INVENIO RDM

Powering open science and collaboration with Invenio

Northwestern University Invenio Team 03 March 2020



@inveniosoftware

OHDSI: open, collaborative science

DRMATICS	VALUES	INNOVATION	Observational research is a field which will benefit greatly from disruptive thinking. We actively seek and encourage fresh methodological approaches in our work.
		REPRODUCIBILITY	Accurate, reproducible, and well-calibrated evidence is necessary for health improvement.
EALTH DATA SC		COMMUNITY	Everyone is welcome to actively participate in OHDSI, whether you are a patient, a health professional, a researcher, or someone who simply believes in our cause.
		COLLABORATION	We work collectively to prioritize and address the real world needs of our community's participants.
		OPENNESS	We strive to make all our community's proceeds open and publicly accessible, including the methods, tools and the evidence that we generate.
		BENEFICENCE	We seek to protect the rights of individuals and organizations within our community at all times.

MISSION: To improve health by empowering a community to collaboratively generate the evidence that promotes better health decisions and better care.



https://www.ohdsi.org/who-we-are/mission-vision-values/

List of OHDSI Working Groups

OHDSI has a variety of ongoing projects lead by Working Group teams. We would be delighted to have your participation. Please contact the team lead to join.

Working group leaders looking for assistance on uploading meeting recordings can find help here:

ma How to upload meeting recordings on the OHDSI wiki

A list of upcoming working group meetings is available here: Working Group Meeting Schedule

<u>v6</u>

CDM

AOMC

Person

Active Workgroups:

- OHDSI Community
- Algorithmic Phenotyping
- Achilles WG
- Architecture WG
- Atlas & WebAPI WG
- CDM and Vocabulary Development WG
- Causal Inference
- Cerner to OMOP
- Chart Review Question Interface Project
- China WG
- Clinical Trials WG
- Devices WG
- Dissemination WG
- FHIR WG
- Genomics WG
- GIS WG
- Gold-Standard Phenotype Library
- Hadoop WG
- Latin OHDSI
- NLP WG
- Maternal & Child Health
- Metadata and Annotations WG
- OHDSI Steering Working Group
- Oncology WG
- PGHD WG
- Patient-Level Prediction WG
- Pharmacovigilance evidence investiga
- Population-Level Estimation WG
- Psychiatry WG
- Quality Measures WG
 THEMIS
- = Transfusion WG
- The Book of OHDSI WG
- Women of OHDSI
- Data Quality



Observation period

Condition occurrence

Drug exposure

Procedure occurrence

Device exposure

Measurement

Survey_conduct

Observation

Specimen

Fact_relationship

Visit detail

Note NLP

tandardized healt

system data

Location

Location history

Care_site

Provider

Standardized derive

elements

Condition_era

Drug_era

Dose era

Cohort Cohort definit

Standardized healt

economics

Cost

Paver plan perio

Standardized metadata

CDM_source

Metadata

Standardized

vocabularies

Concept

Vocabulary

Domain

Concept_class

Concept_relationship

Relationship

Concept synonym

Concept_ancestor

Source_to_concept_map

Drug_strength





Our Community & Data-Sharing Network

- > 2,500 distinct users across six continents have posted to our community forum
- Our community has a distinct range of both stakeholders and disciplines
- > 100 different databases
- > half a billion patient records
- Data from 19 different countries, with > 200 million patient records from outside the U.S.

OHDSI is a multi-stakeholder, interdisciplinary collaborative to bring out the value of health data through large-scale analytics. All our solutions are open-source.

OHDSI has established an international network of researchers and observational health databases with a central coordinating center housed at Columbia University.



https://www.ohdsi.org/

Benefits of <u>opening</u> science...



Greater access to scientific inputs and outputs can increase scientific productivity through reducing duplication, allowing **more research from the same data** and multiplying opportunities for domestic and global participation in the research process.

Open science can **reduce delays in the re-use of scientific research** including articles and data, and promote a swifter path from research to innovation to produce new products and services.

