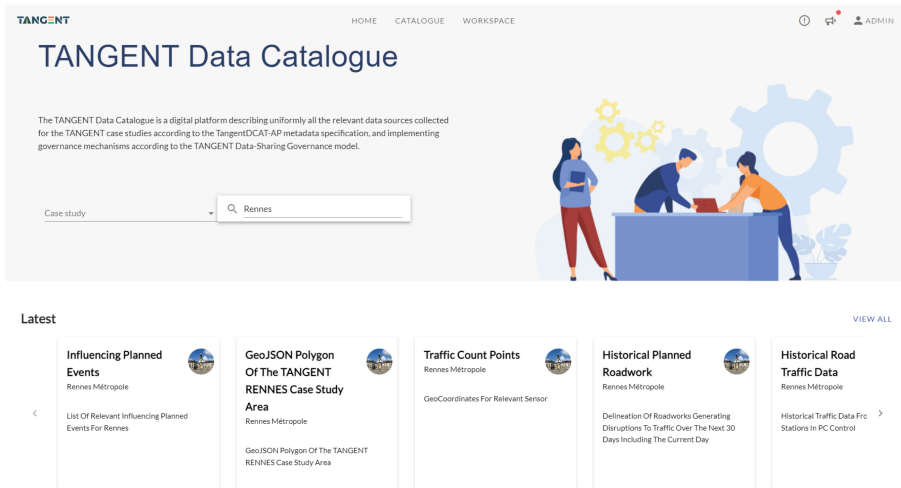


Fig. 2. Homepage of the Data Catalogue.

Figure 2 shows the homepage of the Data Catalogue and some available data sources. Three main areas are available to the user: the homepage, the catalogue and the workspace. The homepage is designed to show the latest changes and the most common interactions with the catalogue, i.e., searching. The catalogue

⁹ <https://w3id.org/mobilitydcat-ap/releases/1.0.0/>.

¹⁰ <https://napcore.eu/release-of-the-mobilitydcat-ap/>.

¹¹ <https://www.linkedin.com/company/deployemds>.

¹² <https://github.com/cefriel/tandcatap>.

area allows the user to browse all available data sources and perform more fine-grained searches. Lastly, the workspace area allows visualising/editing a specific data source, or getting an overview of the state of the data sources owned by the current user. An external identity and access management solution is used to authenticate and authorize different users to access data sources.

Currently, the Data Catalogue contains metadata about 145 data sources.

3.2 Reference Conceptual Model

The *Reference Conceptual Model* (RCM) supports the representation of heterogeneous information from different data sources through a common ontological model to enable shared semantics and interoperability. The model is based on existing data standards to adopt the correct domain terminology and covers the representation of all the entities and properties required to implement meaningful data exchanges among the involved stakeholders.

In the mobility domain, several ontological models have been proposed. However, they cover specific requirements, and it is not possible to identify a generic and well-adopted ontology [18]. For these reasons, we started our work by analysing existing data standards to support the identification of the relevant semantics. Then, following the best practice of reusing existing models [11], relevant ontologies to be included in the RCM were identified.

We started with the analysis of data standards requested by the European Commission (EC) Delegated Regulations (DR) mentioned in Sect. 2.1. The data standards mentioned in those directives are: (i) *DATEX II* (<https://www.datex2.eu/>): the EU standard for the exchange of traffic-related data; (ii) *NeTEx* (<https://netex-cen.eu/>): the CEN Technical Standards for exchanging Public Transport schedules and related data; (iii) *SIRI* (<https://www.siri-cen.eu/>): the CEN technical standard for the exchange of real-time information about the planned, current, or projected performance of public transport operations.

To support the definition of the RCM, the existing ontologies encoding the semantics of the mentioned standards have been analysed.

Considering the DATEX II format, two different models have been identified. The first one, directly developed by the DATEX II organization [17], is a JSON-LD serialisation of the DATEX II conceptual model version 3, divided into five main modules¹³ : *Payload*, *Common*, *LocationReferencing*, *Situation*, *Road Traffic Data*, *Variable Message Sign*. The second model [15] based on DATEX II was developed by the LOD-RoadTran18 project to support the publication of DATEX II data as Linked Open Data. The advantage of the first model is its full coverage with respect to the DATEX II specification, the one-to-one mapping to classes and properties, and the fact that it is directly defined and published by the DATEX II organization; however, this model is defined as an almost automatic conversion of the DATEX II specification. The advantage of the second

¹³ Models used available at <https://datex2.eu/vocab/3>. Additional modules are now available to accommodate the new versions of the DATEX II specification.

