

# BETTER: Better rEal-world healTh-DaTa distributEd analytics Research platform<sup>\*</sup>

Ana León Palacio<sup>1,\*†</sup>, José F. Reyes Román<sup>1,†</sup> and Oscar Pastor<sup>3,†</sup>

<sup>1</sup> Universitat Politècnica de València, Camí de Vera s/n 46022 Valencia, Spain

## Abstract

Over the last few years, data-driven medicine has gained increasing importance in terms of diagnosis, treatment, and research due to the exponential growth of healthcare data. The linkage of cross-border health data from various sources, including genomics, and analysis via innovative Artificial Intelligence (AI) approaches will allow a better understanding of risk factors, causes, and the development of optimal treatment in different disease areas. However, the reuse of patient data is often limited to data sets available in a single medical center. The main reasons why health data are not shared across institutional boundaries rely on ethical, legal, and privacy aspects and rules. Therefore, in order to (1) enable the sharing of health data across national borders, (2) fully comply with the current GDPR privacy guidelines/regulations, and (3) innovate by pushing research beyond state of the art, BETTER proposes a robust decentralized privacy preservation infrastructure which will empower researchers, innovators, and healthcare professionals to exploit the full potential of larger sets of multisource health data through tailored AI tools useful to compare, integrate, and analyze in a secure, cost-effective fashion; with the end goal of supporting the improvement of citizen health outcomes. In detail, this interdisciplinary project proposes the co-creation of three clinical use cases involving seven medical centers located in the EU and beyond, where sensitive patient data, including genomics, are made available and analyzed in a GDPR-compliant mechanism via a Distributed Analytics (DA) paradigm called the Personal Health Train (PHT). The main principle of the PHT is that the analytical task is brought to the data provider (medical center), and the data instances remain in their original location. In this project, two mature implementations of the PHT (PADME and Vantage6), already validated in real-world scenarios, will be fused to build the BETTER platform.

## Keywords

Health Data, Personal Health Train, Artificial Intelligence, Distributed Analytics, PADME, Vantage6

## 1. Introduction

### 1.1. Context and Motivation

Integrating vast arrays of health data from genomics and electronic health records through advanced artificial intelligence (AI) technologies has revolutionized our understanding of diseases, risk factors, and therapeutic strategies [1]. However, the utility of data in medical research is profoundly dependent on its volume and diversity. This is especially true when studying rare diseases and could be solved by sharing data among clinical centers. Nevertheless, sharing patient data across institutions is conditioned by ethical, legal, and privacy concerns [2].

Current data protection regulations, such as the General Data Protection Regulation (GDPR), prohibit data centralization for analysis because of privacy risks, such as the accidental disclosure of personal data to third parties. Overcoming these challenges requires moving from a centralized to a decentralized paradigm that enables the secure and efficient exchange of health information across borders while ensuring compliance with privacy regulations.

<sup>\*</sup> RCIS 2025 Workshops and Research Projects Track. 20 - 23 May, 2025. Seville, Spain.

<sup>1\*</sup> Corresponding author.

<sup>†</sup> These authors contributed equally.

✉ aleon@vrain.upv.es (A. León); jreyes@vrain.upv.es (J. F. Reyes Román); opastor@dsic.upv.es (O. Pastor)

id 0000-0003-3516-8893 (A. León); 0000-0002-9598-1301 (J. F. Reyes Román); 0000-0002-1320-8471 (O. Pastor)

