























## BRIEF REPORT

# FAIR4Health: Findable, Accessible, Interoperable and Reusable data to foster Health Research [version 1; peer review: 3 approved with reservations]

Celia Alvarez-Romero <sup>1</sup>, Alicia Martínez-García<sup>1</sup>, A. Anil Sinaci <sup>2</sup>,  
Mert Gencturk <sup>2</sup>, Eva Méndez <sup>3</sup>, Tony Hernández-Pérez <sup>3</sup>, Rosa Liperoti<sup>4</sup>,  
Carmen Angioletti <sup>4</sup>, Matthias Löbe <sup>5</sup>, Nagarajan Ganapathy <sup>6</sup>,  
Thomas M. Deserno<sup>6</sup>, Marta Almada <sup>7</sup>, Elisio Costa <sup>7</sup>, Catherine Chronaki <sup>8</sup>,  
Giorgio Cangioli<sup>8</sup>, Ronald Cornet <sup>9</sup>, Beatriz Poblador-Plou <sup>10</sup>,  
Jonás Carmona-Pérez <sup>10</sup>, Antonio Gimeno-Miguel<sup>10</sup>, Antonio Poncel-Falcó<sup>11</sup>,  
Alexandra Prados-Torres<sup>10</sup>, Tomi Kovacevic<sup>12,13</sup>, Bojan Zaric <sup>12,13</sup>, Darijo Bokan<sup>13</sup>,  
Sanja Hromis<sup>12,13</sup>, Jelena Djekic Malbasa<sup>12,13</sup>, Carlos Rapallo Fernández<sup>14</sup>,  
Teresa Velázquez Fernández<sup>14</sup>, Jessica Rochat<sup>15</sup>,  
Christophe Gaudet-Blavignac <sup>15</sup>, Christian Lovis <sup>15</sup>, Patrick Weber<sup>16</sup>,  
Miriam Quintero<sup>17,18</sup>, Manuel M. Perez-Perez<sup>17,18</sup>, Kevin Ashley <sup>19</sup>,  
Laurence Horton <sup>20</sup>, Carlos Luis Parra Calderón <sup>1</sup>

<sup>1</sup>Group of Research and Innovation in Biomedical Informatics, Biomedical Engineering and Health Economy, Institute of Biomedicine of Seville, IBiS / Virgen del Rocío University Hospital / CSIC / University of Seville, Seville, 41013, Spain

<sup>2</sup>SRDC Software Research Development and Consultancy Corporation, Ankara, 06800, Turkey

<sup>3</sup>Dept. of Library & Inf Sci. Universidad Carlos III de Madrid, Getafe, 28903, Spain

<sup>4</sup>Department of Geriatric and Orthopedic Sciences, Catholic University of Sacred Heart, Roma, 00168, Italy

<sup>5</sup>Institute for Medical Informatics (IMISE), University of Leipzig, Leipzig, 04107, Germany

<sup>6</sup>PLRI Institute for Medical Informatics of TU Braunschweig and Hannover Medical School, Braunschweig, 38106, Germany

<sup>7</sup>Ucibio Requitme, Faculty of Pharmacy University of Porto. Porto4Ageing, Porto, 4050-313, Portugal

<sup>8</sup>HL7 Europe Foundation, Brussels, 1000, Belgium

<sup>9</sup>Amsterdam UMC, University of Amsterdam, Medical Informatics, Amsterdam Public Health, Amsterdam, 1105AZ, The Netherlands

<sup>10</sup>EpiChron Research Group, Aragon Health Sciences Institute (IACS), IIS Aragón, Miguel Servet University Hospital, Zaragoza, 50009, Spain

<sup>11</sup>EpiChron Research Group, Aragon Health Sciences Institute (IACS), IIS Aragón, Aragon Health Service, Zaragoza, 50009, Spain

<sup>12</sup>Medical Faculty University of Novi Sad, Novi Sad, 21000, Serbia

<sup>13</sup>Institute for Pulmonary Diseases of Vojvodina, Sremska Kamenica, 21204, Serbia

<sup>14</sup>J&A Garrigues, S.L.P., Seville, 41013, Spain

<sup>15</sup>University of Geneva and University hospitals of Geneva, Geneva, 1211, Switzerland

<sup>16</sup>Nice Computing SA Le Mont-sur-Lausanne, Le Mont-sur-Lausanne, 1052, Switzerland

<sup>17</sup>Atos Research and Innovation - ARI. Atos IT., Madrid, 28037, Spain

<sup>18</sup>Atos Research and Innovation - ARI. Atos Spain., Madrid, 28037, Spain

<sup>19</sup>Digital Curation Centre, University of Edinburgh, Argyle House, Edinburgh, EH3 9DR, UK

<sup>20</sup>Digital Curation Centre, University of Glasgow, Glasgow, G12 8QQ, UK

---

First published: 09 Mar 2022, 2:34

---

<https://doi.org/10.12688/openreseurope.14349.1>

**Latest published:** 31 May 2022, 2:34

<https://doi.org/10.12688/openreseurope.14349.2>

## Abstract

Due to the nature of health data, its sharing and reuse for research are limited by ethical, legal and technical barriers. The FAIR4Health project facilitated and promoted the application of FAIR principles in health research data, derived from the publicly funded health research initiatives to make them Findable, Accessible, Interoperable, and Reusable (FAIR). To confirm the feasibility of the FAIR4Health solution, we performed two pathfinder case studies to carry out federated machine learning algorithms on FAIRified datasets from five health research organizations. The case studies demonstrated the potential impact of the developed FAIR4Health solution on health outcomes and social care research. Finally, we promoted the FAIRified data to share and reuse in the European Union Health Research community, defining an effective EU-wide strategy for the use of FAIR principles in health research and preparing the ground for a roadmap for health research institutions to offer access to certified FAIR datasets.