Science, often publicly funded, should be publicly accessible **to promote a greater awareness** among citizens and to build public trust and support for public policies and investments in research. Open science also promotes citizen science in experiments and data collection.

PRDM

https://upload.wikimedia.org/wikipedia/commons/5/5a/UCT_RDM_Why-Open-Science.png

Open access to scientific outputs allow for **greater evaluation and scrutiny** by the scientific community which means more accurate replication and validation of research results. Openness to data contributes to maintain science's self-correction principle.

Science plays a key role in today's knowledge economies

and increased access to research results, including data, can positive impact not only scientific systems but also innovation.

Open science promotes collaborative efforts and faster knowledge transfer for a better understanding of global challenges and wicked problems.



Invenio software powers open science



Open Source framework for large-scale digital repositories

Turn-key Research Data Management repository

Integrated Library System





Free Open Source Software

Invenio is Free Open Source Software supported by a committed community of multidisciplinary institutions.

Code | Docs | Examples



Friendly and Responsive Community

Although Invenio was born at CERN, its community is growing bigger every day. Talk to the team now in our chatroom or forum.

Chatroom | Forum | Get Involved | Events



The "Safe bet"

Invenio community has been around for 20 years. **Solid services** have been built on top of it to ensure long-term confidence.

Live services | Products



How did this collaboration start (and what about Zenodo?!)

What motivated the InvenioRDM project?

- Some organizations tried to reuse the existing open source Zenodo source code
- Other orgs tried to use the Invenio Framework to build a RDM repository from scratch
- Several orgs tried to make the same modifications but had no easy way of sharing their changes
 All these groups came together to create a collaborative open source project and grow a sustainable community.



Zenodo will also run on InvenioRDM by the end of the project period.



We're leveraging Invenio as a strong foundation. Here's why.

- **Research, shared.** Securely share and preserve data records and a wide range of research types with collaborators. Allows easy dissemination to the community.
- Discoverable. Leverages metadata standards and the powerful Elasticsearch full-text search engine retrieves, facets, sorts, and filters your searches with ease.
- Scalable. Invenio is fast. Designed to manage 100+ million records and petabytes of files. All data can be archived independently of the size.
- **Communities.** Create and curate your own community (e.g., workshop, project, lab, or journal).
- A robust community: Large team of developers & active open source community. A SAAS-model for service via TIND (CERN spinoff). Invenio is widely used by <u>many organizations</u> & underlying technology (Python, Flask) widely supported.
- Next-Generation: With InvenioRDM, any organization can launch a turn-key open source next-generation repository
 platform with world-class features to support open and FAIR science. http://ngr.coar-repositories.org/
- Get credit & be cited. Get a DOI to make records easily and uniquely citable. Pre-formatted citation text makes it easy
 to cite your work and be cited. Contributor roles allow you to recognize the whole team.
- **Metrics.** Industry standard usage statistics for record pages with all tracking completely anonymized.
- **FAIR.** Advanced features to make your research Findable, Accessible, Interoperable, & Reusable.
- **Compliance-friendly.** Comply with data sharing mandates* and acknowledge your funders.
- **Easy.** Turn-key research data management platform & index can be easily deployed in the local environment by your team or by a service provider, such as TIND. Customize the look and feel to your local environment.



RDM platforms are critical to help preserve and share research, enable reproducibility, and empower reuse of datasets, protocols, engagement or study materials, & a wide range of other research products.

The InvenioRDM project has two goals:



Repository Platform

Build a turn-key research data management (RDM) repository platform based on **Invenio Framework** and **Zenodo**.



Community

Grow a community of research institutions, private companies and individuals to sustain the platform going forward.



The platform

A few highlights...





https://thenounproject.com/

InvenioRDM stack







PostgreSQL and **MySQL** are powerful relational databases with JSONsupport as well as a strong reputation for reliability, robustness, and performance.



