



Proteomics repositories integration using EUDAT resources

Rafael C Jimenez

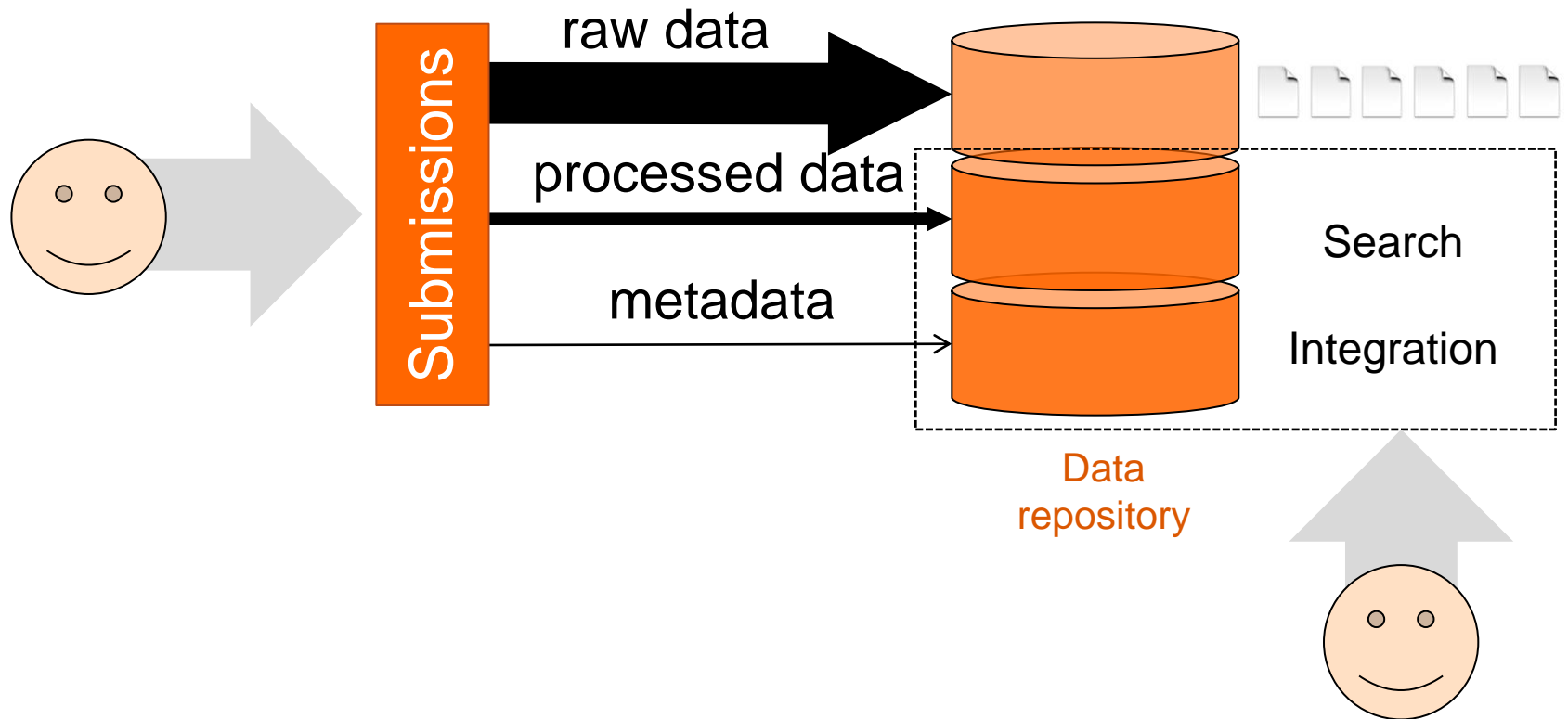
ELIXIR CTO

25 September 2014

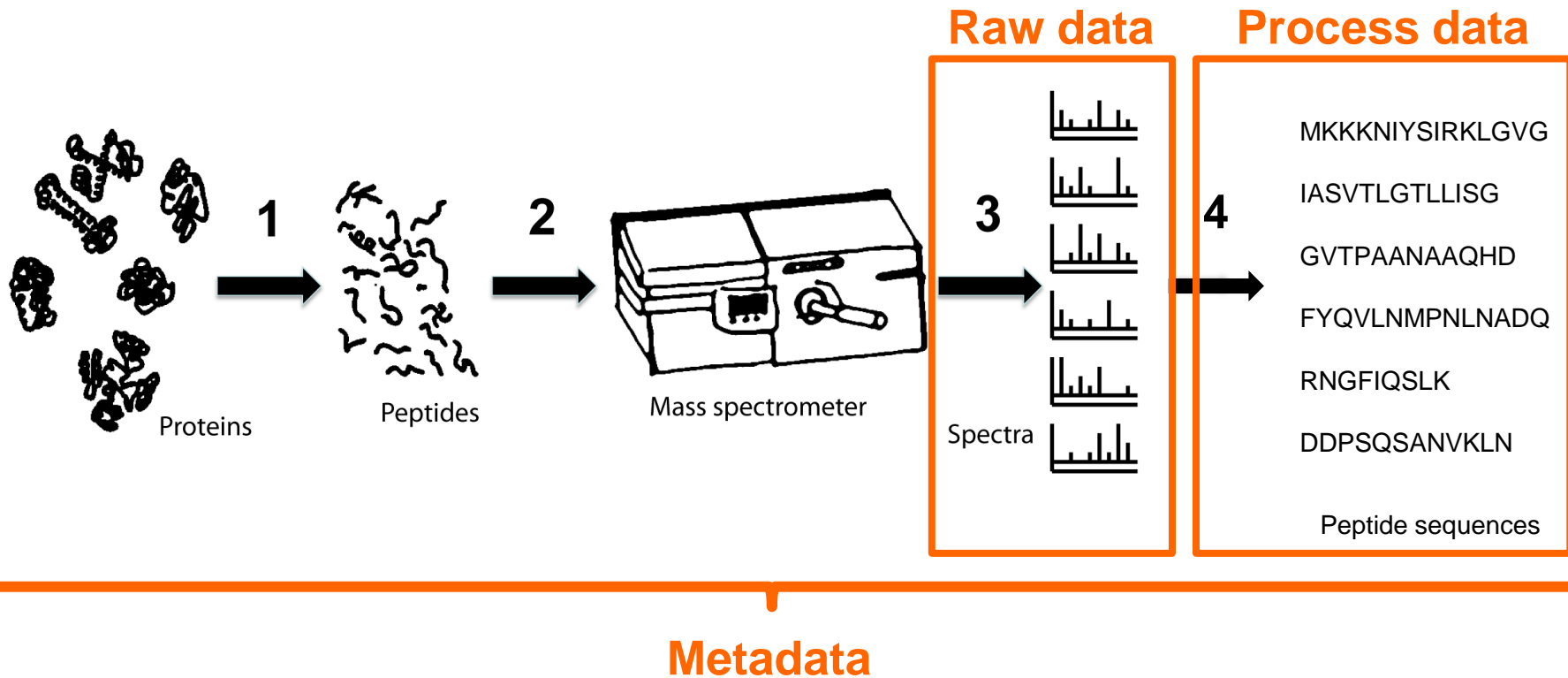


*European Life Sciences Infrastructure for Biological Information
www.elixir-europe.org*

Data submissions



Overview of shotgun proteomics data production

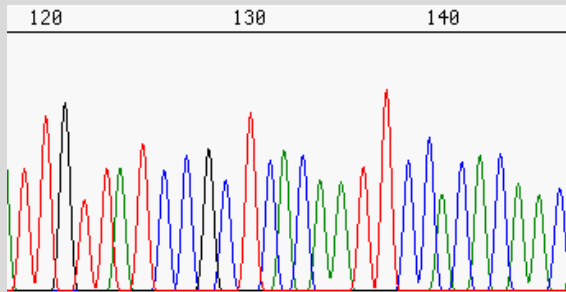


Noble WS, MacCoss MJ (2012) Computational and Statistical Analysis of Protein Mass Spectrometry Data. PLoS Comput Biol 8(1): e1002296. doi:10.1371/journal.pcbi.1002296