This scientific report presents a general overview of the FAIR4Health solution: from the FAIRification workflow design to translate raw data/metadata to FAIR data/metadata in the health research domain to the FAIR4Health demonstrators' performance.

## Keywords

FAIR principles, health research data management, HL7 FHIR, health data, data sharing, data reuse, health research, open science, privacy-preserving computing, machine learning.



This article is included in the [Health Sciences](#) gateway.



This article is included in the [Research on Research](#) gateway.





This article is included in the [Digital Health and Well-being](#) collection.

## Open Peer Review

**Approval Status** ✓ ? ?

|                  | 1                    | 2                    | 3                    |
|------------------|----------------------|----------------------|----------------------|
| <b>version 2</b> |                      |                      |                      |
| (revision)       | ✓                    |                      |                      |
| 31 May 2022      | <a href="#">view</a> |                      |                      |
|                  | ↑                    |                      |                      |
| <b>version 1</b> | ?                    | ?                    | ?                    |
| 09 Mar 2022      | <a href="#">view</a> | <a href="#">view</a> | <a href="#">view</a> |

1. **Anupama Gururaj** , National Institutes of Health, Rockville, USA
2. **Gary Saunders**, EATRIS ERIC, Amsterdam, The Netherlands
3. **Salvador Capella-Gutierrez** , Barcelona Supercomputing Center (BSC), Barcelona, Spain

Any reports and responses or comments on the article can be found at the end of the article.



This article is included in the [Healthcare and Policy](#) collection.

**Corresponding author:** Celia Alvarez-Romero ([celia.alvarez@juntadeandalucia.es](mailto:celia.alvarez@juntadeandalucia.es))

**Author roles:** **Alvarez-Romero C:** Investigation, Methodology, Validation, Visualization, Writing – Original Draft Preparation; **Martínez-García A:** Investigation, Methodology, Validation, Visualization, Writing – Review & Editing; **Sinaci AA:** Investigation, Software, Writing – Review & Editing; **Gencturk M:** Investigation, Software, Writing – Review & Editing; **Méndez E:** Conceptualization, Investigation, Writing – Review & Editing; **Hernández-Pérez T:** Conceptualization, Investigation, Writing – Review & Editing; **Liperoti R:** Investigation, Validation, Writing – Review & Editing; **Angioletti C:** Investigation, Validation, Writing – Review & Editing; **Löbe M:** Investigation, Writing – Review & Editing; **Ganapathy N:** Investigation, Writing – Review & Editing; **Deserno TM:** Investigation, Writing – Review & Editing; **Almada M:** Investigation, Validation, Writing – Review & Editing; **Costa E:** Investigation, Validation, Writing – Review & Editing; **Chronaki C:** Investigation, Writing – Review & Editing; **Cangioli G:** Investigation, Writing – Review & Editing; **Cornet R:** Investigation, Writing – Review & Editing; **Poblador-Plou B:** Investigation, Validation, Writing – Review & Editing; **Carmona-Pérez J:** Investigation, Validation, Writing – Review & Editing; **Gimeno-Miguel A:** Investigation, Validation, Writing – Review & Editing; **Poncel-Falcó A:** Investigation, Validation, Writing – Review & Editing; **Prados-Torres A:** Investigation, Validation, Writing – Review & Editing; **Kovacevic T:** Investigation, Validation, Writing – Review & Editing; **Zaric B:** Investigation, Validation; **Bokan D:** Investigation, Validation, Writing – Review & Editing; **Hromis S:** Investigation, Validation, Writing – Review & Editing; **Djekic Malbasa J:** Investigation, Validation, Writing – Review & Editing; **Rapallo Fernández C:** Conceptualization, Investigation, Writing – Review & Editing; **Velázquez Fernández T:** Conceptualization, Investigation, Writing – Review & Editing; **Rochat J:** Investigation, Validation, Writing – Review & Editing; **Gaudet-Blavignac C:** Investigation, Validation, Writing – Review & Editing; **Lovis C:** Investigation, Validation, Writing – Review & Editing; **Weber P:** Investigation, Writing – Review & Editing; **Quintero M:** Investigation, Software, Writing – Review & Editing; **Perez-Perez MM:** Investigation, Writing – Review & Editing; **Ashley K:** Writing – Review & Editing; **Horton L:** Writing – Review & Editing; **Parra Calderón CL:** Conceptualization, Investigation, Project Administration, Supervision, Writing – Review & Editing

**Competing interests:** No competing interests were disclosed.

**Grant information:** This research was financially supported by the European Union's Horizon 2020 research and innovation programme under the grant agreement No 824666 (project FAIR4Health). Also, this research has been co-supported by the Carlos III National Institute of Health, through the IMPaCT Data project (code IMP/00019), and through the Platform for Dynamization and Innovation of the Spanish National Health System industrial capacities and their effective transfer to the productive sector (code PT20/00088), both co-funded by European Regional Development Fund (FEDER) 'A way of making Europe'.

*The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.*

**Copyright:** © 2022 Alvarez-Romero C *et al.* This is an open access article distributed under the terms of the [Creative Commons Attribution License](#), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

**How to cite this article:** Alvarez-Romero C, Martínez-García A, Sinaci AA *et al.* **FAIR4Health: Findable, Accessible, Interoperable and Reusable data to foster Health Research [version 1; peer review: 3 approved with reservations]** Open Research Europe 2022, 2:34 <https://doi.org/10.12688/openreseurope.14349.1>

**First published:** 09 Mar 2022, 2:34 <https://doi.org/10.12688/openreseurope.14349.1>

## Plain language summary

Health research organizations work more and more with health data. The reuse of health data has significant benefits for society, both financially and for our well-being. However, there are significant barriers to data sharing, which our project has taken steps to overcome. The FAIR principles of Findability, Accessibility, Interoperability and Reusability are intended to influence such institutions to support collaborative use of data. FAIR4Health promotes the application of FAIR principles in health research data derived from public projects. FAIR4Health developed a workflow and tools to support the FAIR principles, and applied these to two case studies, extending across multiple health care sites, which confirmed feasibility.

