



# D3.15 – i4Q GUIDELINES FOR BUILDING DATA REPOSITORIES FOR INDUSTRY 4.0 v2

WP3 – BUILD: Manufacturing  
Data Quality



## Document Information

GRANT AGREEMENT NUMBER	958205	ACRONYM	i4Q
FULL TITLE	Industrial Data Services for Quality Control in Smart Manufacturing		
START DATE	01-01-2021	DURATION	36 months
PROJECT URL	<a href="https://www.i4q-project.eu/">https://www.i4q-project.eu/</a>		
DELIVERABLE	D3.15 – i4Q Guidelines for Building Data Repositories for Industry 4.0 v2		
WORK PACKAGE	WP3 – BUILD: Manufacturing Data Quality		
DATE OF DELIVERY	CONTRACTUAL	31-Dec-2022	ACTUAL 30-Dec-2022
NATURE	Report	DISSEMINATION LEVEL	Public
LEAD BENEFICIARY	ITI		
RESPONSIBLE AUTHOR	ITI		
CONTRIBUTIONS FROM	1-CERTH, 2-ENG, 5-KBZ, 11-UNI		
TARGET AUDIENCE	1) i4Q Project partners; 2) industrial community; 3) other H2020 funded projects; 4) scientific community		
DELIVERABLE CONTEXT/DEPENDENCIES	<p>This document is a public report developed within the context of “Task 3.5 - Manufacturing Data Storage and Use”, that presents a guide for building Data Repositories for Industry 4.0, which corresponds to the i4Q<sup>DRG</sup> solution.</p> <p>This document is the second version of D3.7 – <i>Guidelines for building Data Repositories for Industry 4.0</i>. It also provides feedback to deliverable D3.16 - i4Q Data Repository v2.</p>		
EXTERNAL ANNEXES/SUPPORTING DOCUMENTS	<p>Appendix I of this deliverable provides further details on a questionnaire distributed to pilots and solutions providers to gather specific and technical requirements.</p> <p>Appendix II of this document contains the i4Q<sup>DRG</sup> web documentation.</p>		
READING NOTES	None		

## ABSTRACT

This deliverable presents an overview of the role and importance of data repositories in Industry 4.0 contexts, such as this project. Furthermore, this document explains the challenges and requirements arising when developing data repositories and provides some recommendations on how to address them, using the i4Q<sup>DR</sup> as an illustrative example.

## Document History

VERSION	ISSUE DATE	STAGE	DESCRIPTION	CONTRIBUTOR
0.1	7-Nov-2022	Draft	First Version of Table of Contents, based on final version of D3.7	ITI
0.2	25-Nov-2022	Draft	First draft version: added sections 4 and 5 and updated the rest of the document with the progress made from M18 to M24.	CERTH, ITI
0.3	7-Dec-2022	Internal Review	Internal review	FACTOR, KBZ
0.4	12-Dec-2022	Draft	Second draft version addressing reviewers' comments.	ITI
1.0	30-Dec-2022	Final Doc	Final quality check and issue of final document	CERTH

## Disclaimer

Any dissemination of results reflects only the author's view and the European Commission is not responsible for any use that may be made of the information it contains.

## Copyright message

### © i4Q Consortium, 2022

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 authorised provided the source is acknowledged.

## TABLE OF CONTENTS

Executive summary .....	6
Document structure .....	7
1. Introduction .....	8
2. Guideline on how to tackle data repositories for Industry 4.0 .....	11
3. Design specifications of the i4Q Data Repository .....	13
4. Communication with Data Repositories .....	18
5. Implementation of i4Q <sup>DRG</sup> as web documentation .....	21
6. Conclusions .....	23
7. References .....	24
Appendix I: questionnaire on data storage needs .....	25
Appendix II: web documentation .....	29

## LIST OF FIGURES

<b>Figure 1.</b> Connections and interactions of i4Q <sup>DR</sup> .....	8
<b>Figure 2.</b> Pipeline of Pilot 1 (FIDIA) .....	9
<b>Figure 3.</b> Pipeline of Pilot 3 (Whirlpool) .....	9
<b>Figure 4.</b> Overview of questionnaire to gather needs to be covered by the i4Q <sup>DR</sup> .....	14
<b>Figure 5.</b> Architectural overview of i4Q <sup>DR</sup> .....	17
<b>Figure 6.</b> MinIO Console deployed by the i4Q <sup>DR</sup> for any MinIO scenario .....	19
<b>Figure 7.</b> Mongo Express instance deployed by the i4Q <sup>DR</sup> for any MongoDB scenario .....	20
<b>Figure 8.</b> Snapshot of web application showing information regarding the i4Q <sup>DRG</sup> solution .....	21
<b>Figure 9.</b> Example of highlighting of searched word .....	22
<b>Figure 10.</b> Example of search results .....	22
<b>Figure 11.</b> Questions on necessary type of storage .....	26
<b>Figure 12.</b> Questions on expected data volumes .....	26
<b>Figure 13.</b> Questions on type of interaction and communication interfaces .....	27
<b>Figure 14.</b> Questions on other communication issues and data replication needs .....	27
<b>Figure 15.</b> Questions about authorization, deployment, and performance .....	28
<b>Figure 16.</b> Questions about other specific needs and possible constraints and preferences .....	28

## ABBREVIATIONS/ACRONYMS

<b>ANSI</b>	American National Standards Institute
<b>API</b>	Application Programming Interface
<b>BDA</b>	Big Data Analytics
<b>CNC</b>	Computer Numerical Control
<b>CRUD</b>	Create, Read, Update, and Delete
<b>DR</b>	Data Repository
<b>DRG</b>	Data Repository Guidelines
<b>HA</b>	High Availability
<b>HA+Sec</b>	High Availability with Security
<b>HTTP</b>	Hypertext Transfer Protocol
<b>HTTPS</b>	Hypertext Transfer Protocol (Secure)
<b>IM</b>	Infrastructure Monitoring
<b>IT/OT</b>	Information Technology/Operational Technology
<b>JSON</b>	JavaScript Object Notation
<b>OS</b>	Operating System
<b>PA</b>	Perspective Analysis
<b>PDF</b>	Portable Document Format
<b>PQ</b>	Process Qualification
<b>REST</b>	Representational State Transfer
<b>REST API</b>	RESTful Application Programming Interface
<b>SQL</b>	Structured Query Language
<b>SS</b>	Single Server
<b>SS+Sec</b>	Single Server with Security
<b>TCP</b>	Transfer Control Protocol
<b>TLS</b>	Transport Layer Security

## Executive summary

---

**D3.15** is the second and final version of a guide devoted to the building of Data Repositories for Industry 4.0, which also tackles the estimation of storage capabilities and the building of technological systems that provide easy ways to access manufacturing data and to communicate with industrial platform components and microservice applications. These guidelines correspond to the **i4Q Data Repository Guideline (i4Q<sup>DRG</sup>)** solution.

In this document, firstly the role of data repositories in Industry 4.0 contexts is explained and is presented the importance of the requirements they fulfil. Using illustrative examples, this role explanation is matched with the **i4Q** project and its solutions, especially with the **i4Q Data Repository (i4Q<sup>DR</sup>)** solution.

Furthermore, this deliverable explains the challenges arising when building a data repository for Industry 4.0 and provides an overview on how to approach its development by gathering a set of recommendations addressing these challenges, focusing especially on the analysis and design phases, but also covering some aspects related to the implementation phase.

D3.15 **extends** the information gathered in “*D3.7 – i4Q Guidelines for Building Data Repositories for Industry 4.0*” (submitted in M18) by providing an overview on the possible approaches to achieve the communication of data repositories with other software solutions. Again, the **i4Q<sup>DR</sup>** solution is used as an illustrative example. Furthermore, D3.15 presents the web application that has been implemented in the M18-M24 period to show the most relevant information gathered in this deliverable in the style of standard web documentation.

These guidelines have driven the development of the **i4Q Data Repository (i4Q<sup>DR</sup>)** Solution, which is presented in deliverables “*D3.16 - i4Q Data Repository v2*”.