<http://www.ploscompbiol.org/article/info:doi/10.1371/journal.pcbi.1002296>

Data examples

Raw data

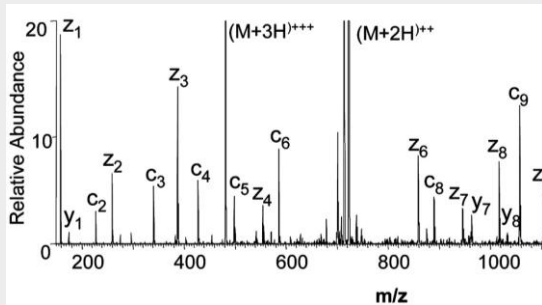


Process data

TTGTTATCCG...

Metadata

DNA
Human
Liver
Mitochondria
W. Smith
...



LPISASHSSK...

Peptide
Mouse
Heart
Nucleus
J. Heinz
...

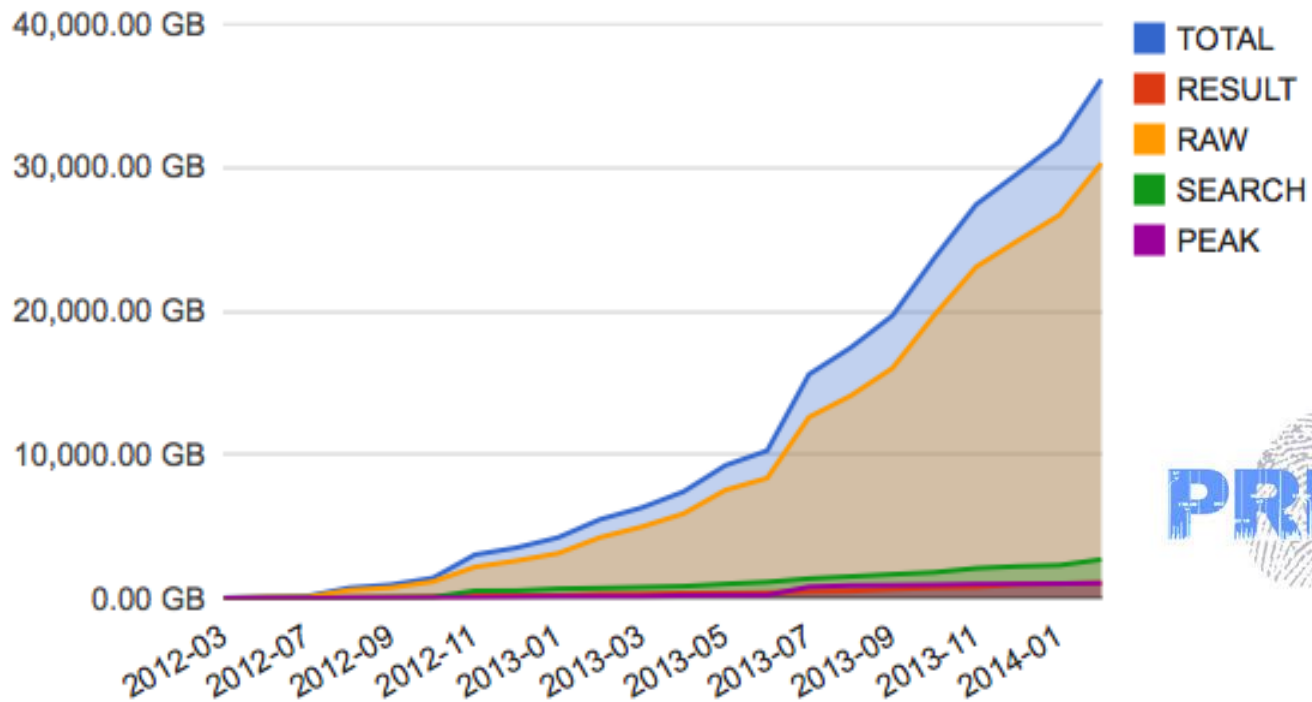
...

