EUROPEAN
UNION
AGENCY
FOR RAILWAYS

Making the railway system
work better for society.

# Roadmap
## *Data and digitalisation Phase 2 - Linked data mainstreaming*

### 1.    Context

**The Agency Data and Digitalization roadmap Phase 1 (2019-2020)** was approved in November 2019. It identified the objectives related to the problems to be addressed, as well as the key principles to be followed.

**Problem statement:**
In common with other organisations the Agency has traditionally taken an application-centric approach (relational databases coupled to applications) in its handling of data, creating isolated digital environments and consequently adding barriers in digital interoperability.
This approach does not support the "once only" principle nor the FAIR principle (findable, accessible, interoperable and reusable) and has a number of negative consequences, in particular:
   › *data models hidden in the application code are likely to trigger an inflation of similar, yet different data models created by IT providers. This* ***adds cost and slows innovation****;*
   › *double reporting causes* ***data inconsistencies*** *and represents a* ***double effort*** *to our stakeholders;*
   › *application code modifications are needed each time there is a new query, resulting in* ***unnecessary IT development costs*** *and vendor lock in. Several one-to-one interfaces between the different information systems in place have been gradually built to meet new information requests or to automate some tasks, leading to a drop in flexibility and increased complexity and costs.*

**The Management Board have given the Agency the task to investigate, through pilot cases, whether linked data can be used to solve the identified problems and whether this approach could be mainstreamed in a second phase.**
The Agency has done so and, based on the successful proof of concept during the pilot phase (see Annex 1) is asking the Management Board to approve the way forward as being linked data mainstreaming.

### 2.    Decision request

**The Agency asks the Management Board to approve the decision that linked data becomes the default setting for any future development of the databases, registers and specifications for data exchange mandated by the EU legal framework, under its remit. Broader application of this concept beyond the realm of ERA is certainly to be expected.**
**Any new initiatives based on the relational database approach could be accepted only as exceptions, subject to a very clear cost/benefit analysis and/or feasibility study. This would allow us to move forward with urgent development needs and would be embedded in a future ICT strategy (2021).**

Based on this MB decision, the Agency will prepare a linked data multi-annual project plan, which will include a sequence of putting in production the already implemented pilots, together with prototyping other

120 Rue Marc Lefrancq | BP 20392 | FR-59307 Valenciennes Cedex
Tel. +33 (0)327 09 65 00 | era.europa.eu
*Any printed copy is uncontrolled. The version in force is available on Agency's intranet.*

1 / 10

potential use cases, with the clear purpose to reach an Agency wide knowledge graph covering all the relevant data concepts, under clear data governance arrangements.

At the level of the Agency, the data governance will be ensured by the recently set up internal Data Governance Steering Committee, in line with the requirements of the EC Internal Audit Service.

The Data Governance Steering Committee can provide regular updates on the progress of the linked data project to the Management Board, as well as to the NRB meetings. Moreover, for the specific registers/applications, the work will involve the dedicated working groups.

## 3. Rationale supporting the decision on Phase 2

All the arguments listed below are based on a real-life prototyping that took place in the course of 2020. This was based on internal Agency work, together with the external expertise provided by DG Digit. The successful implementation of the pilot phase relied also on the valuable input provided by a number of railway actors and their sector organisation, as well as by European public bodies. The Agency would in particular like to express its gratitude to EIM, ProRAIL, SNCF, CEDEX, UIP, JRC and Eurostat for their contributions.

The pilot phase has clearly confirmed the expected linked data advantages, together with the feasibility of this approach and has helped us to dismantle a number of myths associated with linked data.

### 2.1. Linked data virtues confirmed by the pilot phase and benefits to be expected

**(I)  Interoperability and semantics**
› The linked data approach brings data together naturally, through a knowledge graph, which includes the data concepts and the data values, with no need for creating links between databases.
› An important element of the linked data approach is the clarification of the vocabulary by using universal resource identifiers (URIs) in order to univocally identify the data concepts. These are global references, which are not bound to single databases. This means that the domain concepts will have unique definitions despite being multiple labels and formats in the various databases.
› The format in which data exists does not have a relevance/importance either. The purpose of the linked data approach is **not** to harmonise the existing data formats. On the contrary, the approach respects the diversity of existing data formats and does not require that they be harmonised.

