

1. Overview

The Datasets Discovery Index (OmicsDI) provides dataset discovery across a heterogeneous, distributed group of genomics, proteomics and metabolomics data resources spanning eight repositories in two continents and six organisations, including both open and controlled access data resources. The resource provides a short description of every dataset: accession, description, sample/data protocols biological evidences, publication, etc. The search capabilities offer a unique resource to search for omics datasets. This fact converts the DDI in the first resource worldwide to provide search capabilities through multi-omics experiments.

2. Major Partners

2.1 ProteomeXchange (<http://www.proteomexchange.org/>)

The ProteomeXchange Consortium is a collaboration of currently three major mass spectrometry proteomics data repositories, PRIDE (<http://www.ebi.ac.uk/pride/archive/>) at EMBL-EBI in Cambridge (**UK**), PeptideAtlas (<http://peptideatlas.org>) at ISB in Seattle (**US**), and MASSive (<http://massive.ucsd.edu>) at UCSD (**US**), offering a unified data deposition and discovery strategy across all three repositories. ProteomeXchange is a distributed database infrastructure; the potentially very large raw data component of the data is only held at the original submission database, while the searchable metadata is centrally collected and indexed. All ProteomeXchange data is fully open after release of the associated publication.

2.2 MetabolomeXchange (<http://metabolomexchange.org/>)

MetabolomeXchange is a collaboration of 4 major metabolomics repositories, with a total of 10 partners contributing. MetabolomeXchange was inspired by and is implementing similar coordination strategies to ProteomeXchange. The founding partners are MetaboLights at EMBL-EBI (<http://www.ebi.ac.uk/metabolights/>) (**UK**), Metabolomics Repository Bordeaux (<http://www.cbib.u-bordeaux2.fr/>) (**FR**), Golm Metabolome Database (<http://gmd.mpimp-golm.mpg.de/>) and the Metabolomics Workbench

(<http://www.metabolomicsworkbench.org/>) (**US**). The Metabolomics Workbench is a NIH funded collaboration of 6 Regional Comprehensive Metabolomics Resource Cores. MetabolomeXchange started accepting metadata submissions in summer 2014, and reached 200 public datasets in March 2015.

2.3 The European Genome-Phenome Archive

The European Genome-Phenome Archive (EGA) provides a service for the permanent archiving and distribution of personally identifiable genetic and phenotypic data resulting from biomedical research projects. Data at EGA was collected from individuals whose consent agreements authorise data release only for specific research use to bona fide researchers. Strict protocols govern how information is managed, stored and distributed by the EGA project. The EGA comprises a public metadata section, allowing searching and identifying relevant studies, and the controlled access data section. Access to the data section for a particular study is only granted after validation of a research proposal through the relevant ethics approval.

2.4. Current Databases

The original DDI project provides access to six different databases with proteomics, genomics and metabolomics data from Europe and United States (see updated list (<http://localhost:8000/app/databases.html>)). The original list includes: PRIDE (proteomics, UK), PeptideAtlas (proteomics, US), MassIVE (proteomics, US), Metabolights (metabolomics, UK), MetabolomeWorkbench (metabolomics, US), EGA (genomics, UK).