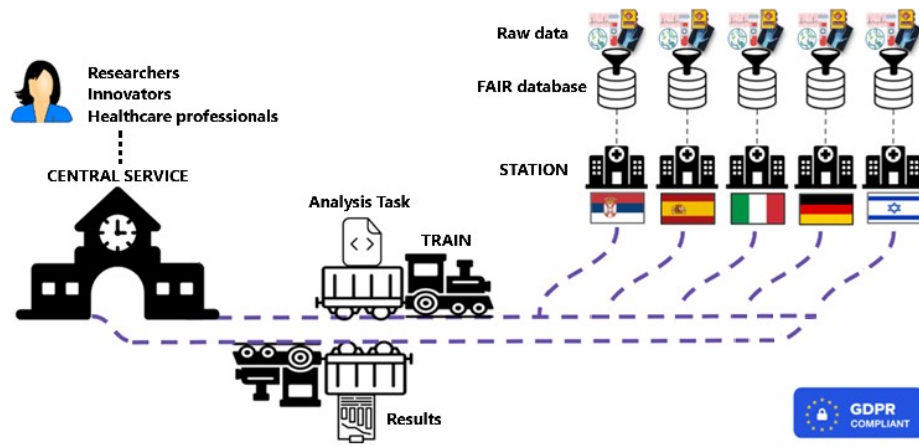


© 2025 Copyright for this paper by its authors. Use permitted under Creative Commons License Attribution 4.0 International (CC BY 4.0).

This shift entails complex technical challenges, including orchestrating decentralized computation, ensuring secure audibility, and harmonizing heterogeneous clinical and genomic datasets. Addressing these challenges requires a coordinated, long-term effort involving the development, deployment, and validation of federated learning infrastructures that can operate in real-world clinical settings.

## 1.2. Current Approach

The BETTER project<sup>2</sup> proposes a decentralized infrastructure that uses Distributed Analytics (DA) through a mechanism known as the Personal Health Train (PHT) [3]. The PHT model ensures that analytical processes are executed at the data provider's site, allowing data to remain securely within its original location. This model can be illustrated using a railway system analogy, where the key infrastructure components include Trains, Stations, and a Central Service (Fig. 1).



**Figure 1:** Main infrastructure of the PHT model adapted to the BETTER project.

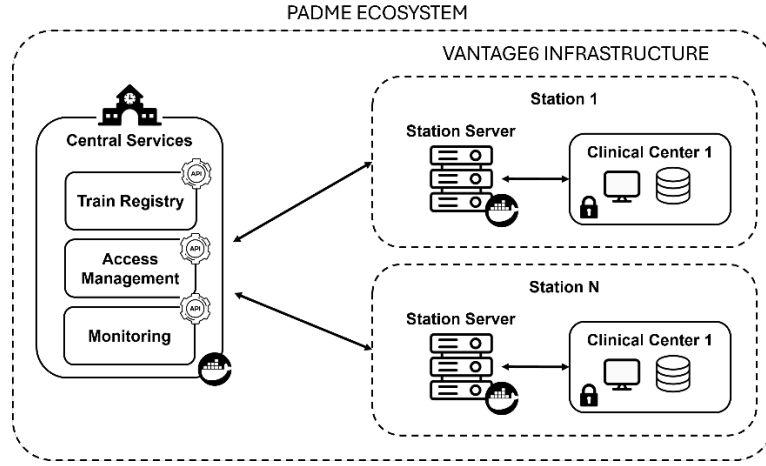
Trains encapsulate code to execute analytical tasks at distributed data nodes, known as Stations. As they travel from one Station to another, they process data locally, leveraging the available information at each stop to incrementally build the final analysis result. A Station is a node (institution, hospital, or department) within the distributed architecture that securely stores confidential data and executes Train operations. Each Station functions as an independent and autonomous unit, managing permission requests to control access to its confidential data. The Central Service includes procedures for Train orchestration, operational logic, business logic, data management, and discovery. The Central Service offers: (1) a metadata repository for efficient data discovery; (2) management tools for Train creation, secure transmission to Stations, orchestration, monitoring, and debugging; and (3) a repository of pre-trained Trains that healthcare professionals can directly apply to their data, enabling them to obtain results from an AI-based method trained iteratively on data from multiple institutions.

## 1.3. The Technology Behind the BETTER Project

The BETTER project integrates two established implementations of the PHT: PADME (Platform for Analytics and Distributed Machine Learning for Enterprises)<sup>3</sup> and Vantage6 (priVAcY preserviNg federaTed leArninG infrastruCTurE for Secure Insight eXchange) [4]. Fig. 2 shows how PADME and VANTAGE6 are connected to create the BETTER infrastructure.