## Introduction

One of the more significant challenges of data-intensive science is to facilitate the breakthrough of knowledge by assisting humans and machines in the discovery, access, integration, and analysis of task-appropriate scientific data and their associated algorithms and workflows, facilitating reproducibility of the research.

The FAIR guiding principles describe distinct considerations for contemporary data publishing environments with respect to supporting both manual and automated deposition, exploration, sharing, and reuse, describing a set of guiding principles to make data Findable, Accessible, Interoperable, and Reusable<sup>1</sup>. Furthermore, the FAIR principles ensure that data are shared to enable and enhance reuse by humans and machines. Although FAIR emerged from a workshop for the life science community, the principles are intended to be applied to data and metadata from all disciplines.

Since their formal release via the [FORCE11](#) community, FAIR principles have been adopted by several funders and governments worldwide. The European Commission data management guidelines were updated in 2017 to introduce the notion of FAIR. The European Open Science Cloud (EOSC) Declaration and recent EOSC Strategic Research and Innovation Agenda (EOSC SRIA) both emphasise the central role of FAIR data.

In addition, it is essential to refer to the report issued by the European Union about the costs of NOT having FAIR data<sup>2</sup>, whose main conclusions are that: i) the cost of NOT having FAIR data is approximately €10.2bn per year for the EU; ii) in addition, the open data economy suggests that the impact on innovation of FAIR could add another €16bn to the minimum cost estimated; and iii) that would make a total of at least €26.2bn per year.

A diverse range of research disciplines are adopting FAIR principles. Several groups have been assessing FAIR uptake to date and the challenges being encountered. In the same way, the [FAIR4Health project](#), which has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement No 824666, promotes the application of FAIR principles to health research data.

## Methods

First of all, we performed a comprehensive analysis of current barriers, facilitators and potential overcoming mechanisms in the EU to implement a FAIR data policy in health research institutions. Information from different perspectives (technical, ethical, security, legal, cultural, behavioral and economic) was gathered to generate guidelines providing an optimal strategy for implementing this policy in EU health research institutions. Concretely, a FAIR4Health public deliverable<sup>3</sup> provided an analytical overview of all the considerations addressed to identify, report and overcome all the barriers that could prevent Health Research Performing Organizations (HRPOs) from opening, sharing and FAIRifying their research data.

Then, FAIR4Health designed a workflow<sup>3</sup> to apply the FAIR principles to health research data, as well as to Electronic Health Record data, based on the FAIRification process of [GO FAIR](#)<sup>4</sup>, but addressing the ethical, legal and technical aspects that health data include due to their sensitive nature by adding new steps in the workflow.

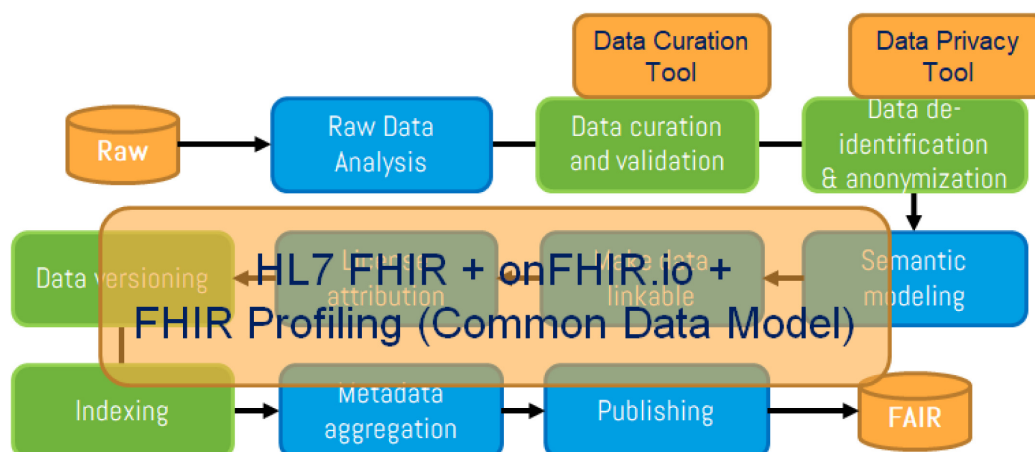
As shown in [Figure 1](#), new steps were included (in green) in the FAIR4Health FAIRification workflow to address these additional aspects through curation, validation and anonymization of sensitive health data. Adapted [GO FAIR steps](#) (in blue) define general actions for raw data analysis, license attribution, linking, semantic modeling, metadata management, and publishing to achieve FAIRness of existing (meta)data.

The requirements of health data were analysed in-depth, and FAIRification tools, based on the use of the [HL7 FHIR standard](#), were developed to obtain FAIR data from raw data resulting from biomedical research.