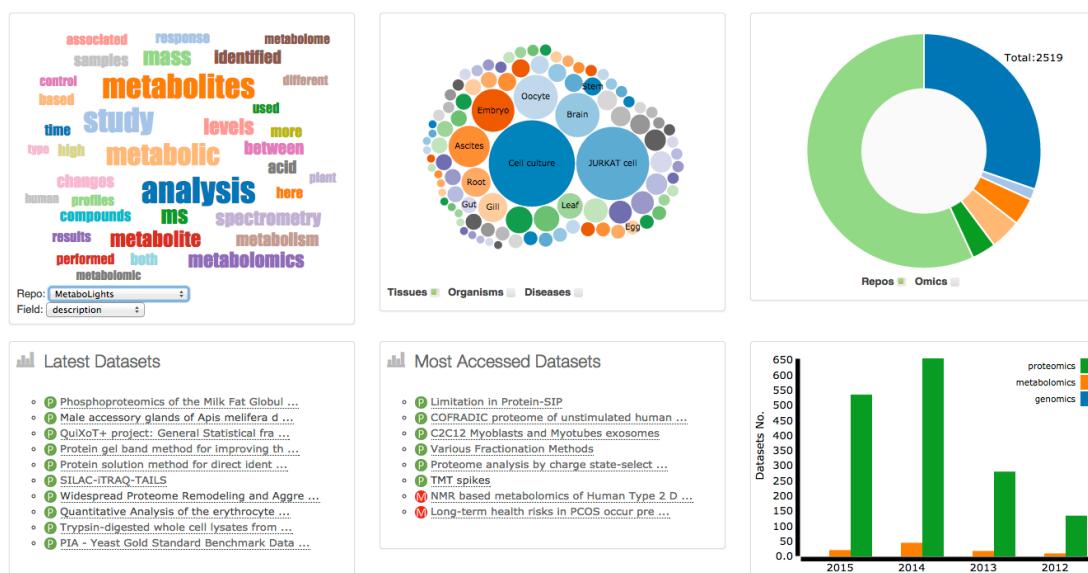
The project is open for new partners and databases (please contact us: pride-support@ebi.ac.uk). If you are interested in the architecture and the metadata that the resource should provide see section [?](#).

3. DDI web application

The main goal of DDI project is to have a way to search interesting datasets across omics repositories. The main web application and web service (see <http://localhost:8000/app/api.html>) allow the user to search and navigate through the DDI datasets. The DDI web application has two main different ways of navigating the data: (i) using the home page navigation blocks or (ii) the search box.

3.1 Navigate the data from home page

The DDI home page provides different blocks to navigate through the datasets, some of them are: 2D WordCloud; the species/organism/diseases bubble chart, repo/omics pie chart, Latest datasets, Most accessed datasets, Datasets per year. All the charts allow the user to search the data using the specific attribute. These boxes also act as an statistic component of the resource : for example the pie chart shows how many datasets for each repository and omics the resource contains.



A **TagCloud or WordCloud** is a visual representation for metadata, typically used to depict keyword metadata (tags) on datasets, or to visualize free form text. The WordCloud is build using the more frequently words for every database/repository. The DDI **WordCloud** can be consider as a two dimensional term representation where the user can select the database and the field they want to look for: description vs database. The user can click the highlight word in the wordcloud to search for this term in the resource.

The **bubble chart block** allows the users to navigate the data using three main categories: Tissues, Organisms, and Diseases. The user can click in the bubble and it will be redirected to the search using the clicked term.

The **Repo/Omics pie chart** and the **Omics vs Year bar chart** allow the users navigate the data using the omics categories (metabolomics, proteomics and genomics). The user can click a bar or the pie and it will be redirected to the search using the clicked term.

The **Latest Datasets** and **Most accessed datasets** blocks provide a list of the datasets by the tow categories.

3.2 Searching datasets

The main search box in DDI allows the user to search datasets using different keywords. The main search redirects the user to the browser page where the user can see the results of the search (see section 3.3).

Databases Discovery Index is an integrated multiple Omics Data Discovery site.

Examples: cancer, Homo sapiens, Orbitrap, Q9HAU5, Phospho, Hela

The DDI search page is better than most of the partners (proteomeXchange, metabolomeXchange) searching the data because the data is also indexed using cross-references to other databases.

3.3.2 Searching using publication details

The user can use the PubMed identifier, Title, Authors or even terms from the publication abstract.

Databases Discovery Index is an integrated multiple Omics Data Discovery site.

Examples: cancer, Homo sapiens, Orbitrap, Q9HAU5, Phospho, Hela

3.3.3 Searching using biological Evidences

The search box allows the end-users to search data using biological evidences such as the list of the proteins identified in the proteomics experiment or the metabolite reported in the metabolomics experiment. For example if the user search for *3-methyl-2-oxobutanoic* in the resource it will found one dataset in Metabolights and five in Metabolome workbench that identified the current molecule.

The screenshot shows the DDI interface with a search bar containing '3-methyl-2-oxobutanoic'. Below the search bar, there are links for Help, About, and Feedback. The main content area displays search results for the term '3-methyl-2-oxobutanoic'. The results are paginated from 1 to 6, with a page size of 10 items per page. The results are listed in a table with columns for Title, Accession, Relevance, and Publication date.