<sup>2</sup> <https://www.better-health-project.eu/>

<sup>3</sup> <https://padme-analytics.de/>



**Figure 2:** BETTER infrastructure.

Both implementations have already demonstrated their effectiveness in various clinical settings, including oncology [5], diabetes [6], and cardiovascular diseases [7]. These works show that federated learning in the healthcare domain is technically feasible.

To ensure the ethical and trustworthy use of artificial intelligence, the BETTER project adopts a comprehensive strategy to address AI bias throughout the development lifecycle. All AI tools are developed following the European Commission’s Trustworthy AI Guidelines, with specific attention to fairness, transparency, and robustness. Co-creation with clinicians ensures that domain knowledge is integrated into the design process, helping to identify and correct biases. Before model training and testing, clinical datasets undergo rigorous curation, including pseudonymization, semantic annotation with standardized ontologies, and data quality checks for completeness, consistency, and integrity. This process is complemented by the FAIRification of metadata to ensure interoperability and reproducibility. Together, these practices establish a robust foundation for the development of clinically reliable AI tools.

#### 1.4. Participants

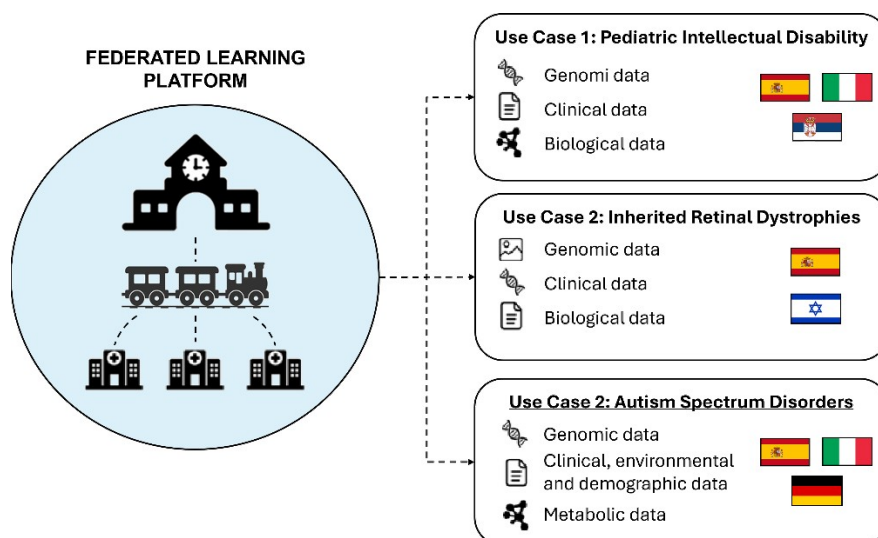
BETTER is a Horizon Europe project with a duration of 42 months and the participation of 14 organizations from eight countries, including seven clinical institutions and seven technological centers:

- *Clinical Institutions:* Klinikum Der Universitaet Zu Koeln (UKK - Germany), Fundació de Recerca Sant Joan de Déu (FDSJD - Spain), Azienda Socio-Sanitaria Territoriale Fatebenefratelli Sacco (BUZZI - Italy), Fundació Docència i Recerca Mutua de Terrassa (Spain), Instituto de Investigación Sanitaria - Hospital Universitario y Politécnico La Fe (Spain), Institut Za Molekularnu Genetiku I Geneticko Inzenjerstvo (IMGGE - Serbia), and Hadassah Medical Organization (HMO - Israel).
- *Technological Centers:* Datrix Spa (Italy), Universiteit Maastricht (UM - Netherlands), Politecnico di Milano (POLIMI - Italy), Universitat Politècnica de València (UPV - Spain), Universitetet i Tromsø - Norges Arktiske Universitet (UiT - Norway), Rheasoft ApS (Denmark), and Noosware Bv (Netherlands).

#### 1.5. Clinical Use Cases

The project aims to apply these innovative DA methodologies to three clinical use cases: Pediatric Intellectual Disability, Inherited Retinal Dystrophies, and Autism Spectrum Disorders. The overarching goal is to harness the full potential of multisource health data, enabling researchers,

healthcare professionals, and innovators to conduct comprehensive analyses, integrate disparate data types, and derive meaningful insights securely and cost-effectively.