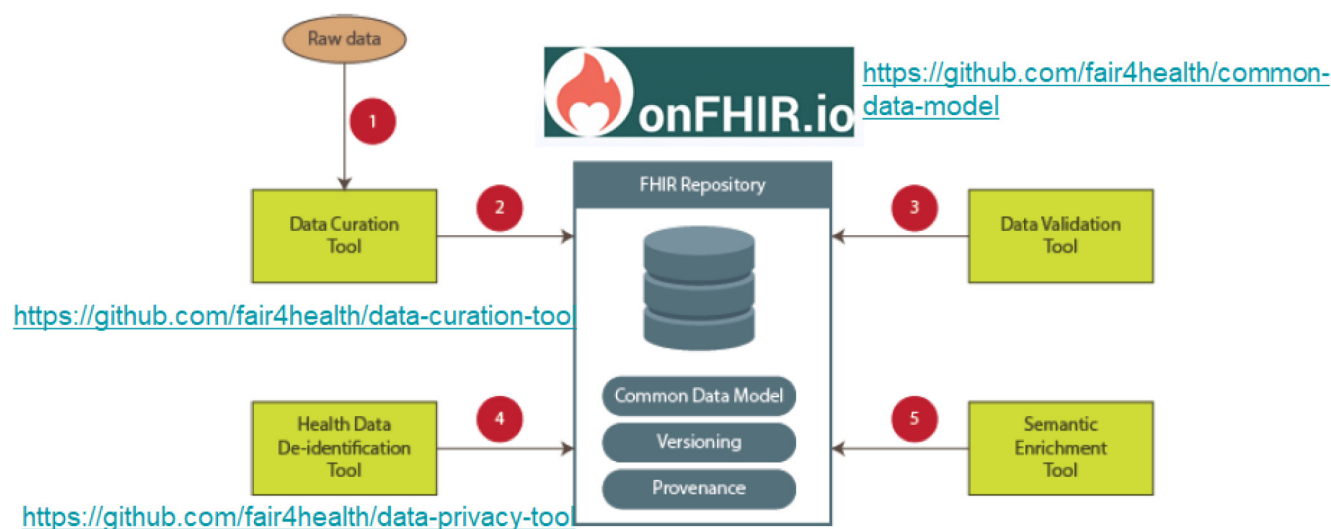
FAIRification tools are standalone, desktop applications developed by the FAIR4Health project to perform "Data curation and validation" and "Data de-identification and anonymization" steps of the FAIRification Workflow in an easier way:

- Data Curation Tool<sup>5</sup> is a highly specialized Extract-Transform-Load tool that can extract data from relational databases and spreadsheets, apply user-defined transformations, and load the transformed resources into an HL7 FHIR repository.
- Data Privacy Tool<sup>6</sup> is responsible for handling the privacy challenges on sensitive health data by applying several data de-identification and anonymization techniques. After the curation process, the Data Manager uses the Data Privacy Tool to de-identify data before making it available to other systems/components as FAIR data. This tool reads and writes de-identified resources back to the HL7 FHIR repository.

[Figure 2](#) shows the architecture implementing the FAIR4Health FAIRification Workflow for health data. In the core of architecture, an HL7 FHIR Repository acts as the health data repository. That way, the FAIR4Health core architecture,



**Figure 1. FAIR4Health workflow to apply FAIR principles in health research data.** The FAIR4Health FAIRification workflow, based on the GO FAIR process (steps in blue), includes new steps (in green) to address the additional considerations for health data through curation, validation and anonymization of sensitive health data. Then, FAIRification tools, based on the use of the HL7 FHIR standard, were developed to obtain FAIR data from raw data.



**Figure 2. FAIR4Health architecture implementing the FAIRification Workflow for health data.** At the core of architecture, an HL7 FHIR Repository acts as the health data repository. The FAIR4Health core architecture, which includes an FHIR Repository and is based on a Common Data Model, is an enabling factor for implementing the steps of the FAIR4Health FAIRification workflow in all aspects of FAIR principles.

including an FHIR Repository and based on a Common Data Model, is an enabling factor for implementing the steps of the FAIRification workflow in all aspects of FAIR principles. In FAIR4Health, onFHIR.io was utilized as the HL7 FHIR Repository deployed within the agents.

On top of these, the FAIR4Health Platform was developed to apply a Privacy-Preserving Distributed Data Mining (PPDDM) framework enabling health research organizations to perform joint data mining operations without exposing any sensitive

patient information to the outside world. In addition, the PPDDM Agent, which is responsible for running the data mining algorithms on top of the FAIRified data for the use cases defined by the user through the FAIR4Health Platform, was developed for training, validation and testing of models for the use cases defined. To achieve its objectives, the PPDDM Agent communicates with the onFHIR.io FHIR Repository within the data source boundaries, and the FAIR4Health Platform to exchange the results and predictive model information in a distributed manner.

The overall architecture of the FAIR4Health solution is shown in Figure 3.

## Results

The main objective of FAIR4Health was to facilitate and encourage the European Union Health Research community to FAIRify, share and reuse their datasets derived from publicly funded research initiatives through the demonstration of the potential impact that such a strategy has on health outcomes and health and social care research.

The FAIR4Health solution was validated with the two pathfinder case studies based on FAIRified data through the PPDDM framework.