Invenio is built using **Python 3**, the **Flask** micro web framework and a suite of the best community-built Python libraries.



InvenioILS UIs are built using React, the well-known JavaScript library.

Invenio is JSON-native and provides RESTful APIs to make it easy to build apps on top of the framework



InvenioRDM roadmap

February

- Milestones: Draft governance and sustainability plan, mock-up feedback from collaborators
- Release
 - Branding customization institutional theming can be applied
 - First iteration of the Data model
 - Improved CLI with improved workflow
 - New documentation site for developers (<u>http://inveniordm.docs.cern.ch</u>)
 - Closer project tracking with enhanced structure and outreach

March

- Milestones: First release for core plugins, review of February release
- Release
 - Search permissions
 - Deposit page
 - Improved record page
 - Data model update
- To see further ahead: <u>https://invenio-software.org/products/rdm/roadmap/</u>



Standing up InvenioRDM

INVENIORDM User	Documentation	Q Search	docs-invenio-rdm
Home Install Preview	Develop Extensions Deploy		
User Documentation Home	Home Welcome to the InvenioRDM project! InvenioRDM is a ready-to-go turn-key Research Da the Invenio framework so you don't have to. This documentation walks you through installing, co InvenioRDM instance. Furthermore, you can custor extensions! Follow the below sections in order to go straight to what you're interested in - each section in	ta Management repository. It puts together onfiguring and deploying your very own nize it for your need and even add your own et a full picture of using InvenioRDM or go is independent.	Table of contents Install Preview Develop Extensions Deploy
	Install ¶ To get started with InvenioRDM, you will want to ins creating and updating your instance. This tool in tur manage it. > Installation Guide	stall invenio-cli , our command line tool for n allows you to easily install an instance and	

1- Install invenio-cli **pip install invenio-cli**

2- Initialize your project invenio-cli init --flavour=RDM

3- Run it cd <project name> invenio-cli containerize

4- Visit https://localhost firefox https://localhost



System requirements

Invenio can run in Docker, on virtual machines, or on physical machines. Invenio can run on a single machine or a cluster of 100s of machines.

It all depends on exactly how much data you are handling and your performance requirements.

Small installation:

- Web/app/background servers and Redis: 1 node
- Database: 1 node
- Elasticsearch: 1 node

Medium installation:

- Load balancer: 1 node
- Web/app servers and background workers: 2 nodes
- Database: 1 node
- Elasticsearch: 3 nodes
- Redis/RabbitMQ: 1 node

Large installation:

- Load balancer: 2 node (with DNS load balancing)
- Web/app servers: 3+ nodes
- Background workers: 3+ nodes
- Database: 2 nodes (master/slave)
- Elasticsearch: 5 nodes (3 data, 2 clients)
- Redis: 3 nodes (HA setup)
- RabbitMQ: 2 nodes (HA setup)



Search and retrieve datasets using standards-based documentation Subjects optional*

Robust search enhanced by:

- Standardized forms of name • (LDAP + ORCiD coming soon)
- Standard subject terms • (MeSH, Library of Congress Subject terms)
- Standardized citation • formats
- Clear levels of access •
- Standard application of • licenses

InvenioRDM @ No ® Catalog your Research	orthwestern University	Medical	Select Medical Subject Heading (MeSH) terms Medical Subject Heading (MeSH 2) terms.				
Search Authors diabetes		Topical	Select Library of Congress (FAST) topical terms				
Bacarelli, Andrea (1) Dyer, Alan R. (1) Hou, Lilang (1) Hou, Cang (1) Liu, Lei (1) Lowe, Lynn P. (1) Metager, Boyd (1) Zhang, Wei (1)	Found 2 results. Hyperglycemia and Advert Metager, Royd & Lowe, Lynn P. & Dyer, Alan I HAPO aim to asses the association between meta gestational diabetes mellitus (GDM). 25,000 prog complete and participated in the study. Women's outcomes were analyzed for any associations or i Diabetes Mellitus (GDM) based of off the data co	rse Pregnancy Stu II. In a state of the set of the set of the set of the set of the set of the set of the set of the set of the set of the set	dy (HAPO) dataset	+ Add another FAST term			
License Onher (Not Open) (2) Resource Type Dutinet (2) -	Methylomics of Prenatal dataset Mov, Lilang & Hu, Gang & Bacarelli, Andrea Building off of The Tianjin GDM Observational St to make a longitudinal study. Blood samples from phenotypes. This analysis investigates the effects children as a result of in store exposure. This proj effected.	Gestational Diabe	etes Mellitus (GDM)	ntion n and hildren			