Table 1. Overview of the TANGENT Reference Conceptual Model

Module	Base Standard	Data Requirements	Reused ontologies
Road Transport Network	Datex II	Road Transport Network (roads, limited access zones, etc.)	GeoSPARQL, Basic Geo, Datex II JSON-LD (location, common)
Road Equipment	Datex II	Road Equipment Position	Datex II JSON-LD (location), DC Terms, Schema.org
Road Traffic Data	Datex II	Road Traffic Measurements (traffic occupancy, speed, flow) Floating Vehicle Data (GPS, mobile, etc.)	Datex II JSON-LD (traffic)
Road Travel Times	Datex II	Road Travel Times (external services, statistics, etc.)	Datex II JSON-LD (location, traffic)
Events	Datex II	Road Transport Network Events (planned) / Incidents (unplanned) Influencing Planned Events (sports, entertainment, etc.) Weather Events	Datex II JSON-LD (situation, location, common)
Weather Data	Datex II	Forecasted Weather Data Weather Data (measurements, e.g., temperature, humidity, etc.)	Datex II JSON-LD (location, traffic, common)
Stop Points	NeTEx	Public Transport Network	Transmodel ontology (commons, journeys), Basic Geo
Schedules	NeTEx	Public Transport Schedules and Lines	Transmodel ontology (commons, journeys, organisations)
Situation Exchange	SIRI	Public Transport Network Events (planned) / Incidents (unplanned)	Basic Geo
Vehicle Monitoring	SIRI	Floating PT Vehicle Data Public Transport Delays	Basic Geo

model is that the LOD-RoadTran18 project followed proper ontology engineering methodologies to define the model; however, it covers only a portion of the DATEX II specification. For these reasons, we selected the first model.

The NeTEx and SIRI standards are based on the Transmodel¹⁴ conceptual model. The Mobility Ontology Catalogue¹⁵ defines a suite of ontologies based on existing standards, including a Transmodel ontology. The Transmodel ontology, firstly defined within the SNAP¹⁶ project, has been extended and reviewed over the years [23] and currently defines five submodules: *Core*, *Commons*, *Fares*, *Facilities*, *Journeys*. The Transmodel ontology does not cover the entire Transmodel but was used to effectively support mappings to NeTEx [27]. To the best of our knowledge, a dedicated SIRI ontology does not exist. However, since SIRI is derived from Transmodel, the current Transmodel ontology can be exploited to represent common concepts and possibly extended to represent the missing ones. We developed and published a dedicated ontology representing concepts and relations mapped from SIRI, adopting an approach aligned to the one leveraged for the definition of the DATEX II JSON-LD ontology from the DATEX II specification. Moreover, we manually curated the SIRI ontology to improve the alignment with the other ontologies adopted in the RCM. The SIRI ontology is available and documented at <https://knowledge.c-innovationhub.com/siri>. The current version of the ontology (v1.0.0) focuses on the modelling of situations affecting the transport network and monitored vehicle journeys.

¹⁴ <https://www.transmodel-cen.eu/standards/>.

¹⁵ <https://w3id.org/mobility>.

¹⁶ <https://snap-project.eu/>.

Based on the performed analysis, the RCM was defined as a suite of ontologies considering the semantics of relevant EU-mandated standards and the already available related ontologies. The DCI Metadata Terms¹⁷ and Schema.org¹⁸ vocabularies are reused by the Transmodel ontology and similarly also in other RCM modules. Table 1 provides a complete overview of the ten different modules defined for the *Reference Conceptual Model*, summarising the considered base standard for concepts and relationships, the data requirements covered by the module, and the reused ontologies.

The latest release of the RCM is published online¹⁹ including the documentation of the different modules and the serialisation of the additionally defined classes/properties. The definition of the RCM has been guided by the requirements elicited in TANGENT for the harmonisation and fusion of data sources for multimodal traffic management. Nevertheless, the RCM is made available and we recommend its reuse and extension to address additional requirements.