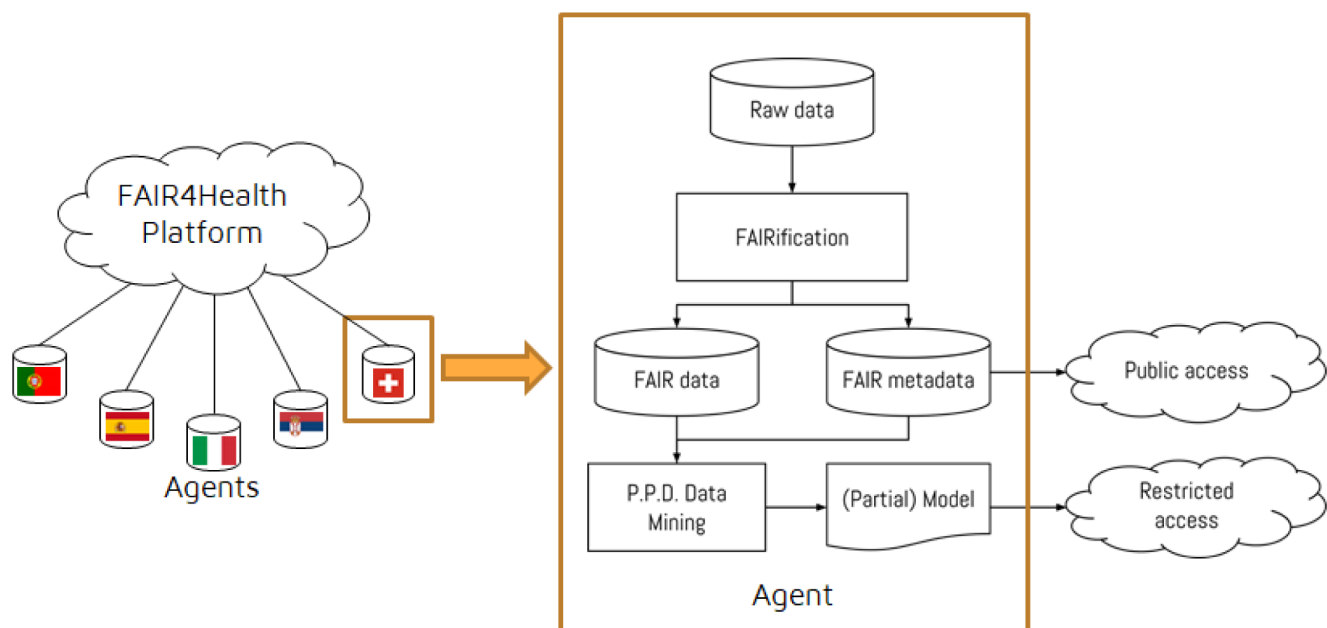
1. Identification of multimorbidity patterns and polypharmacy correlation on the risk of mortality in elderly.
2. Early prediction service for 30-days readmission risk in patients with Chronic Obstructive Pulmonary Disease (COPD).

The goal of these case studies was to test the developed tools in the project. The prototypes were developed making use of federated machine learning methodologies and algorithms implemented upon the FAIR4Health Platform. First, each health research dataset was FAIRified using the FAIR4Health FAIRification tools. Then, the federated machine learning

algorithms were trained and validated with retrospective datasets in both case studies. Finally, a prospective study was performed in the second use case to validate the developed model for prediction.

Concretely, the main goal of the pathfinder case study #1 was to analyze the impact of multimorbidity patterns and polypharmacy on the six-month mortality rate and cognitive impairment among elderly individuals in different health care settings. As a result, a multicentric retrospective observational study was designed in which data were collected from 5 different European cohorts. The population studied consisted of individuals aged 65 years or older with at least two chronic diseases. We used a frequent pattern tree association algorithm<sup>7</sup> implemented in the FAIR4Health Platform to identify the most frequent patterns in five different scenarios. The multimorbidity patterns obtained were consistent with previous studies<sup>8,9</sup>, which show the clinical potential of this method. We could also estimate a strong association between multimorbidity and polypharmacy and each of them with mortality.

COPD is one of the most prevalent chronic diseases. It has been associated with high morbidity and mortality and a high rate of readmission/rehospitalization and therefore associated with high healthcare costs. Thus, the main goal of the pathfinder case study #2 was to develop, validate and assess the accuracy of a clinical decision support tool for predicting 30-day readmission risk in patients suffering from COPD at discharge.



**Figure 3. The overall architecture of the FAIR4Health solution.** FAIR4Health Platform was developed to apply Privacy-Preserving Distributed Data Mining (PPDDM) models enabling health research organizations to perform joint data mining operations without exposing any sensitive patient information to the outside world. PPDDM Agents, which are responsible for running the data mining algorithms on top of the FAIRified data for the use cases defined by the user through the FAIR4Health Platform, were developed for training, validation and testing of models for the use cases defined. To achieve its objectives, the PPDDM Agents communicate with the onFHIR.io FHIR Repository within the data source boundaries, and the FAIR4Health Platform to exchange the results and predictive model information in a distributed manner.



In this line, the pathfinder case study #2 was composed of two phases to reach the main objective. The first one included a retrospective multicenter observational study, including the training and generation of prediction models in the FAIR4Health Platform. In the second phase, a prospective observational study with a 30-day follow-up was performed, from April 2021 to September 2021, to evaluate the accuracy of this tool by collecting data from a selected sample of subjects. The study population consisted of individuals aged 18 and older with a diagnosis of COPD who were admitted to the hospital for this disease. Finally, to assess the prediction risk accuracy associated with the early prediction service for 30-days readmission risk in COPD patients, predictions generated by the FAIR4Health Platform were compared with real-world data. The clinical assessment concluded that from 100 recruited patients, the prediction was correct in 87% of cases (that is, in real-life, the patient was readmitted and the algorithm predicted that there was early 30-days hospital readmission risk; or the patient was not readmitted and the algorithm predicted that there was not early 30-days hospital readmission risk).

Further details of the FAIR4Health pathfinder case studies can be found in the public report on the demonstrators' performance<sup>10</sup>.

## Conclusions/Discussion

FAIR4Health partners achieved the project's objectives and the FAIR4Health use cases were successfully carried out through to the correct implementation of the technologies and performance of the complex FAIR4Health technical solution.

The main aim of the FAIR4Health project was to test the developed tools in the project: 1) application of FAIR principles in health research through the FAIR4Health FAIRification tools; 2) use of federated machine learning techniques; and 3) clinical, technical and functional validation of the FAIR4Health Platform and agents.

Therefore, FAIR4Health partners got positive conclusions from the FAIR4Health use cases. In both use cases, significant cross-cutting data-related issues and challenges were identified and addressed. The task to extract data from EHRs and other kinds of healthcare sources aligning this extraction with a FAIR4Health Common Data Model was not trivial and required a lot of conceptual and technical efforts, because: (i) complexity of the raw data (the source EHRs are commonly very complex including information in several tables in the source databases); (ii) free text used in some fields in the raw data sources; and (iii) differences between the type of the raw data sources. To address the complexity of the raw data, each health research organization from different countries that participated in data extraction involved colleagues who were experts in each source data model. To address the information in free text fields, Natural Language Processing (NLP) techniques were assessed, and finally, in some cases, manual NLP to extract structured information from unstructured information was performed. Due to the differences in the raw data sources, each raw dataset had to be analyzed in depth in collaboration between the clinical partners and the technical

partners. This involved determining the required configuration of the FAIR4Health solution to enable FAIRification of all raw data. Finally, coordinated federated machine learning models were created using all sources.