...

...

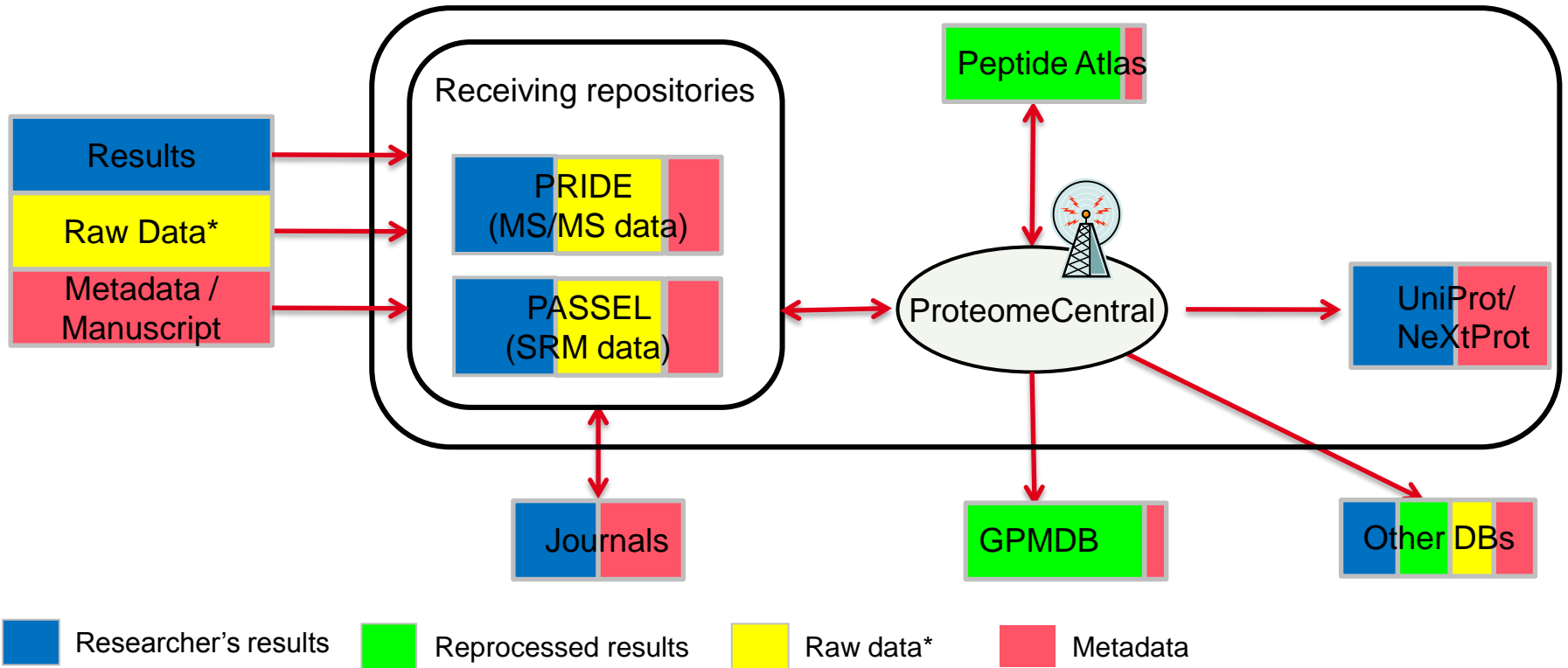
Proteomics data in PRIDE

~85% raw data

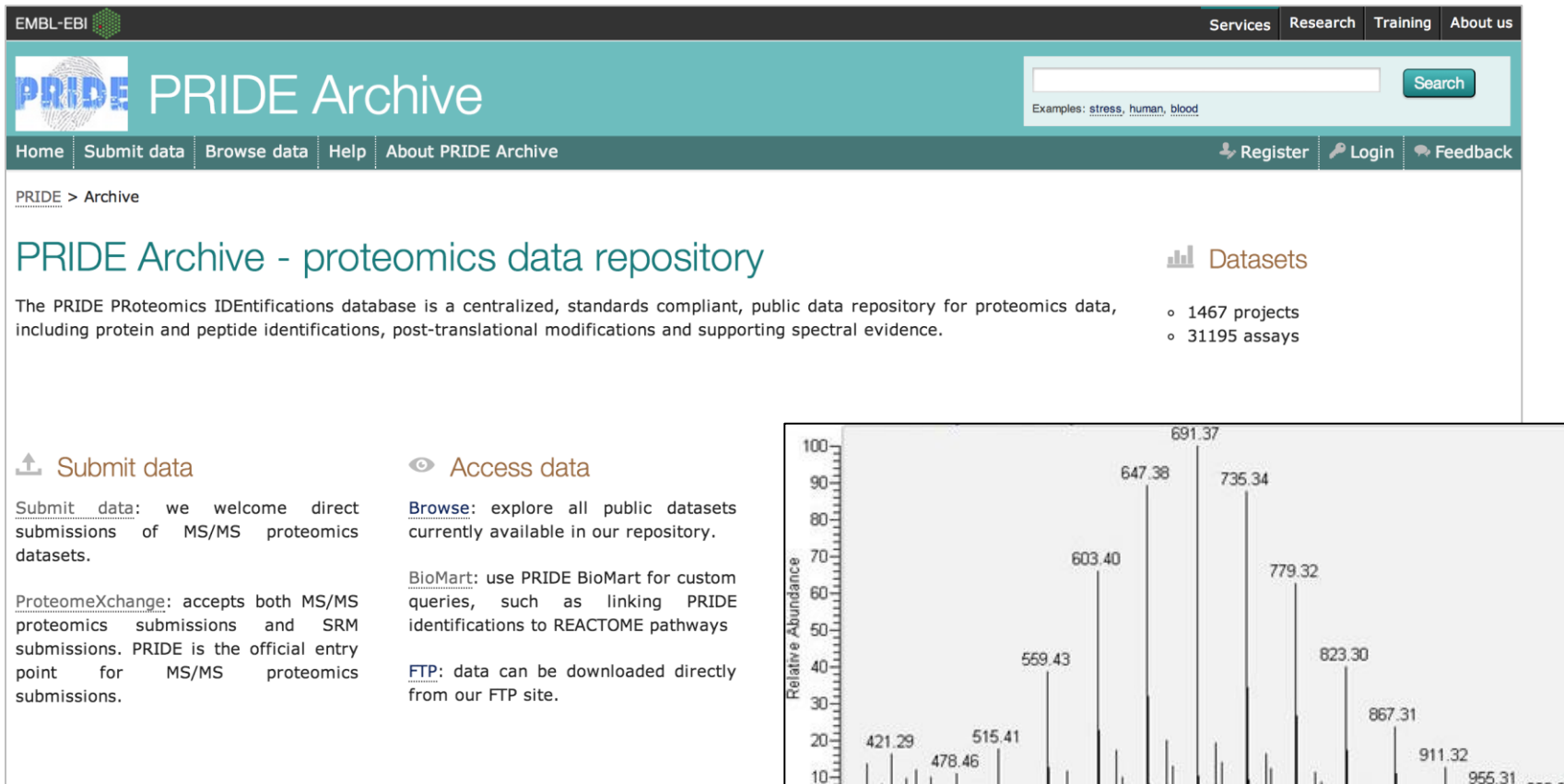


ProteomeXchange

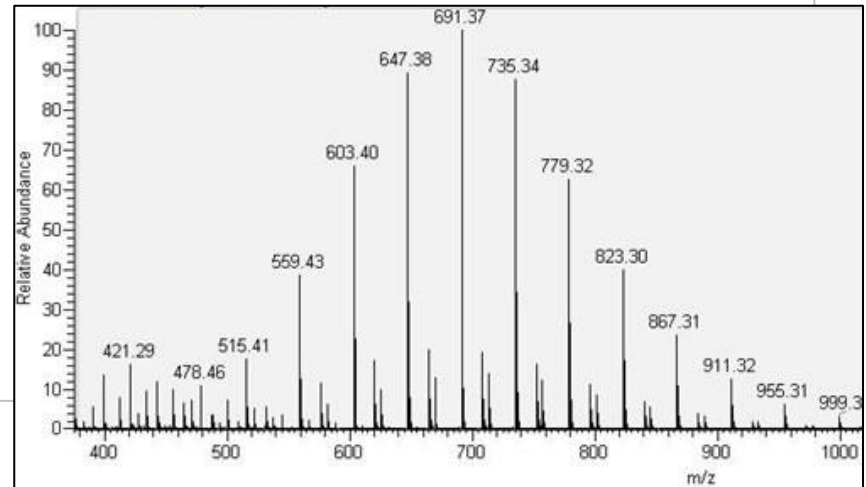
- **Framework to enable standard data submission and dissemination pipelines** between the main existing proteomics resources.



PRIDE (PRoteomics IDEntifications) database



The screenshot shows the PRIDE Archive website. At the top, there is a navigation bar with 'EMBL-EBI' on the left and 'Services', 'Research', 'Training', and 'About us' on the right. Below this is a teal header with the 'PRIDE Archive' logo and a search bar containing the text 'Examples: stress, human, blood'. A secondary navigation bar includes 'Home', 'Submit data', 'Browse