3.3 Semantic Harmonisation and Fusion Pipelines

Transformation pipelines should support the harmonisation and fusion of data sources available by leveraging the *Reference Conceptual Model*. The definition of such pipelines requires the elicitation of harmonization and fusion requirements based on: (a) the analysis of raw data sources stored in the *Data Catalogue*, and (b) the definition of the information required by the other components that are integrated into the solution and the related target output format. We mainly addressed the requirements of two downstream data usage: the need for large-scale historical data for training of machine learning models to support traffic management, the need to access static and (quasi) real-time data to empower a set of innovative applications for traffic managers.

In both cases, the need to overcome data heterogeneity and sparsity requires transformation pipelines. The basic pipeline is composed of a lifting operation, a (set of) graph operations (e.g., to perform data fusion), and a lowering operation.

A flexible and scalable technology for the implementation of the pipelines should provide (i) a set of reusable building blocks that can be configured according to specific requirements, and (ii) a declarative approach to configure the lifting and lowering transformations without developing ad-hoc and hard-to-maintain solutions. Chimera²⁰ [13] is an open-source solution based on Apache Camel to enable the definition of semantic data transformation pipelines with different components for knowledge graph construction, transformation, validation, and exploitation. The advantage of Chimera is its integration with Apache Camel, providing off-the-shelf and production-ready components to implement Enterprise Integration Patterns and to integrate pipelines with heterogeneous systems (e.g., HTTP API, WebSocket, MQTT). For these reasons, Chimera was

¹⁷ <http://purl.org/dc/terms/>.

¹⁸ <https://schema.org/>.

¹⁹ <https://knowledge.c-innovationhub.com/tangent/schema>.

²⁰ <https://github.com/cefriel/chimera>.

selected to implement the semantic any-to-one mapping approach and Apache Camel is leveraged to implement the Data API and smoothly integrate the defined pipelines.

The data transformation from a source format and source semantics to a target ontology, i.e., the lifting process, can be handled by either the RML Component or the Mapping Template Component of Chimera. The lifted data, in the form of an RDF graph aligned with the Reference Conceptual Model, can be manipulated using the operations defined by the Chimera Graph Component. For example, these operations enable the fusion of data and/or the filtering/extraction of certain information. The Mapping Template Component handles the lowering process from RDF to the target data format²¹. For the TANGENT pipelines, we decided to use the Mapping Template Component since it can be applied for both lifting and lowering [26].

The definition of harmonisation and fusion requirements for the solution required the analysis of (i) raw data sources collected in the Data Catalogue, and (ii) the information required by the other components to be integrated into the overall solution. As a result of a collaborative effort among partners in charge of developing downstream data-driven services, we identified the requirements in terms of target output format and the set of data sources to be harmonised. Considering the target output, we identified the need for a harmonised CSV format to support the training of traffic prediction machine learning models, and of JSON schemas to feed the real-time services at runtime via AMQP²² queues. The semantics of the RCM was used as a basis to support both the annotation of columns in CSV and of fields in JSON²³.

All in all, we were able to harmonize and fuse data from 43 data sources adopting different data formats (mainly CSV, XML, JSON) and more than 30 data models across the 4 urban case studies. Indeed, many data sources were based on custom data models, thus requiring dedicated lifting mappings, and we could reuse them only for GTFS²⁴ feeds and data from the same data provider. On the other hand, we could leverage the any-to-one approach to define a single lowering mapping for each target output (10 lowering mappings to JSON schemas). Additionally, via dedicated pipelines, we generated 8 historical datasets targeting a CSV format. These pipelines are based on the same lifting mappings but are configured to regularly (e.g., 1-minute frequency) collect, harmonise and fuse data from real-time data sources. We run these pipelines for over 6 months collecting 92GB of compressed data.

²¹ Regarding the lowering, we initially investigated the possibility of applying JSON-LD frames [20] to convert the RDF Graph to the target JSON Schemas. However, we encountered difficulties in addressing cases in which the structure of the RDF graph does not directly correspond to the structure of the target JSON.

²² <https://www.amqp.org/>.