It is relevant to add other significant conclusions here, related to the application of the FAIR principles in health research:

- Implementation of FAIR principles allowed us to use larger and more heterogeneous datasets in FAIR4Health, increasing the variability of the data, the size of the datasets, and finally, more comprehensive and reliable results/outputs, compared to specific research studies without applying FAIR.
- We could reuse FAIR datasets from other clinical organizations in a secure way, ensuring compliance with General Data Protection Regulation (GDPR), and we could use the clinical datasets in the federated machine learning models. In the FAIR4Health project, we could also consider demographic, environmental, clinical and social information. We achieved greater variability of datasets and inclusion of more variables, compared to research where FAIR datasets are not reused.
- We obtained an increase in the scope of the research and improvements in health research, facilitating the discovery of scientific knowledge through data sharing and data reuse. Likewise, FAIR data reuse provided savings in data collection where much effort is currently invested.
- The implementation of FAIR principles facilitated the reproducibility of the study and access to large volumes of data to make the research more robust. We obtained the increase in secondary use of datasets once FAIR policies were implemented, related to the publication and sharing of FAIR datasets.

## Data and software availability

### Underlying data

No data are associated with this article.

### Extended data

Along with the FAIR4Health software, FAIR metadata related to the FAIRified datasets generated in the FAIRification process, is published in the FAIR4Health GitHub. This is available to the scientific community, and the FAIR4Health consortium continues assessing the possibilities to open publish these metadata in other public repositories. Further information: <https://github.com/fair4health/>.

Data are available under the terms of the [Creative Commons Zero "No rights reserved" data waiver](#) (CC0 1.0 Public domain dedication).

## Author contributions

Celia Alvarez-Romero: Investigation + Methodology + Validation + Visualization + Writing – Original Draft Preparation

Alicia Martínez-García: Investigation + Methodology + Validation + Visualization + Writing – Review & Editing

A. Anil Sinaci: Investigation + Software + Writing - Review & Editing

Mert Gencturk: Investigation + Software + Writing - Review & Editing

Eva Méndez: Conceptualization + Investigation + Writing - Review & Editing

Tony Hernández-Pérez: Conceptualization + Investigation + Writing - Review & Editing

Rosa Liperoti: Investigation + Validation + Writing - Review & Editing

Carmen Angioletti: Investigation + Validation + Writing - Review & Editing

Matthias Löbe: Investigation + Writing - Review & Editing

Nagarajan Ganapathy: Investigation + Writing - Review & Editing

Thomas M. Deserno: Investigation + Writing - Review & Editing

Marta Almada: Investigation + Validation + Writing - Review & Editing

Elisio Costa: Investigation + Validation + Writing - Review & Editing

Catherine Chronaki: Investigation + Writing - Review & Editing

Giorgio Cangioli: Investigation + Writing - Review & Editing

Ronald Cornet: Investigation + Writing - Review & Editing

Beatriz Poblador-Plou: Investigation + Validation + Writing - Review & Editing

Jonás Carmona-Pérez: Investigation + Validation + Writing - Review & Editing

Antonio Gimeno-Miguel: Investigation + Validation + Writing - Review & Editing

Antonio Poncel-Falcó: Investigation + Validation + Writing - Review & Editing

Alexandra Prados-Torres: Investigation + Validation + Writing - Review & Editing

Tomi Kovacevic: Investigation + Validation + Writing - Review & Editing

Bojan Zaric: Investigation + Validation + Writing - Review & Editing

Darijo Bokan: Investigation + Validation + Writing - Review & Editing

Sanja Hromis: Investigation + Validation + Writing - Review & Editing

Jelena Djekic Malbasa: Investigation + Validation + Writing - Review & Editing

Carlos Rapallo Fernández: Conceptualization + Investigation + Writing - Review & Editing

Teresa Velázquez Fernández: Conceptualization + Investigation + Writing - Review & Editing

Jessica Rochat: Investigation + Validation + Writing - Review & Editing

Christophe Gaudet-Blavignac: Investigation + Validation + Writing - Review & Editing

Christian Lovis: Investigation + Validation + Writing - Review & Editing

Patrick Weber: Investigation + Writing - Review & Editing

Miriam Quintero: Investigation + Software + Writing - Review & Editing

Manuel M. Perez-Perez: Investigation + Writing - Review & Editing

Kevin Ashley: Writing - Review & Editing

Laurence Horton: Writing - Review & Editing

Carlos Luis Parra Calderón: Conceptualization + Investigation + Project Administration + Supervision + Writing – Review & Editing

## References

- Wilkinson MD, Dumontier M, Aalbersberg IJ, et al.: **The FAIR Guiding Principles for scientific data management and stewardship.** *Sci Data.* 2016; **3**: 160018. [PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
- European Commission: **Cost of not having FAIR research data - Cost-Benefit analysis for FAIR research data.** 2018. [Reference Source](#)
- FAIR4Health Guidelines for implementing FAIR open data policy in health research.** [Reference Source](#)
- Sinaci AA, Núñez-Benjumea FJ, Gencturk M, et al.: **From Raw Data to FAIR Data: The FAIRification Workflow for Health Research.** *Methods Inf Med.* 2020; **59**(S 01): e21–e32. [PubMed Abstract](#) | [Publisher Full Text](#)
- FAIR4Health Data Curation Tool.** [Reference Source](#)
- FAIR4Health Data Privacy Tool.** [Reference Source](#)
- Han J, Pei J, Yin Y: **Mining frequent patterns without candidate generation.** *ACM sigmod record.* 2000; **29**(2): 1–12. [Publisher Full Text](#)
- Poblador-Plou B, Calderón-Larrañaga A, Marta-Moreno J, et al.: **Comorbidity of dementia: a cross-sectional study of primary care older patients.** *BMC Psychiatry.* 2014; **14**(1): 84. [PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
- Prados-Torres A, Calderón-Larrañaga A, Hancoco-Saavedra J, et al.: **Multimorbidity patterns: a systematic review.** *J Clin Epidemiol.* 2014; **67**(3): 254–266. [PubMed Abstract](#) | [Publisher Full Text](#)
- FAIR4Health Report on the demonstrators performance.** [Reference Source](#)



# Open Peer Review

Current Peer Review Status: ? ? ?