data', 'Help', and 'About PRIDE Archive', along with 'Register', 'Login', and 'Feedback' links. The main content area features the title 'PRIDE Archive - proteomics data repository' and a 'Datasets' section with a bar chart icon and two bullet points: '1467 projects' and '31195 assays'. On the left, there are two sections: 'Submit data' with a description of direct submissions and ProteomeXchange, and 'Access data' with links for 'Browse', 'BioMart', and 'FTP'.



Mass spectrometry

- Origin:**
- 152 USA
 - 108 Germany
 - 67 United Kingdom
 - 53 Switzerland
 - 48 Netherlands
 - 42 China
 - 42 Canada
 - 41 France
 - 36 Spain
 - 33 Belgium
 - 25 Australia
 - 23 Sweden
 - 17 Japan
 - 16 Denmark
 - 13 Norway
 - 12 Finland
 - 12 India
 - 12 Taiwan
 - 10 Italy
 - 9 Republic of Korea
 - 8 Austria
 - 8 Ireland
 - 8 Brazil
 - 7 Singapore
 - 5 Israel
 - 5 Russia ...

Type:

- 273 PRIDE complete
- 501 PRIDE partial
- 47 PeptideAtlas/PASSEL complete

Submissions/year:

- 2012: 102
- 2013: 527
- 2014: 192

Access:

- 38.3% PRIDE public
- 5.3% PASSEL public
- 56% PRIDE private
- 0.4% PASSEL private

Top Species studied by at least 8 datasets:

- 381 *Homo sapiens*
- 100 *Mus musculus*
- 31 *Arabidopsis thaliana*
- 26 *Saccharomyces cerevisiae*
- 16 *Escherichia coli*
- 14 *Rattus norvegicus*
- 12 *Mycobacterium tuberculosis*
- 11 *Drosophila melanogaster*

~ 215 species in total

Data volume:

- Total: >40 TB
- Number of all files: >120,000
- PXD000320-324: ~ 5 TB
- PXD000065: ~ 1.4TB