²³ JSON Schemas are available at <https://github.com/cefriel/tangent-model> and contain a simplified representation of the information modelled in RDF to minimise the amount of exchanged data; if the data should be consumed as JSON-LD, a proper context can be associated with each JSON message.

²⁴ <https://developers.google.com/transit/gtfs/>.

3.4 Data API

The *Data API* provides uniform access to data sources collected through the Data Catalogue and represents the integration point for the overall integrated solution. The Data API aims at solving access issues for both *raw data sources* (raw data as collected and shared by the data publisher) and *harmonised and/or fused data sources* (data produced as the result of a semantic harmonisation and fusion pipeline). Harmonised data sources are represented in the Data Catalogue as different *Distributions*²⁵ of the same data source, i.e., different serialisations of the same information. The same approach is also used for data sources provided in multiple raw formats, e.g., CSV and JSON. The result of a fusion process is instead added to the catalogue as a new record since it represents a new data source. The Data API gives access to two main types of data sources: (i) *datasets* usually directly downloadable from a specific URL, and (ii) *data services* implementing different interaction mechanisms (e.g., a REST API). The Data API implements the API Gateway pattern [22], thus providing a single and coherent entry point for the final user. The user can access a data source by knowing the endpoint at which the Data API is located and the identifier of the data source to be accessed. If the user is authorized through the Data Catalogue, the Data API handles authorization mechanisms for the different data sources in a transparent way and provides access to them. Moreover, in cases where a data source should be filtered according to specific requirements, the Data API can be configured to provide access only to the relevant data (e.g., adding the proper parameters to filter data according to a defined temporal/geographical scope). A dedicated parameter can be used to request a specific distribution of the data source (e.g., the harmonised format). The Data API was implemented using the Apache Camel framework since: (i) it provides all the relevant components to interact with the different data platforms (NAPs, Open data/private portals, etc.) hosting data sources, (ii) it can be easily integrated with the semantic harmonisation and fusion pipelines.

The integration of an external data source published on the Data Catalogue within the Data API requires different steps. First of all, we perform an analysis of the accessibility metadata provided for the Dataset (access URL/download URL) or the Data service (endpoint URL/endpoint documentation). If the data source is not available online, we contact the responsible stakeholder to get access. In this case, we leverage the storage layer of the Data Catalogue to upload the data and we evaluate the need to define a data service. If the data source is already available online, we investigate the expected data access interaction (e.g., used protocol); in particular, we evaluate authorization and authentication mechanisms. We configure the required Apache Camel component to retrieve the data source (e.g., HTTP component), and define the integration logic (e.g., filtering the data of a data source considering the relevant temporal/spatial scope for the case study). Finally, if available, we integrate the semantic harmonisation and fusion pipeline to allow users to request data in their harmonised format.

²⁵ <https://www.w3.org/ns/dcat#Distribution>.

Once a data source is integrated into the Data API, the corresponding metadata to access it are updated in the Data Catalogue. In total, we developed 87 integrations to provide access via the Data API to all the raw/harmonised data sources approved for usage.

4 Evaluation and Lessons Learned

This section discusses an evaluation of the proposed solution considering various perspectives and the lessons learned.

Technical evaluation. The integrated on-cloud deployment consisted of testing and production environments with a cluster of four virtual machines and a set of managed services (e.g., database and load balancer). The runtime data for the four case studies are handled by 32 AMQP message exchanges fed by the Data API. The current data are refreshed with different periodicity depending on the considered data source and are organised with static data in 91 distinct collections on the database. The solution processes 130 messages per minute at peak time and manages around 4 GB of data at a time only to visualise the network's current status. The implemented solution demonstrates the feasibility and advantages of adopting an approach for data harmonisation and fusion based on Semantic Web technologies also to support real-time visualisations and data-driven services in a production-ready environment. The semantic any-to-one approach supported good scalability considering the high number of data sources involved and their substantial heterogeneity. Moreover, we demonstrated how it is possible to define system integrations that seamlessly combine semantic harmonisation and fusion to enable interoperability of data exchanges. We also positively assessed that the executed transformations introduced negligible latency (in the order of milliseconds) considering the update frequency of the data sources (often in the order of minutes). Finally, we highlighted the possibility of leveraging the same solution to generate datasets for training machine learning models. A pipeline defined for harmonisation could be easily configured to collect and aggregate data from real-time data, thus generating a historical dataset of harmonised data. Such an approach, relying on common semantics, reduces enormously the effort needed by data scientists to assess heterogeneous input datasets and facilitates the reuse of training algorithms (e.g., for different cities considering different data sources).