 \times



Data management for reproducibility and Open Access: study-focused resource types



InvenioRDM helps you store, manage and, if needed, share your study's outputs:

- Study-based resource types to manage a large range of assets
- **Reproducibility** is enhanced: store research proposals, datasets, code
- Be **compliant** with data sharing mandates
- **Cite** and **attribute** the work of all contributors to research
- **Reuse** deposited data or measures from other studies



Communities & Collections

Community: Define your research group or other collaborative unit

Collection: Create multiple Collections under the umbrella of the Community. Within Collections, deposit and describe your:

Phenotype Definitions Definitions Characterizations Evaluations Metadata Dissemination Strategy <u>Clinical Studies</u> Research Proposals Protocols Data Management Plans Methods Descriptions Measures Case Reports Datasets and Analyses

Collections bring together related groupings of documentation to communicate process, enable sharing of results, and support publication, compliance, and reproducibility



Collections & Clinical Studies

Store multiple datasets with large numbers of detailed results from each analysis and re-use of data generated by a single study.

Results presented in InvenioRDM are:

- easy to find
- browsable
- publicly available
- citeable

Hone in on the results you seek using InvenioRDM's robust metadata of subject and resource type terms. January 10, 2017

NCT02592655 Individual results for tourniquet study

Migura, Marcus

Individual level results for tourniquet study with clinicaltrials.gov identifier NCT02592655

Files (38.8 kB)						
Name	Size					
NCT02592655 INDIVIDUAL LEVEL RESULTS.xlsx	38.8 kB	2				