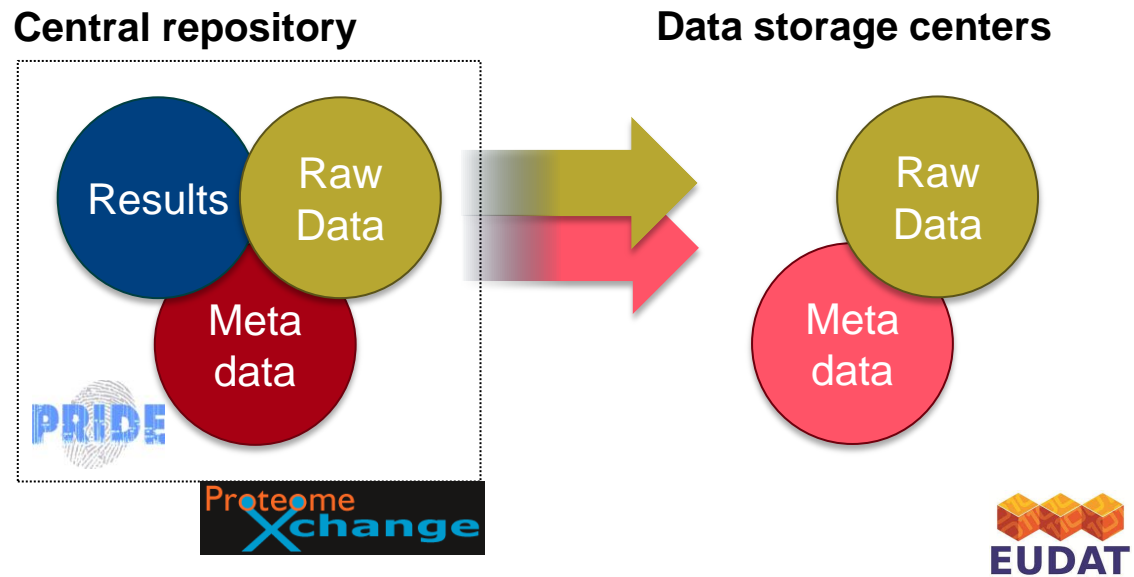


Pilot evolution

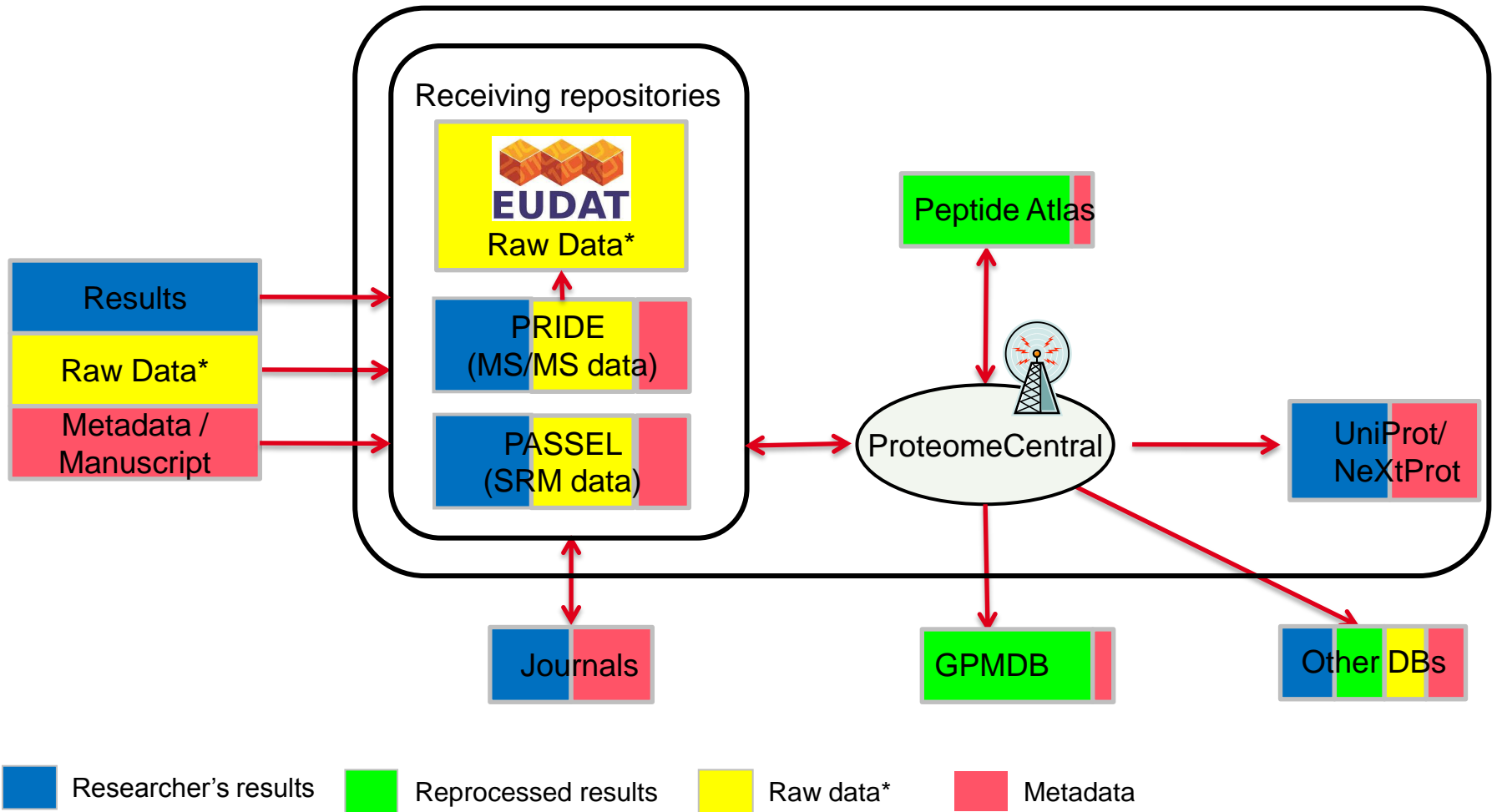
- Use EUDAT
 - Replication of ELIXIR data in EUDAT data centers
 - Delegation of ELIXIR data in EUDAT data centers