This document, “*D3.15 – i4Q Guidelines for Building Data Repositories for Industry 4.0 v2*”, is an update of v1 (D3.7). For this reason, it contains information of the 1st version together with the updates developed in this 2nd version.

## Document structure

---

**Section 1:** provides an overview on the features provided by data repositories in Industry 4.0 contexts, explaining the key role they play. This explanation is illustrated using the *i4Q* project and its solutions (especially the *i4Q<sup>DR</sup>*) as running examples.

**Section 2:** gathers some of the most important challenges and requirements arising when developing data repositories for Industry 4.0. Moreover, it proposes some recommendations and procedures addressing them.

**Section 3:** explains in detail how the recommendations and guidelines provided in Section 2 have been applied to define the design and the development *strategy* of the *i4Q<sup>DR</sup>*. Note, however, that this section does not provide specific details on the implementation of the *i4Q<sup>DR</sup>*. This issue is out of the scope of this deliverable but will be covered by “D3.16 - *i4Q Data Repository v2*” (due on M24).

**Section 4:** describes different approaches to implement the communication of data repositories with other software solutions. For illustrative purposes, the *i4Q<sup>DR</sup>* is used as an example.

**Section 5:** presents the web application implemented to display the most relevant information gathered in this document in a more user-friendly fashion.

**Section 6:** provides the conclusions of this document, summarising its main contributions and providing a brief explanation on the future version of this deliverable.

**Appendix I:** provides information on a questionnaire used to gather specific needs from pilots and other *i4Q* solutions that must be covered by the *i4Q<sup>DR</sup>*.

**Appendix II:** includes the web documentation of the *i4Q Data Repository Guidelines* (*i4Q<sup>DRG</sup>*) solution.

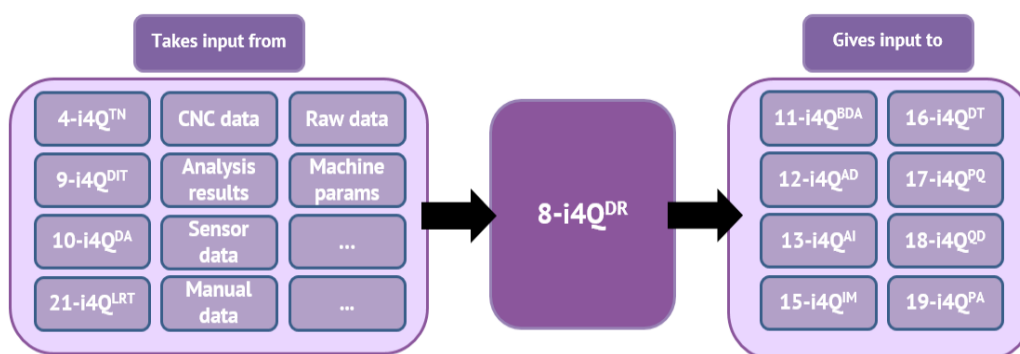


## 1. Introduction

One of the main characteristics of the so-called industry 4.0 is the application of technical solutions to retrieve data from industrial processes so that it can be analysed and processed afterwards. The goal of such exercise is to obtain relevant knowledge, that can be applied to improve these industrial processes and make them more efficient. For instance, by reducing the percentage of defects, or by detecting them in earlier stages of the production process, which typically reduces the economic costs of those defects.

In a typical Industry 4.0 scenario, most part of the data is *generated* by manufacturing devices acting as *sensors* and is dumped into a storage technology or tool. Then, there are several *data analysis and processing tools* that consume this data for different purposes. For instance, there are tools that perform data cleaning processes to remove useless cases, or to extract only the parts of data that really matters for a specific purpose. Other tools, such as dashboards, show the stored information in a more graphical and intuitive way, to facilitate some decision-making processes. Other important data consumers are big data analytics tools, that apply different algorithms to the data and allow, for example, to discover correlations and relationships among different variables of the analysed data.

In the case of the *i4Q* project, the *i4Q Data Repository (i4Q<sup>DR</sup>)* is the solution covering all the data storage and retrieval requirements of the other solutions, which is presented in more detail in deliverable D3.8. Then, there are many other *i4Q* solutions that fall into the category of *data analysis and processing tools* mentioned, such as the *i4Q Data Integration and Transformation Services (i4Q<sup>DIT</sup>)*, the *i4Q Services for Data Analytics (i4Q<sup>DA</sup>)*, and *i4Q Big Data Analytics Suite (i4Q<sup>BDA</sup>)*, which are presented in deliverables D4.1, D4.2 and D4.3, respectively. These solutions take data previously stored in the *i4Q<sup>DR</sup>* as input and produce new one that can be stored into it as well. Another source of data for the *i4Q<sup>DR</sup>* is the *i4Q Trusted Networks with Wireless and Wired Industrial Interfaces (i4Q<sup>TN</sup>)* solution (see deliverable D3.4), which oversees collecting data generated by actuators, sensors, and other IT/OT systems and give it as input for the *i4Q<sup>DR</sup>* and the *i4Q Analytics Dashboard (i4Q<sup>AD</sup>)* solution (described in deliverable D4.4). Figure 1 shows an overview of connection and interactions of the *i4Q<sup>DR</sup>*, including other *i4Q* solutions. More specifically, elements in the box labelled as “Takes input from” are the ones requiring the storage of data into the *i4Q<sup>DR</sup>*, whereas the ones located within the box labelled as “Gives input to” correspond to solutions retrieving data from the *i4Q<sup>DR</sup>*.



**Figure 1.** Connections and interactions of *i4Q<sup>DR</sup>*

Note that the  $i4Q^{DR}$  in some cases can be the final destination of the stored data, but in some other cases, the stored data will be provided as input for another solution that will further process it. This is well illustrated by the pipelines designed for the  $i4Q$  pilots. For instance, in Pilot 1 (see pipeline in Figure 2) data stored into the  $i4Q^{DR}$  will not be retrieved by any other solution, whereas in Pilot 3 (see pipeline in Figure 3) such data will be sent via the Message Broker to other solutions, and directly retrieved by the  $i4Q$  Prescriptive Analysis Tools ( $i4Q^{PA}$ )<sup>1</sup>. As discussed above, storing, and retrieving data is a basic but crucial feature in an Industry 4.0 scenario. Therefore, the development of technical solutions providing that functionality, such as the  $i4Q^{DR}$ , must be addressed very carefully so that this feature is provided in the most appropriate way for each different component that interacts with it.

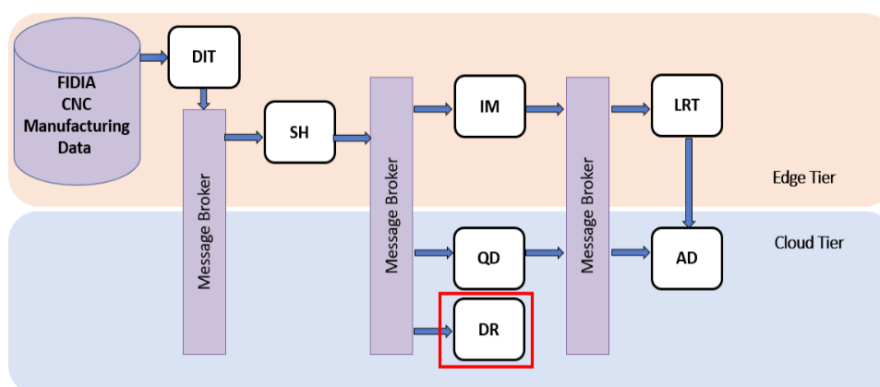


Figure 2. Pipeline of Pilot 1 (FIDIA)