File Home	Insert	Page Layo	out Formulas	Data Revi	ew View	Acrobat	Power Pivot	🖓 Tell me	what you wan	it to do		
💼 🔏 Cut	Cal	libri	* 11 * A		≫·- ₽\	Wrap Text	Numb	er *			Normal	
Copy -									≠ Conditions			
🚽 🚿 Format	Painter B	ΙŪ·	••••••••••••••••••••••••••••••••••••••	* = = =	€ 3 🔛	Merge & Cer	iter * \$ *	% *	Formatting	Table *	Neutral	
Clipboard	r5	F	ont	rs.	Alignment		- G - 1	Number 5				Style
					-							
J56 ÷	: X	✓ f:	r									
	B	C	D	F	F	G	н	1	1	К	1 E	
2	1990 - 1990 - 1990 - 1990 - 1990 - 1990 - 1990 - 1990 - 1990 - 1990 - 1990 - 1990 - 1990 - 1990 - 1990 - 1990 -	-		-					-		-	
3 Summary dat	a for all 19	participar	its.									<u> </u>
4		AGE	BRACH R	BRACH L	DP R	DPL	PTR	PTL	ABLR	ABIL	HAND	1
5 MIN		18	113.00	117.00	135.00	135.00	138.00	137.00	1.02	1.00)	1
6 RANGE		30	43.00	40.00	50.00	44.00	41.00	46.00	0.33	0.34	1	1
7 MAX		48	156.00	157.00	185.00	179.00	179.00	183.00	1.35	1.34	1	
8 MEAN		33.74	131.11	132.74	150.79	150.21	151.95	151.89	1.15	1.15	5	
9 MEDIAN		33.00	131.00	133.00	149.00	150.00	150.00	151.00	1.16	1.15	5	
0 SD		9.25	10.13	10.05	12.08	11.55	10.33	11.61	0.09	0.09	9	
1												
2 Summary dat	a for only t	he 11 par	ticipants to rece	ive two windlas	s tourniquet	s.						
3		AGE	BRACH R	BRACH L	DP R	DP L	PTR	PTL	ABI R	ABI L	HAND	
4 MIN		18	114.00	120.00	136.00	140.00	140.00	137.00	1.02	1.00)	
5 RANGE		30	42.00	37.00	49.00	39.00	39.00	46.00	0.33	0.34	1	
6 MAX		48	156.00	157.00	185.00	179.00	179.00	183.00	1.35	1.34	1	
7 MEAN		36	132.82	135.55	153.64	154.09	154.55	155.64	1.15	1.15	5	
8 MEDIAN		35	131.00	133.00	152.00	151.00	154.00	153.00	1.15	1.15	5	
9 SD		9	10.76	10.76	14.09	12.36	11.60	12.77	0.11	0.11	L	
0												
1 Summary dat	a for the 8	participan	ts to receive on	y one windlsass	tourniquet.	2						
2		AGE	BRACH R	BRACH L	DP R	DP L	PT R	PT L	ABI R	ABI L	HAND	-
3 MIN		18.00	113.00	117.00	135.00	135.00	138.00	138.00	1.04	1.04	1	-
4 RANGE		26.00	30.00	20.00	27.00	24.00	20.00	23.00	0.23	0.19	9	_
5 MAX		44.00	143.00	137.00	162.00	159.00	158.00	161.00	1.27	1.23	3	+
6 MEAN		30.50	128.75	128.88	146.88	144.88	148.38	146.75	1.15	1.13	3	-
7 MEDIAN		31.00	129.00	130.50	148.00	143.50	148.00	145.00	1.17	1.15	5	-
8 SD		9.49	9.36	8.06	7.83	8.27	7.56	7.83	0.08	0.08	3	
0												

Publication date:	
January 10, 2017	
DOI:	
DOI 10.5281/zenodo.237773	
Keyword(s):	
tourniquet NCT02592655	
11	

Creative Commons Attribution 4.0 International

Share

Dataset Open Access

V

Download

Cite as

Migura, Marcus. (2017). NCT02592655 Individual results for tourniquet study [Data set]. Zenodo. http://doi.org/10.5281/zenodo.237773



Properly attribute all contributors to research



InvenioRDM incorporates **contributor roles** for all records. Deposit your SQL code, statistical analysis plan, database code, and other study documentation; receive credit, and group all documents in a Collection



Gonzales, S., O'Keefe, L., Gutzman, K., Viger, G., Wescott, A., Farrow, B., . . . Holmes, K. (n.d.). Personas for the Translational Workforce. *Journal of Clinical and Translational Science*, 1-27. doi:10.1017/cts.2020.2

Collaborators: Work with them and discover new ones

User 1's files





InvenioRDM will allow **private record sharing**, so researchers can:

- Share files with each other, but not anyone else in the university community or the public
- Vet materials collaboratively and privately before switching records to 'public' for open access/data sharing

InvenioRDM will have a social component, allowing researchers to:

- Follow other researchers
- Receive updates when someone they follow deposits something
- Manage requests to access files represented by a metadata-only record



The community





https://thenounproject.com/



The turn-key research data management repository 4 Launching in the summer 2020

A Roadmap

We intend to be ready by summer 2020.

🗩 Talk

Join our project forum and collaborate.

📟 Chat

Find all the partners in our official chatroom.

🗘 Code

Have a look at InvenioRDM code evolution.

Events

InvenioRDM project events for partners



Sneak peak at the future InvenioRDM.



https://inveniosoftware.org/ and click on "RDM"

InvenioRDM collaborators





How can Invenio support the OHDSI community?