Project	Accession	Relevance	Publication date
M Metabolic differences in ripening of wild and mutant cultivars of Solanum lycopersicum			
M Mixed meal tolerance			
M Lung Cancer Cells 4			
M Rat HCR/LCR Stamina Study			
M Simeone Pancreatic Cancer Cells			
M Caloric Restriction vs drugs			

On the left side of the interface, there are several refinement panels:

- Show results for:** Proteomics(0), Metabolomics(6), Genomics(0)
- Repository:** Find your repositories, MetaboLights (1), Metabolomics Workbench (5)
- Refine by:**
 - Organisms:** Find your species, Homo sapiens (3), Mus musculus (1), Rattus norvegicus (1)
 - Tissue:** Find your Tissues, LTQ Orbitrap Velos (Thermo Scientific) (1), Agilent 1200 LC/6530 qTOF MS (3), Agilent Q-TOF 6530 (2)
 - Disease:** Find your Disease, LTQ Orbitrap Velos (Thermo Scientific) (1), Agilent 1200 LC/6530 qTOF MS (3), Agilent Q-TOF 6530 (2)

3.3 Refining the Search results

The search results can be filter or refine using different categories, filters or terms. The DDI web application supports at the moment nine different refinements: (1) omics type, (2) repository/database, (3) organisms, (4) tissue, diseases, (5) modifications (proteomics), (6) instruments and platforms, (7) publication data, (8) technology type.

Show results for

(1) **P Proteomics (1,555)**
M Metabolomics (201)
G Genomics (0)

Repository

(2) **MassIVE (38)**
MetaboLights (94)
Metabolomics Workbench (1)
PeptideAtlas (84)
PRIDE (1,433)

Refine by

(3) Organisms
Find your species
 Homo sapiens (676)
 Mus musculus (212)
 Mycobacterium tuberculosis
 Streptomyces sp. CNB091 (1)

Tissue

(4) **Whole body (29)**
 Cell suspension culture (15)
 Lung (15)
 Root (13)
 Blood serum (12)
 Colon (11)

Disease

(5) **Find your Disease**
 Not available (1,399)
 Disease free (24)
 Breast cancer (12)
 Cervix carcinoma (9)

Modifications in Proteomics

(6) **Find your Modifications**
 No PTMs are included in the dataset (1)
 Unknown modification (6)
 Oxidation (100)
 Deamidated (18)

Instruments & Platforms

(7) **Find your Instruments & Platform**
 Instrument model (665)
 LTQ (105)
 LTQ Orbitrap Velos (255)
 Q Exactive (145)

Publication Date

(8) **Find your publication data**
 2015 (279)
 2014 (607)
 2013 (279)
 2012 (187)
 2011 (112)

Page 1 2 3 4 5 >> 252 Showing 1 - 10 of 2519 Page size 10 20 50 100 Sort by: Title Accession Relevance Publication date

G WTCCC case-control study for Bipolar Disorder
Project description: Not available
Organism: Not available
EGAS000000000001 (EGA)

G WTCCC case-control study for Bipolar Disorder - Combined Controls
Project description: Not available
Organism: Not available
EGAS000000000002 (EGA)

G WTCCC case-control study for Coronary Artery Disease
Project description: Not available
Organism: Not available
EGAS000000000003 (EGA)

G WTCCC case-control study for Coronary Artery Disease - Combined Controls
Project description: Not available
Organism: Not available
EGAS000000000004 (EGA)

G WTCCC case-control study for Coronary Artery Disease, Hypertension, T2D - combined cases
Project description: Not available
Organism: Not available
EGAS000000000005 (EGA)

G Genomewide Association Study of Inflammatory Bowel Disease
Project description: Not available
Organism: Not available
EGAS000000000006 (EGA)

G Genomewide Association Study of Inflammatory Bowel Disease - Combined Controls
Project description: Not available
Organism: Not available
EGAS000000000007 (EGA)

G WTCCC case-control study for Inflammatory Bowel Disease, T1D and RA - combined cases
Project description: Not available
Organism: Not available
EGAS000000000008 (EGA)

G WTCCC case-control study for Hypertension
Project description: Not available
Organism: Not available
EGAS000000000009 (EGA)

G WTCCC case-control study for Hypertension - Combined Controls
Project description: Not available
Organism: Not available
EGAS000000000010 (EGA)

The refine filters works in combination, if the user set two filters they will act at the same time.