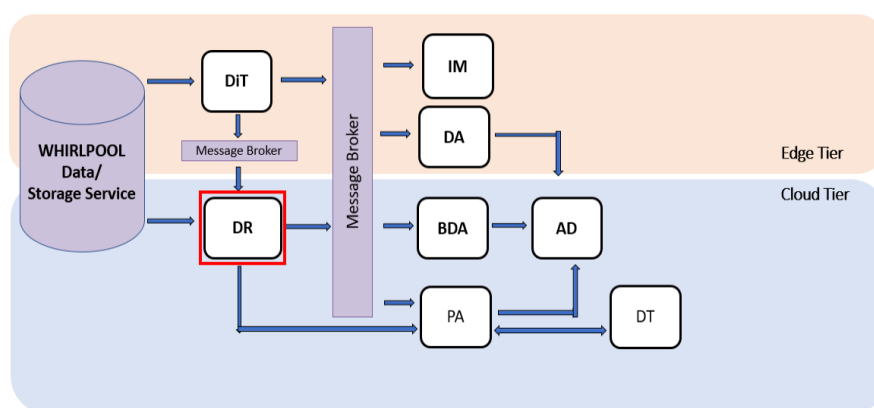


Figure 3. Pipeline of Pilot 3 (Whirlpool)

Indeed, the interaction with other solutions is the main challenge for the development of the  $i4Q^{DR}$ . In this sense, some complexity arises just due to the number of different solutions interacting with the  $i4Q^{DR}$ . However, the highest degree of complexity is due to the fact that solutions retrieving or storing data from or into the  $i4Q^{DR}$  may need different functionalities from it, and, specially, have different technical requirements. For instance, some solutions may only

<sup>1</sup> In both, Σφάλμα! Το αρχείο προέλευσης της αναφοράς δεν βρέθηκε. and Σφάλμα! Το αρχείο προέλευσης της αναφοράς δεν βρέθηκε., solutions are denoted by the acronyms of the solution's code without the word " $i4Q$ ". For example, the  $i4Q$  Data Repository, whose code is  $i4Q^{DR}$ , is denoted by "DR".



need to store data, whereas other may need importing/exporting data from/to other data sources. Moreover, some solutions may require an SQL-style relational database, while others may require storing files.

This document provides the **final** version of the **i4Q** Data Repository Guidelines solution (**i4Q<sup>DRG</sup>**), a guide devoted to the building of Data Repositories for Industry 4.0, which also tackles the estimation of storage capabilities and the building of technological systems that provide easy ways to access industry data and to communicate with industrial platform components and microservice applications. More specifically, this document focuses on the challenges that may arise during the development of a data repository, and some recommendations to address them.

The recommendations gathered in this document have guided the development of the **i4Q<sup>DR</sup>** so far, up to **M24**. Therefore, this solution will be used as an example to illustrate some of the concepts and recommendations explained in this document.

## 2. Guideline on how to tackle data repositories for Industry 4.0

---

Data repositories may be thought of as low-complex technical solutions, since they provide a basic functionality, namely the storage and retrieval of data. However, its development is quite challenging since, as explained in Section 1, they usually interact with many other components with different purposes. Consequently, in this context, the requirements of a data repository, especially the technical, are quite *heterogenous*, since the ones for a given component can be quite different from the ones for another component.

In this sense, a good example is the type of storage technology, which will vary depending on the type of data to be stored. For instance, SQL-based technologies are the suitable ones for structured and relational data, whereas NoSQL storage technologies are more appropriate for data whose structure may change over the time. However, it could also be the case that two components have contradictory requirements. For instance, some components running on premises with constrained resources will require light mechanisms to store just a few bytes of data, whereas other components may require more complex and powerful mechanisms to store large amounts of data.

Therefore, the development of data repositories in that context requires a balance between two opposite concepts. On the one hand, data repositories need to have a *flexible design*, so that they can fulfil heterogenous requirements coming from different dependent solutions. In the other hand, to reduce the complexity of the development, and the maintenance of the data repository once it is deployed, it is necessary to *provide* the required *functionality* in the *most centralised fashion possible*, by extracting as many common functionalities as possible and implement them in the most harmonised possible way. For instance, a data repository should try to offer only one interface to receive data from other solutions to store it, independently of the type of data. In this way, there is only one implementation for that feature, instead of one per storage technology.

The first step to address these challenges when developing a data repository is to identify which components interact with it, either to provide data that needs to be stored, or to retrieve data that has previously been stored. Then, for each solution interacting with the data repository, it is necessary to gather the specific needs regarding several aspects, such as:

- The functionality required (e.g., only CRUD, or something else).
- The volume of data.
- Possible performance constraints, such as the expected number of read or write operations per second.
- The operation mode. For instance, whether the data will be provided in an interactive mode, or in batch, and if will be done on-demand or as a scheduled task.
- The need of replicated instances of the data repository.
- The type of communications and expected interfaces and data formats.
- Who and what components are authorised to access the data repository.

This information must be carefully analysed to make the decisions that will define the main aspects of the data repository's design, namely:

- The most appropriate tools and technologies for the type of data that will be managed.
- The architecture of the data repository, which must be flexible to accommodate different

requirements but, at the same time, should try to fulfil them in a centralised fashion whenever possible, as mentioned above.

- The necessary storage configuration options considering the functional requirements and technical details such as the amount of data that will be handled, and the computing resources that will be available once the data repository is deployed.

For making these decisions it might be helpful to categorise needs and functionalities into those that are specific only to few solutions (or only one) and those that are common to several of them. Addressing the first group of needs might require using concrete tools. However, whenever possible, one should try to use tools that can work for other cases, too. Common needs and functionalities should be addressed in the most general and harmonised way. For instance, using the same technology or tool, whenever possible.

At the implementation level, the most important strategy is to simplify and reduce the complexity of the developments. This can be achieved, for instance, by providing only one implementation of a common functionality, independently of the solution requiring it. This can be the case, for instance, of access control. If the implementation differs depending on the solution the data repository is interacting with, it is recommended to analyse whether there are some sub-functionalities that can be implemented in the same way. For instance, if a solution needs to store data in a relational database, and another one needs to store a file in an object storage tool, the concrete mechanism to do so will be different in both cases. However, as mentioned above, the data repository can offer a common interface to receive the data to be stored, and then run the corresponding storage mechanism.

When different implementations for the same feature are necessary, the recommendation is to follow a similar structure and style to facilitate the identification of the functionality that is being implemented and the update of the code if changes are necessary in the future. An example of this is using similar class or method names.

Finally, another way to simplify and decrease the complexity of developments consists of reducing the number of dependencies and using tools and libraries that do not have too specific requirements. For instance, Docker is compatible with most common operating systems. Furthermore, Docker containers allow to abstract the (virtual) execution of a tool from the operating system (OS) in which Docker is running. This means that a tool running inside a Docker container executed on a Windows OS behaves *exactly the same* as when the tool is running in a Docker container executed on a Linux OS, assuming both containers have the same configuration.

### 3. Design specifications of the i4Q Data Repository

---

This section describes how the design of the i4Q<sup>DR</sup> and the definition of the *strategy* for its development have been approached so far (M18), following the recommendations and guidelines provided in Section 2 to address the challenges and requirements arising when developing data repositories for Industry 4.0. Further details on the implementation of the i4Q<sup>DR</sup> will be explained in deliverable “D3.8 - i4Q Data Repository”(due on M18).

Regarding the interactions of the i4Q<sup>DR</sup> with other components, the i4Q project includes a number of pilots and solutions that manage data one way or another and thus need mechanisms to handle such data. More specifically, the i4Q<sup>DR</sup>, takes as input data from other i4Q solutions as well as raw data from other industrial components (e.g. CNC data or data gathered from sensors), and even data manually provided by a human user. Furthermore, the i4Q<sup>DR</sup> provides data to other i4Q solutions, such as the i4Q<sup>BDA</sup>, the i4Q<sup>IM</sup>, the i4Q<sup>PO</sup>, etc. As explained in Section 2, Figure 1 gathers the interaction and connections of the i4Q<sup>DR</sup> with other components and solutions. Further information on this topic can be found in deliverable “D1.9 –Requirements Analysis and Functional Specification v2”(Section 4.8).