User evaluation. To gather opinions and feedback from real users, we performed a qualitative evaluation involving 43 business stakeholders such as transport operators and authorities. They highlighted the importance of a Data Catalogue with structured metadata descriptions to reduce the current scattering of information on data sources and facilitate their retrieval. Indeed, we experienced this difficulty ourselves during the data collection phase: retrieving updated and consistent information about data sources often required the involvement of different people within the same company and/or third-party solution vendors. Based on the evaluation feedback, we also improved the Data Catalogue by

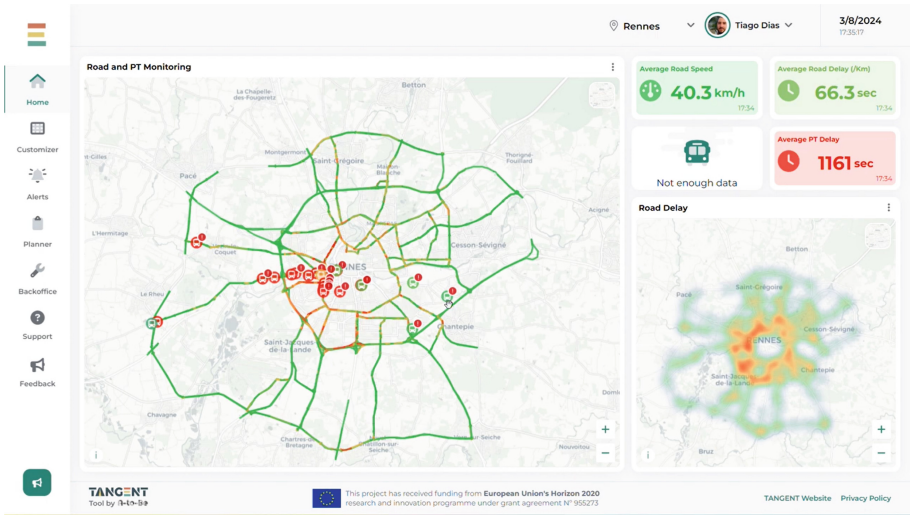


Fig. 3. TANGENT Dashboard for Rennes.

customising filtering operators to further facilitate data discovery. In parallel, we performed the testing of the integrated TANGENT solution (cf. Fig. 3) for the visualization of the current and forecast status of the multimodal network for city and transport authorities, integrating different data-driven services for intelligent incident detection, response plan prescription and the cooperative management of incidents [10]. The stakeholders involved in its evaluation highlighted the advantage of having data from different transport modes available on the same platform in integrated visualizations. The final detailed assessment of the case studies will be included in the future TANGENT deliverables D7.3-D7.6 (due Q4 2024).

Significance Evaluation. The advantages of the proposed solution over potential alternatives are: (i) the collection of structured and actionable metadata in a common machine-readable format that facilitates the findability of heterogeneous data source through the Data Catalogue, and (ii) the reduction of the integration effort for downstream interoperable data-driven services, through the Semantic Harmonisation and Fusion Pipelines and the subsequent Data API: we removed the need for point-to-point integration between different components and we reduced at minimum the custom development for the different deployments in each urban case study.

Uptake and Impact Evaluation. The current uptake of the solution is confined to the involved stakeholders in the four urban case studies, which, however, represent a significant sample of the target market of traffic management solutions and were able to test the solution in their daily operations over an extended period. Indeed, they pointed out the absence of a similar solution in their existing infrastructures and its relevance for both internal governance and traffic management. Moreover, an analysis of the possible exploitation in other cities and the integration of the presented solution with the commercial

offering of A-to-Be are currently ongoing. Concerning the potential impact, we defined a Reference Conceptual Model based on the semantics of existing standards that could be adopted (and possibly extended) to foster interoperability of different solutions for traffic management. Similarly, we demonstrated how to define a metadata extension that supports compatibility with other data portals (leveraging mobilityDCAT-AP) but fulfils additional requirements. Finally, the technological solutions developed for traffic management are not dependent on the mobility domain: with the opportune use of metadata specifications, domain ontologies and the definition of specific mappings, the presented components can be easily configured and adapted to any other domain or market, bringing the same interoperability advantages.