4. Dataset View

The dataset View shows the information for every dataset in the resource. It contains four main information components: (1) the dataset description, (2) the sample/data protocol, (3) the publication information, and (4) the related datasets.

 Mass spectrometry based draft of the human proteome **(1)**

This PXD project contains two projects published on ProteomicsDB (<https://www.proteomicsDB.org>) as integral part of the publication. The first project entitled 'human body map' (<https://www.proteomicsdb.org/#projects/42>) involves the analysis of 36 different human tissues and body fluids. The second project entitled 'Cellzome adopted' includes a collection of raw files which comprises identifications of 'missing proteins'.

Repository:[PXD000865](#) [As supplied by PRIDE] Date: 2014-05-28

Protocol **(2)**

SAMPLE PROTOCOL: For peptide identification, tandem mass spectra were processed using Mascot Distiller and searched against Uniprot using Mascot and a target-decoy strategy. In parallel, raw MS files were also processed by MaxQuant/Andromeda. All search results and tandem mass spectra were imported into ProteomicsDB (<https://www.proteomicsdb.org>) and filtered at 1% PSM FDR and 5% local peptide length dependent FDR.

DATA PROTOCOL: Human tissue specimens were obtained from the bio bank of the TU München following approval of the study by the local ethics committee. Samples were collected within the first 30 minutes after resection, macroscopically resected by an experienced pathologist, snap frozen and stored in liquid nitrogen until use. Body fluids requiring no invasive procedures were provided by volunteers. Proteins were extracted under denaturing conditions and either separated by LDS-PAGE followed by in-gel protease digestion or digested in solution in the presence of chaotropic agents. Peptides were separated by ultra-high pressure nanoscale liquid chromatography on 75 µm x 40 cm reversed-phase columns using gradients of 2–32% acetonitrile in 0.1% formic acid, 5% DMSO. The LC was coupled directly to an Orbitrap mass spectrometer, operating in data dependent mode and using either resonance-type or beam-type collision induced dissociation.

Instruments: [LTQ Orbitrap Velos](#), [LTQ Orbitrap Elite](#)

 Mass-spectrometry-based draft of the human proteome. **(3)**

Proteomes are characterized by large protein-abundance differences, cell-type- and time-dependent expression patterns and post-translational modifications, all of which carry biological information that is not accessible by genomics or transcriptomics. Here we present a mass-spectrometry-based draft of the human proteome and a public, high-performance, in-memory database for real-time analysis of terabytes of big data, called ProteomicsDB. The information assembled from human tissues, cell lines ...[more](#)

PMID: [24870543](#) Publication: [1/ 1](#) 2014-05-01

Related Datasets **(4)**

 [Rat Neural Stem Cell Phosphoproteome](#) [PXD000976](#) (PRIDE) Made public: 2014-05-23

 [CPTAC System Suitability \(CompRef\) Study: Phosphoproteome, PNNL](#) [PXD001000](#) (PRIDE) Made public: 2014-09-16

 [The Pan-Human Library: A repository of assays to quantify 10 000 proteins by SWATH-MS](#) [PXD000953](#) (PRIDE) Made public: 2014-08-08

 [CPTAC System Suitability \(CompRef\) Study: Proteome, PNNL](#) [PXD000966](#) (PRIDE) Made public: 2014-09-16