With respect to the features and functionalities that the i4Q<sup>DR</sup> is required to provide, the main needs common to all pilots and solutions are basic functionalities like storing data and retrieving the stored data. However, there are some additional generic requirements which are quite common in any storage technology, including:

- *Data persistence*: the data is saved to a persistent storage and remains intact until it is altered or deleted on purpose.
- *Availability*: the data is available to any of its legitimate users.
- *Privacy*: the data is protected against unauthorized access.
- *Security*: the data is protected against software and hardware failures.
- *Efficiency*: with a proper use of the available resources.

Moreover, each pilot and solution have specific needs that involve several factors. One of the most important criteria is the *type of the data* to be stored, either relational, document-based, graph-oriented, blob-based, file-based, etc. This criterion on its turn affects other criteria. For instance, the use of enormous blobs or files may make a significant impact on the performance numbers. On the other hand, the same solution may have different needs depending on which pilot it is used. Thus, the i4Q<sup>DR</sup> must be able to offer its services to its clients, under a wide range of needs, scenarios, and settings.

As a second step, needs related to the *interoperability* of the data have been identified. For instance, a pilot or solution generates data with different structures and/or syntax (e. g. relational and document-based data) or data that, for some reason, has to be handled by different storage tools. Another example may be the case of a solution that handles different sets of data and each one has different needs (formats, syntax, access patterns, performance requirements, etc.) and needs to use them in a unified way.

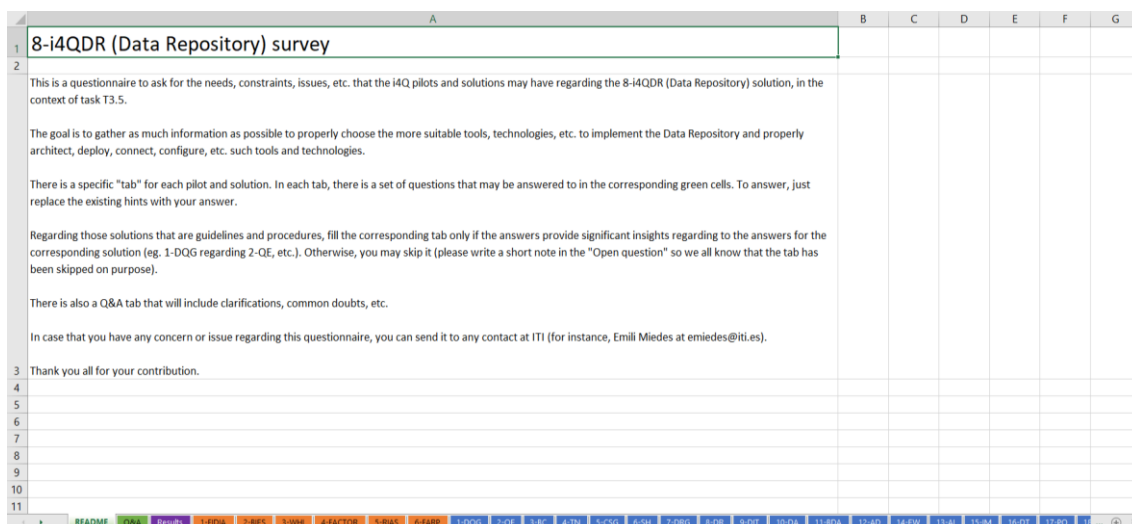
The i4Q<sup>DR</sup> includes some off-the-shelf open-source tools. The selection of the current tools has been performed according to the needs of both the pilot and solution partners of the project. The procedure followed to produce such selection includes the following steps.



1. Prepare a questionnaire to ask about the needs a pilot or solution has regarding the i4Q<sup>DR</sup>.
2. Receive the answers of the questionnaire and summarize them.
3. Identify which tools are specifically required by the partners.
4. Identify other tools that can satisfy the requirements and needs of the partners.

The questionnaire is a Microsoft Excel document available in a private repository<sup>2</sup> (see Appendix I for further details). As shown in Figure 4, which provides an overview of such document, it includes a page for each pilot and solution partner. These pages contain a set of questions to be answered by each partner. The questions include these issues:

- The types of storage mechanism needed by the pilot or solution (relational, structured, JSON-based, key-value-oriented, etc.).
- Any information available so far regarding the volume of the data to be handled (for instance, in orders of magnitude).
- The type of interaction that may be expected (interactive, batch, stream-oriented, etc.).
- The type of communication interfaces that may be used (HTTP, HTTPS, TCP, etc.).
- Any other requirement related to communication.
- Any requirement or need related to data replication.
- Any requirement regarding data privacy and security (access control, anonymization, etc.) that may apply.
- Any specific deployment needs if any (on bare machines, Docker, Kubernetes, etc. and also on-premises vs cloud).
- Any requirement or need regarding the performance.



**Figure 4.** Overview of questionnaire to gather needs to be covered by the i4Q<sup>DR</sup>

<sup>2</sup> Survey to solution providers and pilots to gather i4Q<sup>DR</sup> requirements: [https://knowledgebiz.sharepoint.com/:x/r/sites/i4Q/\\_layouts/15/Doc.aspx?sourcedoc=%7B4A34149D-B1D1-4146-B4EF-0E44056E200E%7D&file=Survey\\_8-i4QDR.xlsx&cid=79408e25-3af5-4d1e-a129-6e25abf22f3e](https://knowledgebiz.sharepoint.com/:x/r/sites/i4Q/_layouts/15/Doc.aspx?sourcedoc=%7B4A34149D-B1D1-4146-B4EF-0E44056E200E%7D&file=Survey_8-i4QDR.xlsx&cid=79408e25-3af5-4d1e-a129-6e25abf22f3e)

After analysing the survey's answers and carefully evaluating the needs of the partners, the following tools were selected:

- *Cassandra*<sup>3</sup>, a wide-column NoSQL distributed database server, appropriate for managing massive amounts of data.
- *MinIO*<sup>4</sup>, which offers a high-performance, S3 compatible object storage. It is used to store files and is compatible with any public cloud.
- *MongoDB*<sup>5</sup>, a JSON document-oriented database server.
- *MySQL*<sup>6</sup>, a SQL relational database server.
- *Neo4J*<sup>7</sup>, a graph database server that allows storing data relationships.
- *Redis*<sup>8</sup>, an in-memory data structure store, used as a distributed, in-memory key-value database, cache, and message broker, with optional durability.

In addition, the following tools have been included later, to satisfy additional needs and requirements brought out by some partners:

- *MariaDB*<sup>9</sup>, a SQL relational database server, very similar to MySQL.
- *PostgreSQL*<sup>10</sup>, a SQL relational database server.
- *TimeScaleDB*<sup>11</sup>, a time-series SQL database, which is an extension of PostgreSQL.

Both tools offer similar features, in the sense that they are SQL relational databases. However, it is not possible to use only one of them or replace them by MySQL to have only one tool for relational data storage because they were especially requested by some pilots to facilitate the integration with other technological systems.

Moreover, according to the needs of the pilot and solution partners, it has been decided to offer a variety of configurations or *scenarios* of each tool. These scenarios can be seen as different ways of deploying the tools to meet some requirements related to aspects such as performance or security. Generally speaking, for each tool, the following scenarios are considered:

- A first basic *single server (SS)* scenario, which offers the tool in its most basic version. It leverages a single instance of the tool, ready to be used. It is worth noting that such single instances do not consider any security issue beyond its default configuration. Thus, an "SS" scenario is only recommended for the development tasks of a pilot or solution.
- A *single server with TLS security (SS+Sec)* scenario, that offers the tool with a security configuration based on the use of TLS (with x509 certificates). This scenario is suitable for production settings in which a single instance of the tool suffices to provide a secure

---

<sup>3</sup> Cassandra. <https://cassandra.apache.org>

<sup>4</sup> MinIO. <https://www.min.io>

<sup>5</sup> MongoDB. <https://www.mongodb.com>

<sup>6</sup> MySQL. <https://www.mysql.com>

<sup>7</sup> Neo4J. <https://www.neo4j.com>

<sup>8</sup> Redis. <https://redis.io>

<sup>9</sup> MariaDB. <https://mariadb.org>

<sup>10</sup> PostgreSQL. <https://www.postgresql.org>

<sup>11</sup> TimeScaleDB. <https://www.timescale.com/>





and stable service.