Lessons Learned. We now discuss some of the lessons learned referencing the five challenges in Sect. 2. Regarding the *Locate* challenge, we experimented the difficulty of obtaining good quality metadata. Structured descriptions of existing data sources are often not available, and the Data Catalogue not only enforced common metadata descriptors but also provided users with useful guidelines to collect high-quality metadata (e.g., guided and dynamic forms for metadata insertion). Concerning the *Access* challenge, we experienced the advantage of having an integrated solution for data findability and access. Indeed, often data portals act as simple metadata catalogues, only referencing existing data sources and without taking into account difficulties in accessing the actual data (e.g., authorization, missing documentation of data services, etc.); our Data Catalogue also incorporates and hides the complexity of data harmonisation, giving access to a uniform data API. Considering the *Harmonise* challenge, it is often hard to define common semantics to support different use cases. For this reason, it is important to leverage the semantics already encoded in existing standards without reinventing the wheel and to facilitate the adoption of ontologies by domain experts. Regarding the *Integrate* challenge, it is often difficult to integrate certain types of data (e.g., geographical data considering different location referencing methods or identifiers of stop stations across different transport modes). In these cases, we managed to implement complex transformations by defining mapping rules that integrate custom functions, and by leveraging data fusion with external data sources that specify the correct correspondence between values. Finally, concerning the *Extract* challenge, we demonstrated that a reference ontology supports the generation of harmonised outputs also in formats different from RDF (in our case, CSV for model training and JSON schemas for runtime interactions). This not only facilitates the integration with systems unable to process RDF, but also reduces the size of exchanged data (and, consequently, latency) by avoiding possibly verbose RDF representations.

5 Conclusions

This paper has comprehensively described the proposed solution to address data interoperability challenges within a complex traffic management scenario and its validation within four European cities. The solution guarantees interoperable

descriptions of the data sources and applies the any-to-one centralized approach for semantic interoperability, enabling data exchange with unambiguous and shared meaning.

The proposed solution consists of four components: the Data Catalogue, for sharing uniform data source descriptions according to the TangentDCAT-AP metadata profile and for enforcing governance mechanisms; the Reference Conceptual Model, a reference ontology defining common semantics for multimodal traffic data; the Semantic Harmonisation and Fusion Pipelines, to fulfil heterogeneous data integration requirements; and the Data API, a uniform mechanism to access all data sources. Semantic Web technologies proved their efficacy in addressing data interoperability challenges and providing a production-ready to integrate data in downstream services and applications.

Supplemental Material Statement: Public deliverables are available on the TANGENT website at <https://tangent-h2020.eu/deliverables/>. Implementation reports and related artifacts (e.g., source code) for the TANGENT solution are part of confidential deliverables and can not be shared. The Reference Conceptual Model is published online at <https://github.com/cefriel/tangent-model>, TangentDCAT-AP and controlled vocabularies at <https://github.com/cefriel/tandecatap>. The Chimera framework used for the implementation of the pipelines is available at <https://github.com/cefriel/chimera>. A video describing the TANGENT solution can be visualised at <https://youtu.be/lrwu79lIx4k?feature=shared>.

Acknowledgments. The presented research was partially supported by the TANGENT project (Grant Agreement 955273), which is funded by the European Commission under the Horizon 2020 Research and Innovation Programme. The authors would like to acknowledge Antonia Azzini, Alessio Carenini, Andrea Fiano, and Gianluca Gizzi from Cefriel for their contribution to the development of the TANGENT solution.