 [Silkworm testis LC-MS/MS](#) [PXD000909](#) (PRIDE) Made public: 2014-07-28

 [Phospho-iTRAQ](#) [PXD001574](#) (PRIDE) Made public: 2015-01-30

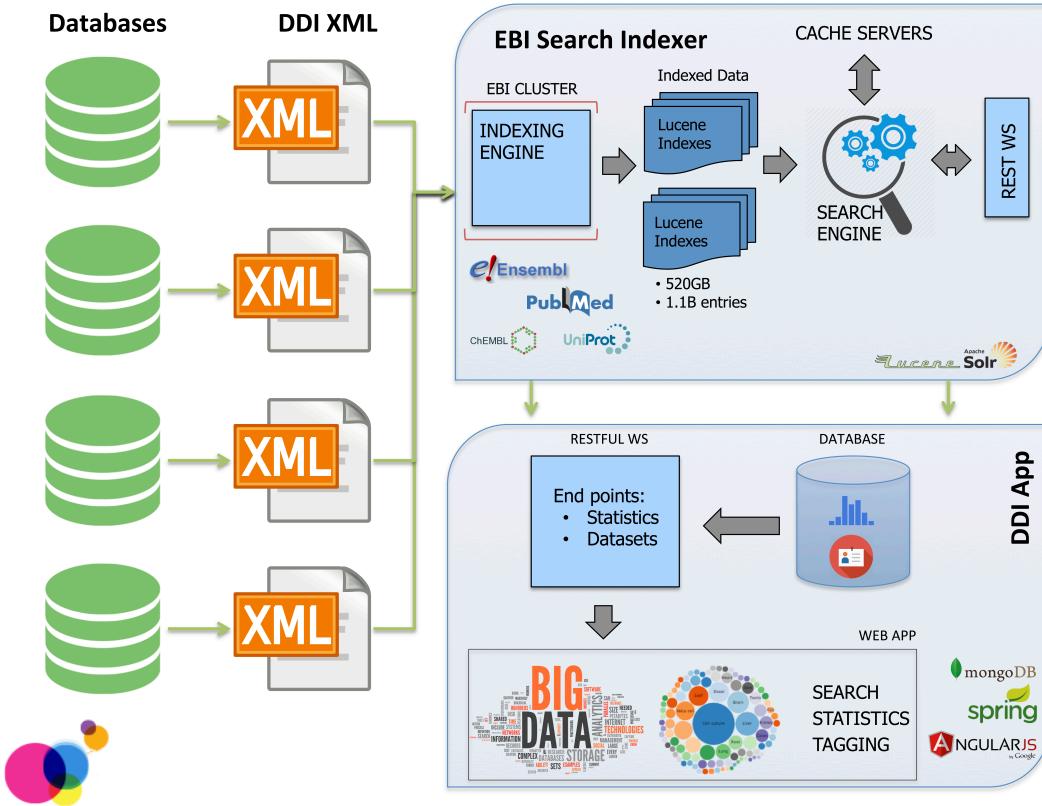
 [Quantitative phosphoproteomics unveils temporal dynamics of thrombin signaling in human endothelial cells](#) [PXD000597](#) (PRIDE) Made public: 2014-02-06

 [Insight into the protein composition of immunoglobulin light chain deposits of eyelid, orbital, and conjunctival amyloidosis.](#) [PXD000743](#) (PRIDE) Made public: 2014-05-29

5. DDI Architecture

The DDI project design is a modular architecture with three main components (Figure 1): (i) Database/Repository Adapters, (ii) EBI Search Indexer, (iii) DDI Application. In summary, the Database/Repository Adapters provides all the libraries and readers to translate the original repository/databases information into a universal and highly customizable XML file for each dataset (see section

5.1). The EBI Search Indexer is a Lucene-based framework that enables to index all the metadata and biological evidences for each dataset (see section 5.2). The DDI application component provides a MongoDB, web-services and web application for search and access the different datasets.



5.1 Database/Repository Adapters

The database/repository Adapter components are a set of readers and libraries that translate the data from the original repositories to a common XML-based file. The DDI-XML is a fully customizable XML file used to represent the data from the different datasets, it contains as mandatory for each dataset:

- Dataset ID
- Dataset Title
- Type of experiment
- Publication date
- Submitter details

A full description of the XML and the exporters can be found in Github (<https://github.com/BD2K-DDI>). If you need to add a new data resource, that is a new domain, get in touch first with the DDI team (peide-support@ebi.ac.uk).

5.2 EBI Search Indexer

In summary, the indexing pipeline:

- Based on Apache Lucene
- An automatized indexing cycle takes place on daily basis.
- Check every night, for every domain if new data is available.
- The process benefits from the indexing parallelization.
- All domains together: ~750GB, more than 300GB indexes
- While the software releases are less frequent, the data to index is checked daily