- A *high availability (HA)* scenario, which offers the tool in a high availability mode. This typically involves defining some *cluster* or *replica set* environment, composed by instances of the tool that work in a cooperative mode (for instance, a cluster of replicas of a database server). Again, these instances do not offer any security mechanism beyond those that are offered by default and are, thus, only recommended for development purposes.
- Finally, a *high availability with TLS security (HA+Sec)* scenario, which extends the "HA" scenario with a security configuration based on the use of TLS. This is the recommended option to offer a secure, fault-tolerant, highly available service.

For each scenario, *software artifacts* are provided, including configuration files, Docker Compose<sup>12</sup> orchestration files, and shell-scripts for Bash. Moreover, additional common configuration files and shell scripts are provided. The purpose of such artifacts is to deploy the selected storage tools according to the features of one of the four scenarios described above. Deliverable "*D3.8 - i4Q Data Repository*" (section 3.1 "Current implementation") provides additional details about these artifacts. Moreover, it provides a summary about the status of the development of the scenarios of each tool.

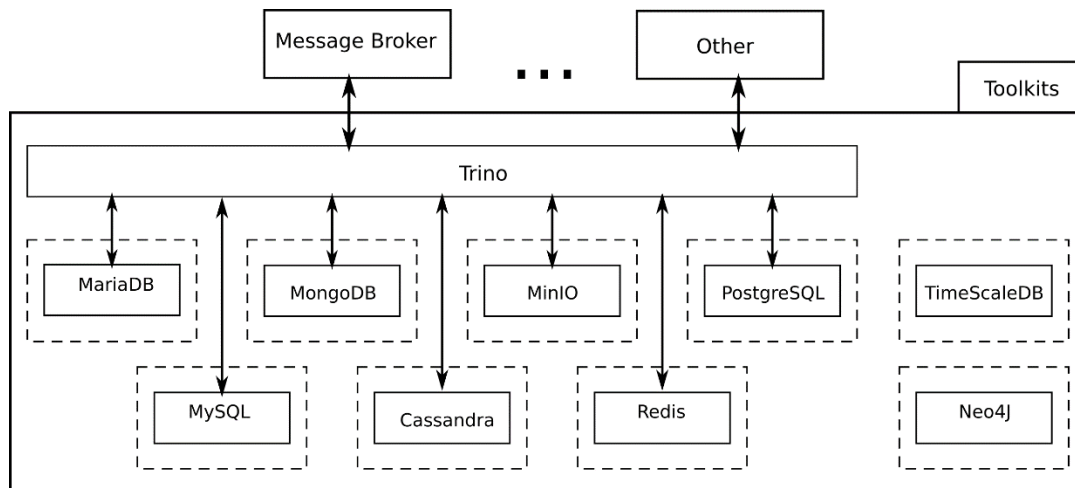
Furthermore, the *i4Q<sup>DR</sup>* includes another tool, Trino<sup>13</sup>, which is a highly parallel and distributed ANSI SQL-compliant open-source query engine. These features enable a very efficient performance, even in the case of many simultaneous queries. Trino offers a relational-like view of a number of data storage tools including Cassandra, MariaDB, MongoDB, MySQL, PostgreSQL, Redis and others.

From an architectural perspective, Trino will be deployed as a layer on top of the different artifacts mentioned above, as illustrated in Figure 5. Note that there are two storage technologies that do not appear below the Trino layer. On the one hand, there is Neo4J, as this storage tool is not currently supported by Trino. On the other hand, there is TimeScaleDB, a storage technology that has been included during M24 to cover a possible requirement for the generic pilot defined in T6.7. The list of Trino connectors does not include this technology. However, the plan is to check whether the connector for PostgreSQL could be used. Therefore, in these two cases, access to the storage tool will not be performed via Trino and, instead, it may require a direct connection or a dedicated access mechanism.

---

<sup>12</sup> Docker Compose documentation. <https://docs.docker.com/compose/>

<sup>13</sup> Trino. <https://trino.io>



**Figure 5.** Architectural overview of i4Q<sup>DR</sup>

In the context of the i4Q<sup>DR</sup>, besides its efficient performance, the most important benefits of using Trino are related to interoperability aspects. Firstly, it enhances interoperability at an internal level. In this regard, Trino adds a level of abstraction that allows the pilots and solutions to have a relational SQL-based view of different data storage tools, no matter what type of data they handle. This also allows to build *federated views* of data that is managed by a heterogeneous set of tools, which in turns allows to build more complex and interoperable software components.

Secondly, Trino also enhances the interoperability of the i4Q<sup>DR</sup> at an external level since it facilitates its interaction with other i4Q solutions. More specifically, Trino allows the i4Q<sup>DR</sup> to offer a common communication interface for all the selected storage tools. For instance, Trino offers a REST API which will allow the interaction of the i4Q<sup>DR</sup> with other solutions by means of the HTTP protocol.

## 4. Communication with Data Repositories

---

Just like for any software component, an important feature of data repositories is the communication and interaction with other software components, as mentioned in Section 1. This action typically involves the exchange of data among the two parties. In the specific case of data repositories, such a communication usually takes the form of a solution willing to store some data, or to retrieve it.

There are several approaches to the implementation of the communication of data repositories with other software components. One of the possibilities is to develop a general-purpose backend application offering a REST API to other components. In this approach, communications with the data repository are performed only at an internal level between the backend application and the data repository. Note that in this case external software components are not allowed to communicate directly with the data repository. Instead, they send HTTP requests to the endpoints defined in the REST API to get data from the repository, or store new data into it, or to update some fields of data already stored in the repository.

The approach of implementing REST APIs to allow interaction among different services is quite common in web programming. However, it might not be the most efficient one in the case of data repositories. Especially, when huge amounts of data are involved, and/or when the need to exchange data is very frequent, as it happens in industrial scenarios, where machines, sensors and other devices often generate data in a continuous way. In this type of scenarios, it is more appropriate to use tools specifically designed and implemented for these requirements. One of the most well-known tools used for this use cases is Apache Kafka<sup>14</sup>.

Apache Kafka is an open-source platform for event streaming capable of handling huge amounts of messages in real-time. It is used in all sorts of real-time data streaming applications as it is highly scalable while providing high data throughput with low latency.

As explained in [1] a deployment of Kafka consists of server and client applications that communicate via a high-performance TCP network protocol. Clients read, write, and process streams of messages (also called record or events in the documentation), whereas servers oversee the management and persistence of the messages exchanged among clients. There are two types of clients: producers, the ones that publish (write) messages to Kafka, and consumers, which subscribe to (read and process) these messages.

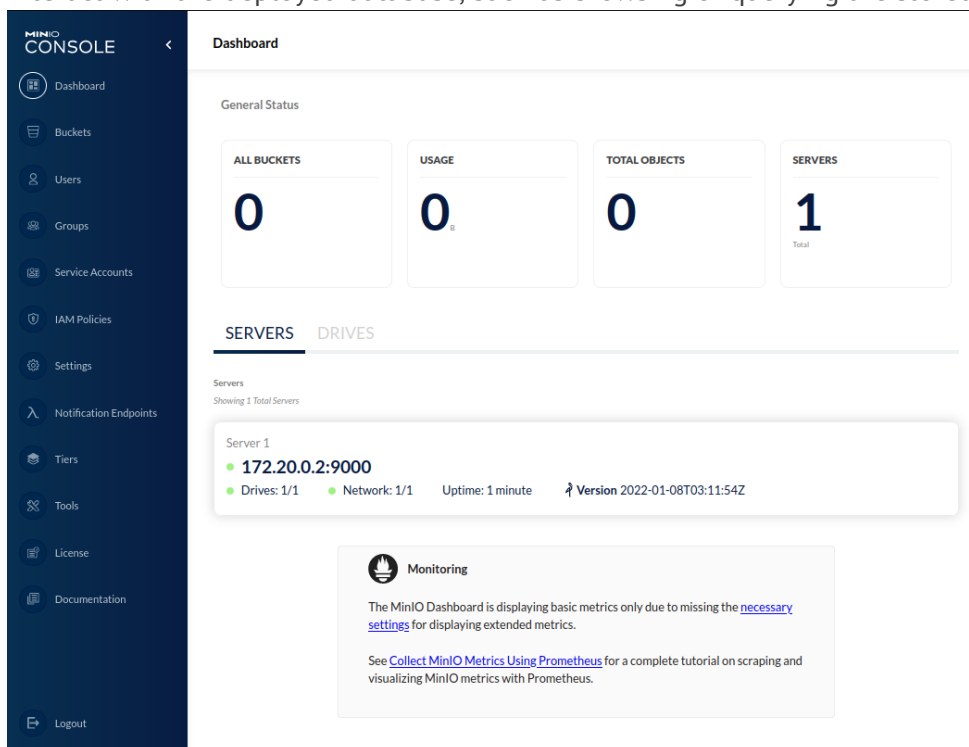
The communication in Kafka is conducted using Kafka topics. A topic is an abstraction that acts as an intermediary between producers and clients. More specifically, producers publish messages in a certain topic and clients can later consume these messages via subscribing to that specific topic. The data format of the messages can be either JSON or Avro which is a more efficient and compact way of exchanging messages, with a data model similar to JSON.