References

1. Albertoni, R., Browning, D., Cox, S., Gonzalez Beltran, A., Perego, A., Winstanley, P.: Data Catalog Vocabulary (DCAT) - Version 2. <https://www.w3.org/TR/vocab-dcat-2/>
2. Azzini, A., Comerio, M., Metta, S., Scrocca, M.: The TANGENT governance model for mobility data sharing. In: Transport Transitions: Advancing Sustainable and Inclusive Mobility - Proceedings of the 10th TRA Conference. LNM, Springer (2024), (to appear)
3. Carenini, A., Fiano, A., Scrocca, M., Comerio, M., Celino, I.: Enabling cross-border travel offers through National Access Point federation via metadata harmonisation. In: Proceedings of the 3rd International Workshop Semantics And The Web For Transport. CEUR Workshop Proceedings, vol. 2939. CEUR, Online, September (2021). <http://ceur-ws.org/Vol-2939/paper6.pdf>, iISSN: 1613-0073
4. Comerio, M., Fiano, A., Grassi, M., Scrocca, M.: Mobility data harmonisation: the TANGENT solution. In: Transport Transitions: Advancing Sustainable and

- Inclusive Mobility - Proceedings of the 10th TRA Conference. LNM, Springer (2024), (to appear)
5. Comerio, M., et al.: TANGENT project deliverable D2.1: data requirements and available data sources (2022). <https://tangent-h2020.eu/deliverables/>
 6. Comerio, M., et al.: TANGENT project deliverable D2.2: data sharing governance model (2022). <https://tangent-h2020.eu/deliverables/>
 7. European Commission: Delegated Regulation No 886/2013 supplementing directive 2010/40/EU of the european parliament and of the council with regard to data and procedures for the provision, where possible, of road safety-related minimum universal traffic information free of charge to users (2013). <https://eur-lex.europa.eu/legal-content/EN/TXT/PDF/?uri=CELEX:32013R0886&from=EN>
 8. European Commission: delegated regulation No 2015/962 supplementing directive 2010/40/EU of the european parliament and of the council with regard to the provision of EU-wide real-time traffic information services (2015). <https://eur-lex.europa.eu/legal-content/EN/TXT/PDF/?uri=CELEX:32015R0962&from=EN>
 9. European Commission: delegated regulation No 2017/1926 supplementing directive 2010/40/EU of the european parliament and of the council with regard to the provision of EU-wide multimodal travel information services (2017). <https://eur-lex.europa.eu/legal-content/EN/TXT/PDF/?uri=CELEX:32017R1926&rid=6>
 10. Dias, T., Vieira da Silva, A., Trigueiro Moura, L., Matos, D.: Software patterns and architectural decisions for a next-generation traffic management system in the TANGENT project. In: Transport Transitions: Advancing Sustainable and Inclusive Mobility - Proceedings of the 10th TRA Conference. LNM, Springer (2024), (to appear)
 11. Garijo, D., Poveda-Villalón, M.: Best practices for implementing FAIR vocabularies and ontologies on the web. In: Cota, G., Daquino, M., Pozzato, G.L. (eds.) Applications and Practices in Ontology Design, Extraction, and Reasoning, Studies on the Semantic Web, vol. 49, pp. 39–54. IOS Press (2020). <https://doi.org/10.3233/SSW200034>, <https://doi.org/10.3233/SSW200034>
 12. Gendre, P., et al.: CHOUETTE: an open source software for PT reference data exchange. In: ITS Europe (2011)
 13. Grassi, M., Scrocca, M., Carenini, A., Comerio, M., Celino, I.: Composable semantic data transformation pipelines with Chimera. In: Proceedings of the 4th International Workshop on Knowledge Graph Construction co-located with 20th Extended Semantic Web Conference. CEUR Workshop Proceedings, vol. 3471. CEUR, Heronissos, Greece (2023). <https://ceur-ws.org/Vol-3471/paper9.pdf>, iSSN: 1613-0073
 14. Guerreiro, G., Figueiras, P., Silva, R., Costa, R., Jardim-Goncalves, R.: An architecture for big data processing on intelligent transportation systems. an application scenario on highway traffic flows. In: 2016 IEEE 8th International Conference on Intelligent Systems (IS), pp. 65–72 (2016). <https://doi.org/10.1109/IS.2016.7737393>
 15. Gutierrez, J., Samper, J.J., Delgado, A., Rocha, J., Petr, B., Zuzana, P.: LOD-RoadTran18: supporting the cross-border use of road traffic data with linked open data based on DATEX II. In: Proceedings of the 11th Euro American Conference on Telematics and Information System, pp. 1–8. EATIS 2022, Association for Computing Machinery, New York, NY, USA (2022). <https://doi.org/10.1145/3544538.3544651>
 16. Hohpe, G., Woolf, B.: Enterprise Integration Patterns: Designing, Building, and Deploying Messaging Solutions. Addison-Wesley Professional (2004)