---

## Version 1

Reviewer Report 19 April 2022

<https://doi.org/10.21956/openreseurope.15486.r28796>

© 2022 Capella-Gutierrez S. This is an open access peer review report distributed under the terms of the [Creative Commons Attribution License](#), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.



**Salvador Capella-Gutierrez** 

Department of Life Sciences, Barcelona Supercomputing Center (BSC), Barcelona, Spain

Alvarez-Romero and colleagues provided a fascinating article on how FAIR data in health research can benefit researchers and patients. This effort is part of the recently finished FAIR4Health project and illustrates its progress.

I have to congratulate the authors on the extended FAIRification process for health research data. I appreciate the data versioning as an essential step toward reproducibility of any result.

I have some comments aiming to foster the discussion and clarify some aspects of this manuscript:

1. Is there any specific reason to focus on HL7 FHIR as an interoperability standard across participating centres? Despite HL7 FHIR, I'd expect at least a mention of other standards to facilitate health-related research, e.g. OMOP. This aspect is especially relevant when working with observational studies, e.g. cohort-based research. Thus, the interesting point is to know how extensible this work is and the technical implications of making such an effort.
2. Can you provide additional details on the privacy-preserving mechanisms? If it has already been published somewhere else, a couple of sentences summarising it with the reference should be enough.
3. I wonder about the use of onFHIR platform. I assume it has been developed by one of the FAIR4Health partners. Then it makes sense to use it. I appreciate the fact of using open-source software. Can similar tools/platforms replace it? How interoperable are the outputs of this platform? I'm thinking of interested parties willing to use/leverage FAIR4Health outcomes having their solutions.
4. I'd suggest delineating/introducing the use-cases earlier in the text.
5. I find figure #3 very informative. However, it led me to ask myself about potential

mechanisms for accessing data produced in the consortium. If there is any formal mechanism to access or request access to those datasets, I think they can be included in this figure.

6. Looking at use-case #1, you mentioned that implemented algorithms are available as part of the platform in the text. Perhaps you can include a link to it, e.g. a specific repository in the FAIR4Health GitHub Organization.
7. Still looking at use-case #1, you mentioned: "... a strong association between multimorbidity and polypharmacy and each of them with mortality." Have you performed any statistical analysis here? Perhaps it is good to include them as part of the use-case #1 discussion.
8. I like the description of use-case #2 and appreciate that there is an extensive report (68 pages) describing it. However, readers would appreciate a short explanation of the main findings for this use-case. Otherwise, it seems a bit disconnected.
9. You have briefly discussed the possibilities of using FAIR data for distributed ML across different sites. Can you explain the minimal computational requirements for carrying on these analyses? I think it is essential for readers to realise that it is not enough to have FAIR data at their sites but also the computational capabilities to conduct such analyses.
10. I like the description of the efforts to FAIRify health research data. However, I'm missing the language axis when using NLP technologies. NLP models and resources are language-dependent, which means that the final results are partially affected by them. As you are working with data from 5 different European countries, can you share your experiences on these aspects?

Looking forward to your comments on those aspects.

**Is the work clearly and accurately presented and does it cite the current literature?**

Yes

**Is the study design appropriate and does the work have academic merit?**

Yes

**Are sufficient details of methods and analysis provided to allow replication by others?**

Partly

**If applicable, is the statistical analysis and its interpretation appropriate?**

Partly

**Are all the source data underlying the results available to ensure full reproducibility?**

Yes

**Are the conclusions drawn adequately supported by the results?**

Yes

**Competing Interests:** No competing interests were disclosed.

**Reviewer Expertise:** FAIR principles

**I confirm that I have read this submission and believe that I have an appropriate level of expertise to confirm that it is of an acceptable scientific standard, however I have significant reservations, as outlined above.**

Author Response 23 May 2022

**Celia Alvarez**, Group of Research and Innovation in Biomedical Informatics, Biomedical Engineering and Health Economy, Institute of Biomedicine of Seville, IBiS / Virgen del Rocío University Hospital / CSIC / University of Seville, Seville, Spain

Dear reviewer, All comments and recommendations have been addressed. Improvements and expanded descriptions of some parts of the article have been included in the new version, along with contributions suggested by the other reviewers.

Thank you very much for the opportunity to improve the manuscript and for all the comments, very timely to improve the scientific quality. Best regards.

**Competing Interests:** No competing interests were disclosed.

Reviewer Report 08 April 2022

<https://doi.org/10.21956/openreseurope.15486.r28777>

© 2022 Saunders G. This is an open access peer review report distributed under the terms of the [Creative Commons Attribution License](#), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.



**Gary Saunders**

EATRIS ERIC, Amsterdam, The Netherlands

This is a strong brief report of two use cases utilising the FAIR4Health project solution in addressing two different use cases of health data demonstrating the benefits in the application of the FAIR principles to such data that result in strong conclusions. Furthermore, the authors have published and made available all software that is referenced in the report so that it is available to the community.