In this project the approach followed by the i4Q Data Repository (i4Q<sup>DR</sup>) is to support different ways of communication, so that other i4Q solutions or external software components can use the one that suits best their needs. The supported approaches are as follows:

---

<sup>14</sup>Apache Kafka: <https://kafka.apache.org/>

1. First, it is possible to interact directly with the server application that can be deployed for each one of the storage technologies. Such communication can be made by means of scripts or programs, possibly using libraries provided by third parties for that purpose, which depend on the chosen programming language.
2. Secondly, for some storage technologies, the *i4Q<sup>DR</sup>* also deploys an instance of a specific graphical client application. This is the case, for instance, of MinIO or MongoDB scenarios, when the *i4Q<sup>DR</sup>* deploys an instance of MiniO<sup>15</sup> console (see Figure 6), and Mongo Express<sup>16</sup> (see Figure 7), respectively. These applications allow a more user-friendly mechanism to interact with the deployed database, such as browsing or querying the stored data.



**Figure 6.** MinIO Console deployed by the *i4Q<sup>DR</sup>* for any MinIO scenario.

3. However, they are not intended for an intensive use, such as for inserting big amounts of new data.
4. Since the *i4Q<sup>DR</sup>* deploys an instance of Trino, the third possibility is to interact with the *i4Q<sup>DR</sup>* by using the mechanisms provided by Trino for that purpose. More specifically, it is possible to send queries to Trino and receive results, or otherwise interact with Trino and the connected data sources by means of the official clients provided at Trino's official website<sup>17</sup>, or by some others developed by the community<sup>18</sup> for platforms such as Python, and these can in turn be used to connect applications using these platforms. Furthermore,

<sup>15</sup> Minio Console website: <https://min.io/docs/minio/linux/administration/minio-console.html>

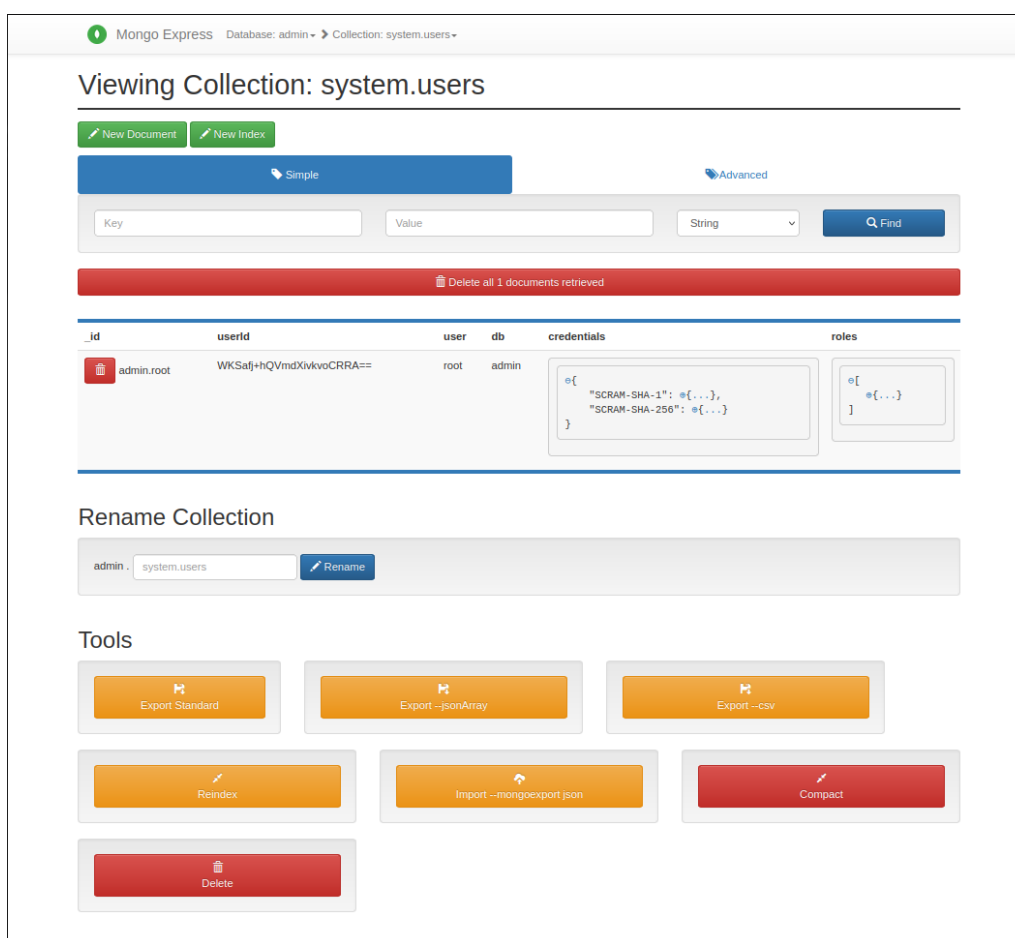
<sup>16</sup> Mongo Express GitHub repository: <https://github.com/mongo-express/mongo-express>

<sup>17</sup> Available Official clients for Trino: <https://trino.io/docs/current/client.html>

<sup>18</sup> Information on other clients for Trino provided by the community: <https://trino.io/resources.html>

Trino provides a client REST API<sup>19</sup> that allows submitting queries by making HTTP requests. However, the preferred method to interact with Trino is using the existing clients above mentioned.

Finally, in the context of this project, there is another possible way for communicating with the *i4Q<sup>DR</sup>*: by means of the *i4Q* Message Broker, an additional *i4Q* component developed to achieve solution interoperability by establishing a communication channel to exchange information and data among *i4Q* solutions. The Message Broker is responsible to provide a fast and secure way of inter-solution communication through data streaming and is based on Apache Kafka, as distributed by the Confluent platform. Further details regarding the integration of the *i4Q<sup>DR</sup>* with the *i4Q* Message Broker are provided in deliverable *D3.16 - Data Repository v2*, due at the end of M24.



**Figure 7.** Mongo Express instance deployed by the *i4Q<sup>DR</sup>* for any MongoDB scenario.

<sup>19</sup> Trino's REST API: <https://trino.io/docs/current/develop/client-protocol.html>

## 5. Implementation of i4Q<sup>DRG</sup> as web documentation

In order to improve the readability of the information gathered in this document, a simple web application has been implemented to present its content as web documentation, similarly as the i4Q Manufacturing Line Reconfiguration Guideline (i4Q<sup>LRG</sup>) solution.

The main technologies used to implement such a web application are Sphinx<sup>20</sup> and Docker<sup>21</sup>. On the one hand, Sphinx is a library that was originally created for Python documentation, but that can be used to document software projects in a variety of languages. More specifically, it converts reStructuredText<sup>22</sup> files into HTML websites. In the other hand, to allow an easy deployment of the web application, it is provided as a Docker container that can be easily created and bootstrapped.