**(II)  Agility and reusability**
› It is easy *to add new data sources to an existing knowledge graph.*
› It is easy to add new concepts (meanings), while always building on what has already been developed in terms of existing vocabulary (URIs). For example, the current ontology/vocabulary covers the current RINF and ERATV data concepts related to the route compatibility check use case and can be extended in the future with new data concepts (e.g. related to ATO, route book etc.)
› Due to the less rigid architecture of a knowledge graph, changes in the data model (i.e. single value to multi-value modification) are more easily, quickly and seamlessly propagated to the application layer.

**(III)  Efficiency**
› The linked data approach does not require any change in the existing infrastructure
› A one-off effort to map concepts and data in a knowledge graph can serve several new queries that a single database would not allow. While the pilot was focused on the route compatibility check use case, by having mapped the RINF and ERATV data in the knowledge graph, other possible new use cases/queries can be addressed without having to re-run the mapping work. If new queries require certain additional data concepts which are not yet in the knowledge graph, it is enough to integrate those additional data concepts in the knowledge graph.

120 Rue Marc Lefrancq | BP 20392 | FR-59307 Valenciennes Cedex
Tel. +33 (0)327 09 65 00 | era.europa.eu

2 / 10

*Any printed copy is uncontrolled. The version in force is available on Agency's intranet.*

› Exposing and linking the data values from different data sources through the knowledge graph is an opportunity to enhance our data quality, without the need to run standalone expensive projects dedicated to data quality improvement. The pilot phase proved a good level of maturity of linked data technologies to support datasets and notably railway topology data and vehicle type data. The concepts and the relationships between them can be further developed to reach a complete and robust ontology in the railway domain.

## 2.2. Linked data myths which we are ready to challenge

### Myth no. 1 Linked data does not respect historic IT investments
› Linked data does in fact respect historic IT investments.
› Creating a link data layer does not entail a change in the application code. In fact, when running the pilot for the route compatibility check, the team did not have access to the codes of either RINF or ERATV applications
› The linked data approach can use data in whatever format it currently exists.
› It just prevents additional inefficient IT investments for querying separate datasets, changing data models or designing new ones, integrating separate applications etc.

### Myth no. 2 Linked data is costly
› The linked data approach is based on open standards and there are plenty of open source tools.
› When mapping the concepts used at enterprise level, the effort to be made for clarifying the concepts and the links between them needs to be done anyhow (also for the relational database approach).
› It generates savings on IT CAPEX and OPEX and mitigates the vendor lock-in.
› It worked for a rather complex pilot with a rather limited volume of resources.

### Myth no. 3 There is limited in-house expertise
› Linked data is more about vocabularies and concepts than IT. It's in fact good domain experts that it requires.
› It was relatively easy for our team to learn how to query the data from the knowledge graph.

### Myth no. 4 Other actors might be unable to use linked data
› Other actors such as the Infrastructure Managers can develop this expertise, but this is not compulsory – linked data can also work with the data that they have in the existing formats.
› Taking the example of RINF, which is currently fed by IMs uploading the required data as .xml files, the switch to the linked data approach would mean that eventually the knowledge graph can be populated in different ways than uploading a file, for example by crawling data published on their own website.
› The knowledge graph prepared by the Agency proved to be compatible with the linked data approaches taken in other domains, with no additional effort for ensuring their integration.

120 Rue Marc Lefrancq | BP 20392 | FR-59307 Valenciennes Cedex
Tel. +33 (0)327 09 65 00 | era.europa.eu
*Any printed copy is uncontrolled. The version in force is available on Agency's intranet.*

3 / 10

**Annex 1: Outcomes of the pilot phase**

The current legal framework specifies that after vehicle authorization, the Railway Undertaking is responsible for the route compatibility check between the vehicle and the intended routes of operation within the area of use. Once this check is done, the RU do not need to repeat it unless there is a change in the infrastructure or a change in the vehicle.