I support the publication of this report, and have only two minor corrective suggestions:

1. Please can the authors include the total number of individuals that were analysed in use case 1? It would also be good to have this number broken down to the 5 distinct populations, however this may not be possible due to privacy concerns. As the use case discusses the application of federated ML it would be good to have some indication of sample size that was used in the analysis.

2. The NLP techniques that are discussed in the conclusion should be described earlier in the manuscript. There is some nice text in the discussion section not only of the NLP techniques but also of the complex application of the FAIR4Health Common Data Model to the use cases discussed in the report. It would be good to have this properly described in the Methods section of the manuscript as not only will this aid transparency in the processes, and therefore reproducibility of the results described, but also aid the community in the potential application of the CDM to future datasets.

In summary, this is a nice, neat and concise manuscript describing the work in an easy-to-read, follow, and digest manner.

**Is the work clearly and accurately presented and does it cite the current literature?**

Yes

**Is the study design appropriate and does the work have academic merit?**

Yes

**Are sufficient details of methods and analysis provided to allow replication by others?**

Partly

**If applicable, is the statistical analysis and its interpretation appropriate?**

Partly

**Are all the source data underlying the results available to ensure full reproducibility?**

No source data required

**Are the conclusions drawn adequately supported by the results?**

Yes

**Competing Interests:** No competing interests were disclosed.

**Reviewer Expertise:** FAIR data, big data management, federated data access and analysis, health care data.

**I confirm that I have read this submission and believe that I have an appropriate level of expertise to confirm that it is of an acceptable scientific standard, however I have significant reservations, as outlined above.**

Author Response 23 May 2022

**Celia Alvarez**, Group of Research and Innovation in Biomedical Informatics, Biomedical Engineering and Health Economy, Institute of Biomedicine of Seville, IBiS / Virgen del Rocío University Hospital / CSIC / University of Seville, Seville, Spain

Dear reviewer, All comments and recommendations have been addressed. Improvements and expanded descriptions of some parts of the article have been included in the new version, along with contributions suggested by the other reviewers.

Thank you very much for the opportunity to improve the manuscript and for all the comments, very timely to improve the scientific quality. Best regards.

**Competing Interests:** No competing interests were disclosed.

Reviewer Report 04 April 2022

<https://doi.org/10.21956/openreseurope.15486.r28776>

© 2022 Gururaj A. This is an open access peer review report distributed under the terms of the [Creative Commons Attribution License](#), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited. The author(s) is/are employees of the US Government and therefore domestic copyright protection in USA does not apply to this work. The work may be protected under the copyright laws of other jurisdictions when used in those jurisdictions.



**Anupama Gururaj**

National Institute of Allergy and Infectious Diseases, National Institutes of Health, Rockville, MD, USA

This paper describes a project focusing on the need to make healthcare data FAIR and reusable in the European Union context. The study addresses an emerging research need within the healthcare system and will be of interest to clinical research personnel in various healthcare settings. This team has performed two pilot studies that outline a general framework that can be adopted by other systems. The manuscript is concise and provides an overview of the project and project outcomes including tools developed and use cases addressed.

#### Major Criticisms:

1. In the abstract, the authors mentioned access to certified FAIR datasets. It is not clear in the rest of the manuscript how and who has provided such certification to the metadata of the two datasets that are shared. Please provide further details or remove mention of certification.
2. In parts, the paper overstates some aspects of the methods and conclusions. As an example, in the methods section, it is stated that the FAIR4Health provides an overview of all considerations to overcome all barriers to open research data sharing by HRPOs. Likewise, in the conclusions/discussions section, the authors state that they have seen an increase in secondary use of the datasets without providing any usage metrics. My recommendation is to tone down the absolute statements.
3. It would be useful to the reader to get a better understanding of the Common Data Model that is used to harmonize the data from the various sources. Please provide a description of the model in the paper.
4. Other large-scale efforts such as OHDSI are community-led and have leveraged distributed



analytics for answering scientific questions. Please compare and contrast the FAIR4Health initiative to such efforts in the discussion.

5. It would be useful to the community to hear about lessons learnt and challenges that have not yet been solved through the life of the project. Also, please comment on the scalability of the approach since a lot of manual effort and coordination was a part of the project.

Minor comments:

1. Some of the sentences (such as the starting sentence of the second paragraph of the introduction (The FAIR guiding principles...)) are very long and make the paper hard to read. It is recommended to simplify or split these sentences.
2. Figure 1 is hard to read and follow since there are overlapping boxes. Please fix the figure to increase clarity.

**Is the work clearly and accurately presented and does it cite the current literature?**

Partly

**Is the study design appropriate and does the work have academic merit?**

Yes

**Are sufficient details of methods and analysis provided to allow replication by others?**

Yes

**If applicable, is the statistical analysis and its interpretation appropriate?**

Not applicable

**Are all the source data underlying the results available to ensure full reproducibility?**

Yes

**Are the conclusions drawn adequately supported by the results?**

Partly

**Competing Interests:** No competing interests were disclosed.

**Reviewer Expertise:** Biomedical data sharing and management, clinical informatics

**I confirm that I have read this submission and believe that I have an appropriate level of expertise to confirm that it is of an acceptable scientific standard, however I have significant reservations, as outlined above.**

Author Response 23 May 2022

**Celia Alvarez**, Group of Research and Innovation in Biomedical Informatics, Biomedical Engineering and Health Economy, Institute of Biomedicine of Seville, IBiS / Virgen del Rocío University Hospital / CSIC / University of Seville, Seville, Spain

Dear reviewer, All comments and recommendations have been addressed. Improvements and expanded descriptions of some parts of the article have been included in the new version, along with contributions suggested by the other reviewers.

Thank you very much for the opportunity to improve the manuscript and for all the comments, very timely to improve the scientific quality. Best regards.

**Competing Interests:** No competing interests were disclosed.

---