Figure 8 shows a snapshot of the web application displaying the most relevant information gathered in this document. Note that there is a side menu on the right part of the screen, which provides a more user-friendly navigation throughout the content.

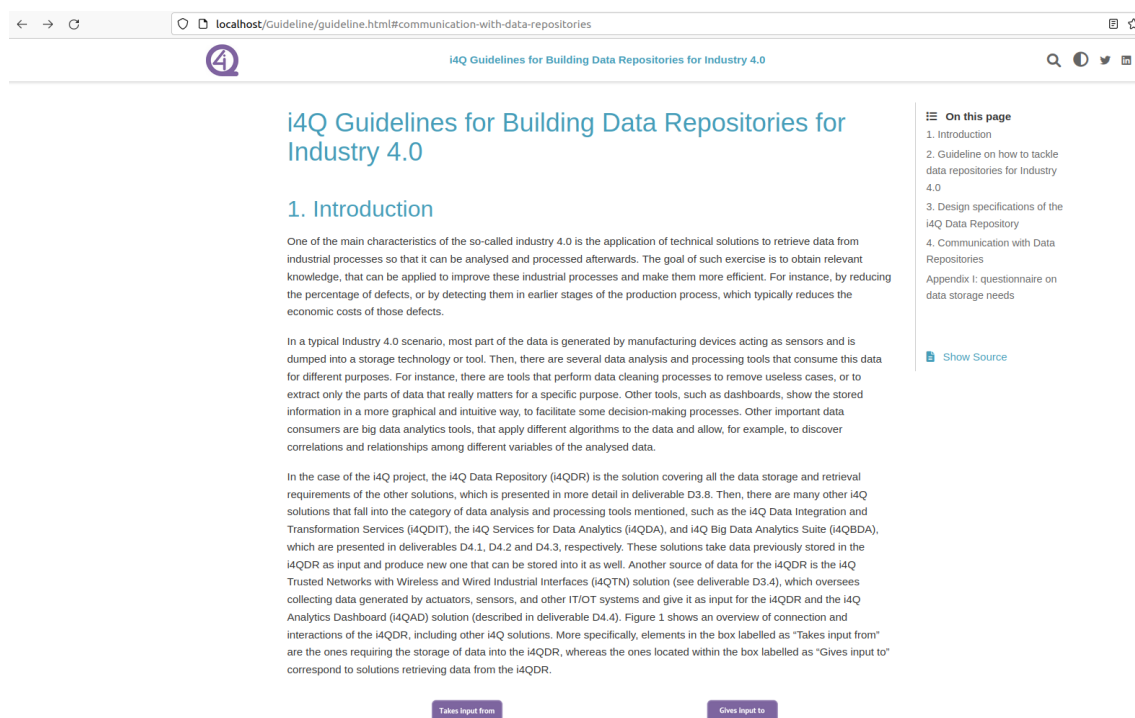


Figure 8. Snapshot of web application showing information regarding the i4Q<sup>DRG</sup> solution.

Furthermore, Sphinx allows searching for words along all the text. This can be done by clicking on the “loupe” icon, in the right-top corner of the screen. If the user looks for a word, the application shows the sections where that word appears, as shown in Figure 10. Once the user

<sup>20</sup> Sphinx website: <https://www.sphinx-doc.org/en/master/>

<sup>21</sup> Docker website: <https://www.docker.com/>

<sup>22</sup> reStructuredText is a file format for textual data used primarily in the Python programming language community for technical documentation.



clicks in one of those results, she is redirected to the beginning of that section, and the searched word is highlighted in yellow (see Figure 9).



Figure 9. Example of highlighting of searched word

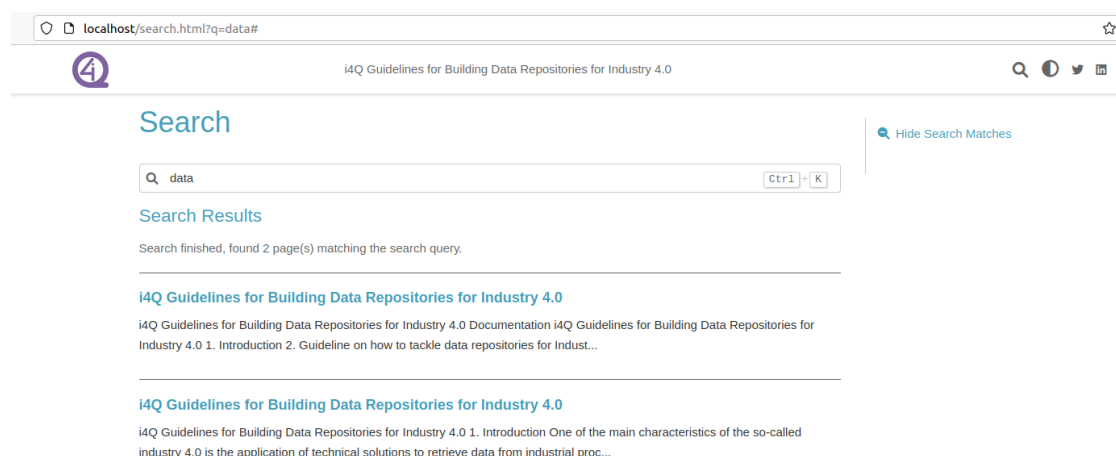


Figure 10. Example of search results.

The source code of this web application is available at an internal GitLab repository<sup>23</sup>. Further information on how to deploy it is provided in Appendix II: web documentation.

<sup>23</sup> Internal GitLab repository of i4Q DRG solution: <https://gitlab.com/i4q/drg>

## 6. Conclusions

---

This document provided an overview of the importance of data repositories in Industry 4.0 contexts. In Section 1 this explanation was illustrated with the *i4Q<sup>DR</sup>*, the technical solution that will cover the data management requirements of any solution or pilot in the context of the *i4Q* project.

Moreover, this deliverable gathers the main problems and challenges arising when developing such data repositories. In this regard, a non-exhaustive list of recommendations addressing these problems and challenges has been proposed with the purpose of serving as a guideline for the development of data repositories. Indeed, Section 3, as an illustrative example, describes in detail how these recommendations have been applied in practice during the design of the *i4Q<sup>DR</sup>* solution, and the definition of its development strategy.

Communication with data repositories from/to other software solutions is an important issue tackled in this document. In this regard, Section 4 provided an overview of the most common approaches to implement this issue, and briefly explains how this has been done in the case of the *i4Q<sup>DR</sup>* solution.

In order to provide the information gathered in this deliverable in a more user-friendly fashion, a web application has been implemented to show its most relevant content in the style of standard web documentation, as explained in Section 5.

This deliverable is the second and final version of the series of deliverables regarding the *i4Q Guidelines for Building Data Repositories for Industry 4.0*. Basically, it extends its previous version (D3.7) with the content provided in Sections 4 and 5. This document provides feedback to deliverable “D3.16 - *i4Q Data Repository v2*”, due on M24.



## 7. References

---

- [1] «Apache Kafka official website,» [Online]. Available: <https://kafka.apache.org>. [Accesso: November 2022].

## Appendix I: questionnaire on data storage needs

---

This appendix provides further information on the questionnaire distributed to the responsables of pilots and other i4Q solutions. The purpose of this questionnaire was gathering more specific and technical details on their data storage needs. The information provided by the partners that filled in this questionnaire was used to select the most suitable tools and technologies to implement the i4Q<sup>DR</sup>.