**Figure 3:** Type of data to be analyzed on each use case.

As shown in Fig. 3, seven medical centers will integrate, validate, and utilize the digital tools built on top of the BETTER platform, where multiple data sources from different centers can be fused and exploited to improve clinical outcomes.

### 1.5.1. Use Case 1: Integration of Genomic and Phenotypic Data from Pediatric Rare Diseases to Decipher Pathways of Intellectual Disability

Intellectual disability (ID) is a common disorder characterized by significant limitations of cognitive functions and adaptive behavior, with onset before age 18. This use case aims to (1) evaluate and correlate the phenotypic, genomic, multi-omic, and clinical parameters between early-diagnosed and later-diagnosed patients; (2) Improve diagnosis by identifying new genetic biomarkers that can be used in newborn screening protocols; (3) Develop new tools based on Digital Twins Model to define new diagnostic biomarkers, pathways, and therapeutic molecular targets.

To such an aim, clinical data, brain images, genomic data (whole-exome, whole-genome sequences), and biological data (cellular and molecular pathways) will be integrated. The participants in this use case are the Hospital Sant Joan de Deu (medical leader), IMGGE, the Children's Hospital Vittore Buzzi, and the Politecnico di Milano (technological leader).

### 1.5.2. Use Case 2: Accelerate Inherited Retinal Dystrophies Diagnosis using AI

Inherited Retinal Diseases (IRDs) are a group of disorders characterized by the generally progressive death or dysfunction of photoreceptors and retinal pigment epithelium (RPE) cells, leading to loss of visual function, sometimes leading to legal blindness. An early molecular diagnosis is necessary to confirm the clinical diagnosis and offer adequate care to patients. In addition, developing new genetic analysis tools that allow the precise identification of the molecular cause of disease is essential to improve the understanding of the pathophysiological mechanisms at the base of the symptoms and open the doors to future therapies. This study aims to (1) identify pathogenic genes and variants responsible for the IRDs, and (2) define existing genotype-phenotype correlations to better understand the prognosis of patients and improve their clinical management.

To such an aim, genomic data (gene panels, clinical exome, whole exome, whole genome), clinical reports, and images will be integrated. The participants in this Use Case are the Hospital

Universitario La Fe de Valencia (medical leader), the Hadassah Medical Center, and the Polytechnic University of Valencia (technological leader).

### **1.5.3. Use Case 3: Predicting the Risk of Self-Harm and Suicidal Behaviors in Patients with Autism Spectrum Disorders**

Autism Spectrum Disorders (ASD) are neurodevelopmental disabilities characterized by social, communication, and behavioral challenges. Children and adolescents with ASD are at a substantially higher risk of self-injurious and suicidal behavior compared to the general population (up to 9 times). However, the causes of this increased risk remain largely unknown, and there is little knowledge about the potential role of phenotypic, metabolic, genomic, and environmental factors. This use case aims to (1) identify predictive phenotypic, genomic, and environmental risk factors of suicidality and self-injury in ASD individuals; (2) Personalize prevention intervention plans to reduce self-injury and suicidality in each ASD individual; and (3) Develop monitoring strategies to recognize signs of vulnerability in ASD individuals that will lead to prevent strategies at an earlier stage and thereby further reduce risk of self-harm and suicidal intentions.

To such an aim, clinical, metabolic, environmental and demographic data, patient interviews, and genomic data (epigenome and whole genome sequencing) will be integrated. The participants in this Use Case are Hospital Universitario Mutua Terrassa (medical leader), Children's Hospital Vittore Buzzi, and the Klinikum Der Universitaet Zu Koeln (UKK). UKK also participates as the technological leader.

## **2. Project Objectives**

The BETTER project has five main objectives (1) Overcome cross-border barriers to health data integration, access, FAIRification, and preprocessing; (2) Ensure health data fusion and integration; (3) Deploy a distributed analytics framework for cross-border data processing and analysis; (4) Development of distributed tools leveraging artificial intelligence capabilities; (5) Include ethical, legal and societal aspects (ELSA) in the AI lifecycle.