17. Jäderberg, J.: DATEX II and Linked Open Data - Session 4 (2021). <https://www.youtube.com/watch?v=5XkMefwptyE>. Accessed 16 July 2024
18. Katsumi, M., Fox, M.: Ontologies for transportation research: a survey. *Transp. Res. Part C: Emerging Technol.* **89**, 53–82 (2018). <https://doi.org/10.1016/j.trc.2018.01.023>
19. Landaluce, H., et al.: The TANGENT project architecture: towards new traffic management approaches. In: *Transport Transitions: Advancing Sustainable and Inclusive Mobility - Proceedings of the 10th TRA Conference*. LNM, Springer (2024), (to appear)
20. Longley, D., Kellogg, G., Champin, P.A.: JSON-LD 1.1 Framing. <https://www.w3.org/TR/json-ld11-framing/>
21. Parliament, E.: Directive 2010/40/eu on the framework for the deployment of intelligent transport systems in the field of road transport and for interfaces with other modes of transport (2010). <https://eur-lex.europa.eu/legal-content/EN/TXT/PDF/?uri=CELEX:32010L0040&from=EN>
22. Richardson, C.: API Gateway pattern. <http://microservices.io/patterns/apigateway.html>. Accessed 16 July 2024
23. Ruckhaus, E., Anton-Bravo, A., Scrocca, M., Corcho, O.: Applying the LOT methodology to a public bus transport ontology aligned with Transmodel: challenges and results. *Semant. Web* **14**(4), 639–657 (2023). <https://doi.org/10.3233/SW-210451>, publisher: IOS Press
24. Sadeghi, M., Carenini, A., Corcho, O., Rossi, M., Santoro, R., Vogelsang, A.: Interoperability of heterogeneous systems of systems: from requirements to a reference architecture. *J. Supercomput.* **80**(7), 8954–8987 (2024). <https://doi.org/10.1007/s11227-023-05774-3>
25. Scrocca, M., Azzini, A., Bureš, P., Comerio, M., Lubrich, P.: Towards napDCAT-AP: roadmap and requirements for a transportation metadata specification. In: Şimşek, U., Chaves-Fraga, D., Pellegrini, T., Vahdat, S. (eds.) *Proceedings of Poster and Demo Track and Workshop Track of the 18th International Conference on Semantic Systems*. CEUR Workshop Proceedings, vol. 3235. CEUR, Vienna, Austria (2022). <https://ceur-ws.org/Vol-3235/paper22.pdf>, iISSN: 1613-0073
26. Scrocca, M., Carenini, A., Grassi, M., Comerio, M., Celino, I.: Not everybody speaks RDF: knowledge conversion between different data representations. In: *Proceedings of the 5th International Workshop on Knowledge Graph Construction*. CEUR Workshop Proceedings, vol. 3718. CEUR, Hersonissos, Greece (2024). <https://ceur-ws.org/Vol-3718/paper3.pdf>, iISSN: 1613-0073
27. Scrocca, M., Comerio, M., Carenini, A., Celino, I.: Turning transport data to comply with EU standards while enabling a multimodal transport knowledge graph. In: *Proceedings of the 19th International Semantic Web Conference*, vol. 12507, pp. 411–429. Springer, Cham (2020). https://doi.org/10.1007/978-3-030-62466-8_26
28. Van Assche, D., Delva, T., Haesendonck, G., Heyvaert, P., De Meester, B., Dimou, A.: Declarative RDF graph generation from heterogeneous (semi-)structured data: a systematic literature review. *J. Web Semant.* 100753 (2022). publisher: Elsevier
29. Van Nuffelen, B.: DCAT Application profile for data portals in Europe (DCAT-AP). <https://joinup.ec.europa.eu/collection/semantic-interoperability-community-semantic-solution/dcat-application-profile-data-portals-europe/release/211>. Accessed 16 July 2024
30. Vetere, G., Lenzerini, M.: Models for semantic interoperability in service-oriented architectures. *IBM Syst. J.* **44**(4), 887–903 (2005). <https://doi.org/10.1147/sj.444.0887>