We have not followed an application-centric approach but a data centric one. No new information system collecting information from all EU vehicles types and collecting the complete operational characteristics of all desirable EU train routes has been implemented.

In accordance to the data centric, reusable and interoperable per design approach, we have:

- reused the information from current registers at Agency premises, RINF and ERATV
- used semantic technologies to achieve data interoperability between those information systems which use different data models.
- the existing legacy databases (ERATV and RINF) that have been mapped to the Data Centric model (ERA Knowledge Graph using open standard W3C protocols.
- extracted new knowledge from the combination of both information datasets
- other datasets were loaded directly into de Data Centric model (ERA knowledge graph)
- identified and easily spotted data quality issues

We have not used any application code from ERATV or RINF as input to the exercise but just the data hosted in their respective databases. In this specific first Pilot 1 first stage, we have annotated/mapped each concept involved in the route compatibility check process as per TSI-OPE annex D.

This mapping process is done using a W3C standard protocol and consists of matching each concept to an existing concept in the ERA vocabulary. This vocabulary is semantically expressed, shareable and hosted in the EU Vocench accessible platform. In a nutshell, the load to a data centric architecture has been fast. It has entitled us to obtain a first flavour in relation to what our EU legal databases may contribute to the EU railway operation scenario.

By linking the two data sets and exposing the data, we have spotted easily large amount of missing information (see *no data* label in Fig 1). Spotting the misalignment is the first necessary step to take to achieve consistent data. This exercise has entitled us to provide the first step towards the automation of the information exchange needed to perform the Route Compatibility Check

Actors Involved: ERA, DG DIGIT, ProRail, EIM, SNCF, UIP

Current status: The ERA knowledge graph 1.0.0 is the outcome of the mapping – linked to the route compatibility use case – of the RINF and ERATV datasets. Using this RDF data, the route compatibility check is performed. The demo can be found in the deliverables below.

Deliverables (hyperlinks included):

1. *Route Compatibility Check demo, Source code*
2. *ERA knowledge graph GraphDB, EC Data Platform*
3. *ERA Vocabulary , Ontology*
4. *ERA data mappings*

More details can be found @ERA youtube channel: https://youtu.be/KnntVq3Zxzk

120 Rue Marc Lefrancq | BP 20392 | FR-59307 Valenciennes Cedex
Tel. +33 (0)327 09 65 00 | era.europa.eu
*Any printed copy is uncontrolled. The version in force is available on Agency's intranet.*

4 / 10

1. *Route Compatibility Check demo*, *Source code*

Detailed description of the demo (User Guide Help Page)

The prototype of the Route Compatibility Check application and service works on the Knowledge Graph of the EU Agency for Railways (ERA). The objective is to check if a selected vehicle type can travel the route from operational point A to operational point B. Each route is composed of sections of lines (tracks) with different technical parameters. Each track is between two operations points. The small icons that appear on the map show the location of the operational points and details about them when you hover.



Fig 1. Print Screen from Agency´s demo on route compatibility [Route Compatibility Check demo]

The railway topology data is georeferenced and is published in *tiles* through a semantically annotated API (i.e. an API whose responses include metadata that semantically describe how clients can interact with it). This API follows the same principles of the Slippy Map approach, used by map applications (e.g. Google Maps and OSM) to serve only the data related to a specific region that an application is trying to visualize at a certain moment. This means that when you pan or zoom on the map, the application will request to the API the data for the region seen in your screen, according to the Slippy Map specification. The API will receive this request and will proceed in turn to perform a SPARQL query to the KG data store for these specific data.

Many of the items in the comparison table of each track on the selected route are clickable. These are not just links to pages but links to the edges and nodes of the knowledge graph. Edges are in the column Properties. The properties link in a meaningful way a node, such as a track, to another node.