### **2.1. Objective 1: Overcome Cross-Border Barriers to Health Data Integration, Access, FAIRification, and Preprocessing**

The main aim of this first objective is to guide medical centers in collecting patients' data following a common schema to promote interoperability and the reuse of datasets in scope. This includes collecting legal, ethical, and data protection authorizations and using well-established and widely understood ontologies. Data pseudonymization will be performed as a default preprocessing step to mitigate the risk of personal data leaks. Finally, a BETTER station will be installed at each medical center, ensuring access to the relevant local datasets.

### **2.2. Objective 2: Ensure Health Data Fusion and Integration**

To gain the maximum from data, one of the important steps is integrating multiple data sources to produce more consistent, accurate, and useful information than any single data source. The ambition is to fuse several dimensions, including laboratory analysis, medical reports, drug therapy, imaging, genomics, socio-demographic, geographical, and medical questionnaires.

BETTER uses standardized ontologies (e.g., NCIT, LOINC, ICD-11) and a shared metadata schema to ensure semantic alignment of data across sites. This enables data fusion and integration using the proposed distributed framework in two directions: within a single medical center (local data fusion) and across centers (distributed data fusion) by leveraging each other's historical datasets. Local data fusion involves integrating data from multiple sources within a single institution. This type of data fusion is useful when the data sources are heterogeneous, such as genomic, clinical, and phenotypic data of the same patient. Distributed data fusion integrates data from multiple institutions, a fairly

novel discipline that includes removing potential biases due to different collection protocols or techniques.

### **2.3. Objective 3: Deploy Distributed Analytics Framework for Cross-Border Data Processing and Analysis**

The ambition of this objective regards the deployment of BETTER, a privacy-by-design infrastructure, to all medical centers connecting FAIR data sources and allowing federated data analysis and machine learning. To effectively exploit multiple datasets via AI, a common schema and ontology should be applied.

### **2.4. Objective 4: Development of Distributed Tools Leveraging Artificial Intelligence Capabilities**

To properly answer clinical needs and push data analysis boundaries beyond state-of-the-art, tailored tools must be developed to exploit DA and AI within each use case. The tools will be developed using a co-creation methodology where medical end-users closely collaborate with researchers and technology providers, enabling the development of new concepts.

### **2.5. Objective 5: Include Ethical, Legal and Societal Aspects (ELSA) in the AI Lifecycle**

Most data science projects do not co-create or co-develop using a methodology that includes the ethical, legal, and societal aspects (ELSA) involved in the data science lifecycle. In this objective, the BETTER project will develop ELSA-awareness tools and methods for co-creating and co-developing AI models and apply them to the proposed use cases. This will ensure the appropriateness and clinical effectiveness of the developed AI tools while considering the safety, value, and sustainability of the AI.

## **3. Impact and Expected Outcomes**

Overcoming the current barriers to data sharing and utilization, the BETTER project opens the way to more accurate diagnoses, tailored treatments, and a deeper understanding of complex diseases. The project's target groups are healthcare professionals, researchers, innovators, health policymakers, and citizens.

BETTER promotes a hands-on and experience-building approach towards the implementation of cross-border data-sharing partnerships in the area of real-world health data. This contribution paves the way for a European medical center data sharing and analysis network. The expected outcomes for the BETTER project are:

- A public release of the platform implementation, which will reinforce two open-source projects, namely PADME and Vantage6.
- Publication of the FAIRification pipelines, data catalogs, and ontologies to unleash the potential of data exchange and reuse.
- Release of synthetic datasets to the community using generative AI techniques. This is particularly valuable for developing and benchmarking models in data-restricted scenarios.
- Finally, cross-border health data secure exchange and reuse require a solid and compliant legal, ethical, and data protection framework. By enhancing existing templates, BETTER will consolidate and publish a documentation folder useful and applicable for future initiatives.

## 4. Current Project Status and Future Work

The BETTER project is currently in its 14th month of work, having accomplished the following milestones:

- The ethical approvals and other documents required to execute the use cases have already been approved.
- All clinical and technological partners have already installed the required technology, and tests are already underway to ensure the correct connection to the platform.
- The FAIR database and the ETL process for each medical center are already finished.