Some Use Cases

Our multi-institution health equity

collaborate with our communitybased partners and credit these

materials from community health

events, project materials, training

materials, annual reports, and lay

InvenioRDM helps us to be better

collaborators and the community.

project uses InvenioRDM to

partnerships. We can share

summaries of research.

partners, accountable to

We're managing a large multi-site project, harmonizing data from numerous sources and managing research projects. We want to create communities of practice to integrate theories, data, techniques, and tools.



I lead a large basic science research group. We use InvenioRDM to support reproducible science by packaging combined with big data mining, a desire to process collected data using the latest bioinformatics tools.



a way to pre-register protocols or research proposals, search on demographics of participants in similar studies, get insights into recruitment, share portions of study for compliance.

My team wants to find out about clinical trial opportunities to offer patients all options for treatment. It is important to us to openly share the latest research with patients. InvenioRDM communities give us a way to make these materials openly available and packaged in a cohesive and attractive manner. As resources are updated, we can upload the new versions and track access.



I'm an early career researcher just getting started on my research career. I need to "put my best foot forward" to showcase my work and demonstrate my expertise and collaborations. Invenio gives me a way to make all of my research efforts findable and the metrics are helpful for reporting. and highlighting my impact to my leadership.

Our institute wants a way to publish and disseminate content such as our handbook, lay summaries, and more. We want to credit all contributors and produce an attractive and interactive resource that can be easily updated.



am a clinical researcher. I need

FAIR: OHDSI & InvenioRDM



InvenioRDM's records are made **findable** through each being issued a Digital Object Identifier (DOI), and through their metadata being indexed and made searchable immediately.

OMOP database summaries can be published in InvenioRDM as findable descriptor records to reference the database for reproducibility and citation nteroperable



InvenioRDM leverages metadata encoding (JSON) and vocabulary (FundRef, OpenAIRE, COAR Resource Types, etc.) standards to ensure maximum **interoperability** for records describing digital assets.

OMOP similarly ensures interoperability through its CDM and standardized vocabulary, and the OHDSI community goes beyond this work by providing a platform to enable an interoperable understanding of the analysis methods for healthcare data.



Metadata in InvenioRDM are **accessible** because they are retrievable using a standardized communications protocol which is free and universally implementable.



OMOP data can be mapped through similar open protocols through SQL interfaces, though largely for secure querying. Results of analyses in multiple OMOP databases can be cataloged in InvenioRDM, and these records retrieved through the open protocol OAI-PMH.







OHDSI's Metadata Working Group is actively working toward attaching provenance information to OMOP records.



Links

- Official InvenioRDM site: https://inveniosoftware.org/products/rdm/
- Roadmap: https://inveniosoftware.org/products/rdm/roadmap/
- GitHub: <u>https://github.com/inveniosoftware/invenio-app-rdm</u>
- Documentation: <u>https://invenio-app-rdm.readthedocs.io/en/latest</u>
- **Project Boards:** <u>https://github.com/orgs/inveniosoftware/projects</u>
- RFC (Request for Comments): <u>https://github.com/inveniosoftware/rfcs</u>

Northwestern's Proof of Concept: http://bit.ly/inveniordm-at-nu

- Test Login: gla3975
- Password: InvenioRDM@NU_2019

Install your own instance! https://inveniordm.docs.cern.ch/



With thanks...

Teams

- The Invenio team @ CERN & RDM collaborators (here)
- Galter Health Sciences Library & Learning Center
- Northwestern University Clinical and Translational Sciences
 Institute (NUCATS)
- CTSA Program Center for Data to Health (CD2H) team
- The NU Institute for Innovations in Developmental Sciences
- Confederation of OA Repositories (COAR)

Support

Work presented here is supported in part by: CERN Knowledge Transfer Fund CD2H: U24TR002306 (NCATS) NUCATS: UL1TR001422 (NCATS)

All of the InvenioRDM project partners



Sara Gonzales



Guillaume Viger



Lisa O'Keefe



Matt Carson