120 Rue Marc Lefrancq | BP 20392 | FR-59307 Valenciennes Cedex
Tel. +33 (0)327 09 65 00 | era.europa.eu
*Any printed copy is uncontrolled. The version in force is available on Agency's intranet.*

5 / 10

This other node can be a value, for example T_81527_358_81534 - maximum temperature -> 40, or another node identified with URI, for example T_81527_358_81534 - train detection system -> - track circuit - , wheel detector. Each property has meaning defined by the ERA vocabulary.

With them -both humans and machines- can "understand" that a node is of specific type, for example track, operational point, or vehicle type, or the meaning of a statement such as the examples above. If you click on a property you can see its description as a set of triples which is the pertinent subset of the triples from the vocabulary (this is the knowledge model, or ontology, which is also a knowledge graph).

If you click on a Track, or Vehicle, you'll see that part of the ERA knowledge graph, which describes them again in the same way, as triples of subject-predicate-object (called alternatively resource-property-value).

If you click on a clickable parameter value in the columns Track or Vehicle, you'll go to a description of a reference data node. For example, clicking on track circuit, you'll be able to navigate the graph, just by clicking, to all the tracks and vehicle types which train detection system is track circuit. From any node, you can explore the knowledge graph node by node just by clicking on the neighbouring nodes.

*2. ERA knowledge graph GraphDB, EC Data Platform*

The ERA knowledge graph is hosted in two different platforms (GraphDB and Virtuoso from EC Data Platform). This duplicity validated also the portability and consequently the vendor independence of the technique used. Both platforms offers a SPARQL endpoint. Different queries are already pre-recorded (Fig.2) and others can be created on demand by the end user (Fig.3). Particularly, many queries to spot data quality issues have been constructed. In Fig. 4, we can observe Operational Points detached from the section of line. Many data quality issues have appeared particularly when adding and operational use to the data set. For example, in Fig. 1, empty boxes labelled as "no data" appeared.



120 Rue Marc Lefrancq  |  BP 20392  |  FR-59307 Valenciennes Cedex
Tel. +33 (0)327 09 65 00  |  era.europa.eu
*Any printed copy is uncontrolled. The version in force is available on Agency's intranet.*

6 / 10

Fig. 2 ERA knowledge graph expressed in GraphDB hosting 6 Million triple statements.

120 Rue Marc Lefrancq | BP 20392 | FR-59307 Valenciennes Cedex
Tel. +33 (0)327 09 65 00 | era.europa.eu
*Any printed copy is uncontrolled. The version in force is available on Agency's intranet.*

7 / 10

Fig.3  Details of the ERA types of vehicle which are compatible with Austrian tracks (gauging profile)



Fig.4  Details of the Spanish Operational Points unreachable from the railway line (error in the data provision to RINF )

120 Rue Marc Lefrancq  |  BP 20392  |  FR-59307 Valenciennes Cedex
Tel. +33 (0)327 09 65 00  |  era.europa.eu
*Any printed copy is uncontrolled. The version in force is available on Agency's intranet.*