Currently, the medical and technological leaders of each use are working on the definition of the analytical tasks, based on the datasets and their specific goals. Once the analytical tasks are defined, the technological partners can start designing the AI algorithms (trains) according to the federated learning paradigm to obtain the results. The infrastructure is being tested to ensure that all nodes (technological and medical) are correctly connected and trains can move between stations and execute simple tasks. The different interfaces to get statistics about the data are also under development. In addition, since no real data is currently available, some synthetic datasets are being generated to perform the tests without compromising security and privacy.

## 5. Conclusion

The BETTER project aims to provide a decentralized paradigm to enable the secure and efficient exchange of health information across borders while ensuring compliance with privacy regulations. Using the Personal Health Train model implemented by two open implementations (PADME and Vantage6), the project will explore the feasibility of federated learning in three clinical use cases (Pediatric Intellectual Disability, Inherited Retinal Dystrophies, and Autism Spectrum Disorders). BETTER is a 42-month project financed by the European Union's Horizon Europe research and innovation program, with 14 participants from eight countries, including seven clinical institutions and seven technological centers. The overarching goal is to harness the full potential of multisource health data, enabling researchers, healthcare professionals, and innovators to conduct comprehensive analyses, integrate disparate data types, and derive meaningful insights securely and cost-effectively.

## Acknowledgements

This work is part of the Horizon Europe project BETTER. The BETTER project has received funding from the European Union's Horizon Europe research and innovation program under grant agreement No. 101136262. <https://www.better-health-project.eu/>.

## References

- [1] S. Welten, L. Neumann, Y. U. Yediel, L. O. B. da Silva Santos, S. Decker, O. Beyan, Dams, A distributed analytics metadata schema, *Data Intelligence* 3 (2021) 528–547.
- [2] T. M. Deist, F. J. Dankers, P. Ojha, M. S. Marshall, T. Janssen, C. Faivre-Finn, C. Masciocchi, V. Valentini, J. Wang, J. Chen, et al., Distributed learning on 20 000+ lung cancer patients—the personal health train, *Radiotherapy and Oncology* 144 (2020) 189–200.
- [3] O. Beyan, A. Choudhury, J. van Soest, O. Kohlbacher, L. Zimmermann, H. Stenzhorn, M. R. Karim, M. Dumontier, S. Decker, L. O. B. da Silva Santos, A. Dekker, Distributed Analytics on Sensitive Medical Data: The Personal Health Train, *Data Intelligence* 2 (2020) 96–107. doi:10.1162/dint\_a\_00032.

- [4] A. Moncada-Torres, F. Martin, M. Sieswerda, J. Van Soest, G. Geleijnse, Vantage6: an open source privacy preserving federated learning infrastructure for secure insight exchange, in: AMIA annual symposium proceedings, volume 2020, 2021, p. 870.
- [5] S. Theophanous, P.-I. Lønne, A. Choudhury, M. Berbee, A. Dekker, K. Dennis, A. Dewdney, M. A. Gambacorta, A. Gilbert, M. G. Guren, et al., Development and validation of prognostic models for anal cancer outcomes using distributed learning: protocol for the international multi-centre atomcat2 study, *Diagnostic and prognostic research* 6 (2022) 14.
- [6] C. Sun, J. van Soest, A. Koster, S. J. Eussen, M. T. Schram, C. D. Stehouwer, P. C. Dagnelie, M. Dumon-tier, Studying the association of diabetes and healthcare cost on distributed data from the Maastricht study and statistics netherlands using a privacy-preserving federated learning infrastructure, *Journal of Biomedical Informatics* 134 (2022) 104194.
- [7] B. Scheenstra, A. Bruninx, F. van Daalen, N. Stahl, E. Latuapon, M. Imkamp, L. Ippel, S. Duijsings-Mahangi, D. Smits, D. Townend, et al., Digital health solutions to reduce the burden of atherosclerotic cardiovascular disease proposed by the carrier consortium, *JMIR cardio* 6 (2022) e37437.

## **Declaration on Generative AI**

The author(s) have not employed any Generative AI tools.