The questionnaire is a Microsoft Excel document containing different sheets. Besides some sheets providing an overview on the survey and useful information on how to fill in it, there is a sheet for each pilot and solution. These sheets contain a set of questions on different issues, as briefly explained in Section 3. Since this questionnaire is stored in a private repository, several figures have been included in this section to show the concrete questions that were asked to the partners. More specifically:

1. *Figure 11 shows the questions about the type of storage that is necessary.*
2. *Figure 12 gathers the questions on the expected volume of data to be handled*
3. *Figure 13 reflects the questions on the expected type of interaction with the i4Q<sup>DR</sup> (interactive, batch, stream-oriented, etc.), and the communication interfaces with the solution that may be used (HTTP, HTTPS, TCP, etc.).*
4. *Figure 14 shows questions regarding other communication issues, and the needs related to data replication.*
5. *Figure 15 gathers the questions in relation to authorization issues, deployment of the solution, and its expected performance.*
6. *Figure 16 shows the questions asked to collect information on whether the partners have specific needs regarding tools or technologies to use, possible constraints and preferences about the premises where the solution will be deployed, or other concrete needs, such as data anonymisation or any other proposed by the partners.*

1	<b>1 Storage types</b>	
2		
3	The goal is to know what tools and technologies are needed by the pilot or solution.	
4		
5	<b>1.1</b>	Select which of the following storage types are needed. If specific tools are needed (eg. "MongoDB"), write down their names. Also specify if the need of such tool is strict or flexible (eg. can be replaced by an equivalent one).
6		
7	1.1.1	SQL storage (eg. MySQL, Oracle, MS SQL Server, etc.): <input type="text" value="Yes/No/I don't know/Specify which one"/>
8		
9	1.1.2	Specific SQL dialect: <input type="text" value="Yes/No/I don't know/Specify which one"/>
10		
11	1.1.3	Structured non-SQL storage (eg. Cassandra, HBase, etc.): <input type="text" value="Yes/No/I don't know/Specify which one"/>
12		
13	1.1.4	JSON storage (eg. MongoDB): <input type="text" value="Yes/No/I don't know/Specify which one"/>
14		
15	1.1.5	Key-value storage (eg. etcd, Redis, memcached, etc.): <input type="text" value="Yes/No/I don't know/Specify which one"/>
16		
17	1.1.6	Other storage (eg. graph databases, object-oriented, etc.): <input type="text" value="Yes/No/I don't know/Specify which one"/>
18		
19	1.1.7	Storage for regular files: <input type="text" value="Yes/No/I don't know/Specify which one"/>
20		
21	1.1.8	A storage for very large files: <input type="text" value="Yes/No/I don't know/Specify which one"/>
22		
23	1.1.9	Document Management System-like storage (with content indexing, searches, categories, tags etc.): <input type="text" value="Yes/No/I don't know/Specify which one"/>
24		
25	1.1.10	Other: <input type="text" value="Specify which ones"/>
26		
27	<b>1.2</b>	Other comments: <input type="text"/>
28		

**Figure 11.** Questions on necessary type of storage

28	<b>2 Expected data volumes</b>	
29		
30	The goal is to have an approximate idea of the volumes of data the tools and technologies should be able to handle, in order to choose the proper tools but also to properly estimate the resources that may be needed (servers, main memory, disk capacities, bandwidth, etc.).	
31		
32	2.1	Do you have an estimation (or exact numbers) of the volume of the data managed by the pilot or solution? <input type="text" value="Specify which (*)"/>
33		
34	2.2	Do you have an estimation (or exact numbers) of the rate at which their volume increases daily/weekly/monthly/yearly/etc.? <input type="text" value="Specify which (*)"/>
35		
36		
37		
38		
39		
40		
41	<b>3</b>	<b>Operation modes</b>

**Figure 12.** Questions on expected data volumes

39

40 **3 Operation modes**

41 The goal is to know how the pilot or solution should access the Data Repository.

42

43

44 3.1 Select which of the following operation modes

45 are needed (select as many as needed):

46 3.1.1 Interactive on-demand (eg. a human user

47 interactively executes queries, updates, Yes/I don't know

48 3.1.2 Interactive (eg. a software agent executes

49 queries, updates, etc.) Yes/I don't know

50 3.1.3 Batch (eg. a scheduled tasks executes a

51 series of queries, updates, etc.) Yes/I don't know

52 3.1.4 Streaming needs (eg. a stream of data is

53 received from a sensor) Yes/I don't know

54 3.1.5 Publish/subscriber behavior (eg. data is

55 published to queues, topics, etc. and

56 asynchronously consumed at any later

57 3.1.6 Other Specify which ones

58 3.2 Other comments:

59

60 **4 Communication interfaces**

61

62 4.1 The storage has to be used through an HTTPS

63 channel (eg. REST API, non REST, etc.): Yes/No/I don't know

64 4.2 The storage has to be used through a binary

65 TCP channel: Yes/No/I don't know

66 4.4 Other comments:

67

68

69

4-FACTOR 5-RIAS 6-FARP 1-DQG 2-QE 3-BC 4-TN 5-CSG 6-SH 7-DRG 8-DR 9-DIT

**Figure 13.** Questions on type of interaction and communication interfaces

67

68

69

70 **5 Other communication issues**

71

72 5.1 Security constraints to consider (eg. Public Key

73 Infrastructure issues, etc.): Specify which ones

74 5.2 Any in-the-middle agent or environment

75 needed (a proxy, a VPN, any other tunnel, etc.): Specify which ones

76 5.3 Other comments:

77

78 **6 Replication**

79

80 6.1 The data may or must be replicated: Yes/No/I don't know

81 6.2 There are synchrony replications constraints

82 (eg. as much synchrony as possible, decoupled

83 master-slave, etc.):

84 6.3 There are synchrony architecture constraint

85 (eg. multi-primary required, a single primary is

86 enough, etc.):

87 6.4 Other comments:

88

89

90

91

4-FACTOR 5-RIAS 6-FARP 1-DQG 2-QE 3-BC 4-TN 5-CSG 6-SH 7-DRG 8-DR

Listo Accesibilidad: es necesario investigar

**Figure 14.** Questions on other communication issues and data replication needs

88	7	Authorization	
89	7.1	The data must be accessed by users with different access levels (eg. administrators, privileged users, plain users, etc.):	Yes/No/I don't know
90	7.2	Which is the "permission grain" (eg. coarse - database level- vs fine -row level, column level, item level-)?:	Specify which ones
91	7.3	Other comments:	
92			
93			
94			
95	8	Deployment	
96	8.1	The storage must be installed in a specific environment (eg. Kubernetes, Docker, VirtualBox, etc., specific OS like Linux, Windows, etc.):	Specify which ones
97	8.2	Other comments:	
98			
99			
100			
101	9	Expected performance	
102	9.1	There are performances requirements (eg. number of reads or writes per second):	Specify which ones
103	9.2	Other comments:	
104			
105			
106			
107			
108			

4-FACTOR 5-RIAS 6-FARP 1-DQG 2-QE 3-BC 4-TN 5-CSG 6-SH 7-DRG 8-DR

Listo Accesibilidad: es necesario investigar

**Figure 15.** Questions about authorization, deployment, and performance

109	10	Specific tools, technologies, etc.	
110	10.1	Please list any other tool, technology, resource, etc. that may be needed by the pilot or solution regarding the Data Repository	1) 2) 3) 4) 5) ...
111			
112	11	Open question	
113	11.1	Please write down any other issue, concern, etc. regarding the relation between the Data Repository and the pilot or solution	
114			
115	12	Last minute	
116	12.1	Do you have any restrictions on even preferences about data storage in case of being on the cloud or on-premise?	Specify which ones
117	12.2	Is there any requirement for data anonymization?	Specify which ones
118			
119			
120			
121			
122			
123			
124			
125			
126			
127			
128			
129			
130			
131			
132			
133			
134			

4-FACTOR 5-RIAS 6-FARP 1-DQG 2-QE 3-BC 4-TN 5-CSG 6-SH 7-DRG 8-DR

Listo Accesibilidad: es necesario investigar

**Figure 16.** Questions about other specific needs and possible constraints and preferences



## Appendix II: web documentation

---

**i4Q Data Repository Guidelines** (i4Q<sup>DRG</sup>) web documentation can be accessed online at: [http://i4q.upv.es/7\\_i4Q\\_DRG/index.html](http://i4q.upv.es/7_i4Q_DRG/index.html). This documentation includes instructions on how to deploy and use the web application presented in Section 5.