8 / 10

## 3. ERA Vocabulary , Ontology



Fig.5 Extract of visualised Ontology



Fig.6 ERA Vocabulary overview

*Any printed copy is uncontrolled. The version in force is available on Agency's intranet.*
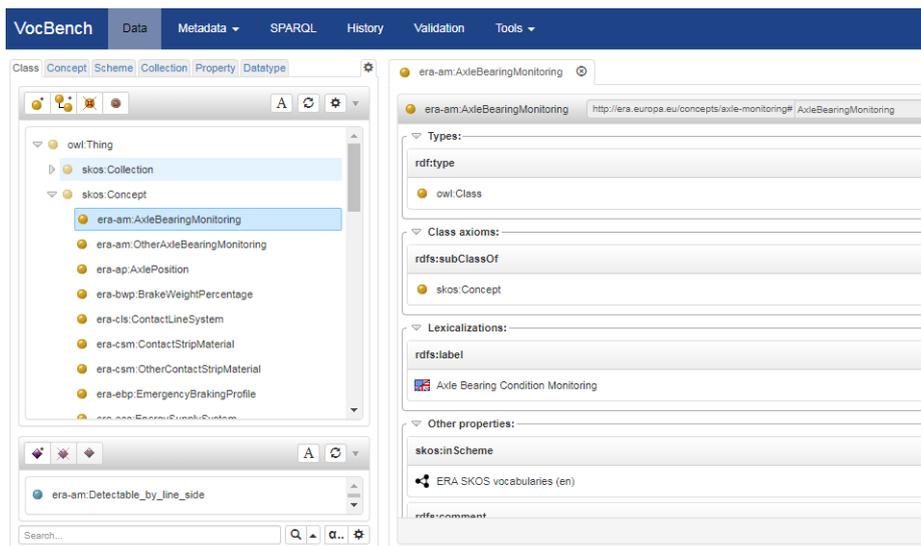
Fig.7 Version of ERA Reference Data imported in VocBench VocBench

The **sub-pilot 1.2 (Expose information on usage of NOI TSI compliant and non-compliant wagons on quieter routes)** consists of adding two external data sources to the ERA Knowledge graph to expose the data related to noise related restrictions in the ECVVR register plus the data related to noise and breaking characteristics from UIP RSRD database. These new pieces of information are neither public as the previous private content of the ERATV and RINF datasets, nor fully internal to the Agency.

Data source 1 provenance is UIP RSRD private database. Data catalogue and sample of data consisting of vehicle number, keeper data, Brake Special Characteristics, CompositeBrakeBlockFitted, Quieter Routes exemption country.

Data Source 2 provenance is ERA register on vehicle information.  ECVVR. This register includes the vehicle number, optionally the Type of Vehicle as per ERATV, noise restrictions as per the list of harmonised national restriction codes (2.7.1 Can be used in all quieter routes-TSI Noise compliant –Silent-retrofitted without testing, 2.7.2 Can be used in all quieter routes-TSI Noise compliant –Silent (tested against a TSI NOI), 2.7.3 Can be used in all quieter routes-TSI Noise compliant –Very quiet (tested against a TSI NOI), 2.7.3 Can be used in all quieter routes-TSI Noise compliant –exempted in accordance with TSI NOI). 2.7.5 Can be used in quieter routes only in this MS-Covered by specific case. 2.7.6 Can be used in quieter routes in this MS- Fitted with historic CBBs. 2.7.7 Can´t be used in quieter routes.

Data Source 3 is already in the knowledge graph. It is the property Quieter route from a Section of Line from the RINF register.

The user in this Sub Pilot is able to type a vehicle number and the route from Operational Point A to Operational Point B. If the vehicle number has an ERA vehicle type the complete route compatibility check is carried out in accordance to its type and adding the information on noise compliance from both source of information. If the vehicle number does not have an ERA vehicle type the compatibility is only carried out in relation to the previously mentioned noise parameters.

Under **Pilot 2 (Harmonized way of sharing the rail topology with enough detail and accuracy to support signalling activities in a semantic machine readable manner as to reduce the data integration current costly effort between IMs and ERTMS suppliers)**, which is due February 2021, the Railway topology description from RINF data model -currently expressed in ERA Knowledge graph- evolves to a higher level of detail- by integrating it to other data sources with railway topology information. These data sources include: IM Spoor ProRail – nanoscopic model, ADIF data provided by CEDEX, routable GeoData from ESTAT. In parallel, the Agency´s request for research on this topic is currently addressed by S2R Lynx4rail project on conceptual data model and will cover the ATO over ETCS rail topology data exchange between IMs and ERTMS suppliers.

120 Rue Marc Lefrancq  |  BP 20392  |  FR-59307 Valenciennes Cedex
Tel. +33 (0)327 09 65 00  |  era.europa.eu
*Any printed copy is uncontrolled. The version in force is available on Agency's intranet.*

10 / 10