- Adopt EDUAT
 - Replication of ELIXIR data in ELIXIR data centers using EUDAT technology

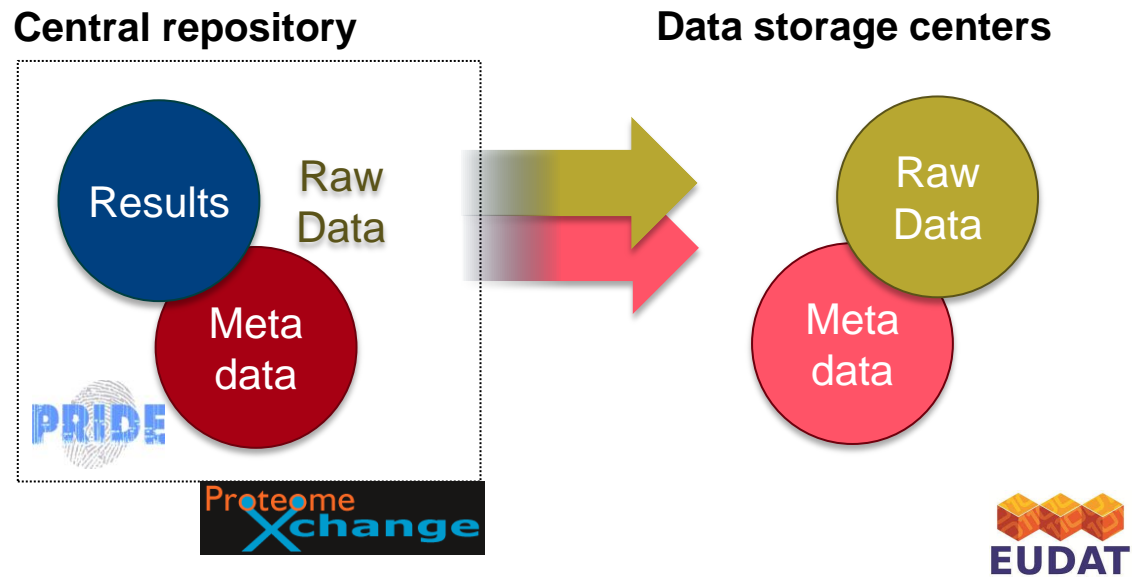
Replication of ELIXIR data in EUDAT data centers



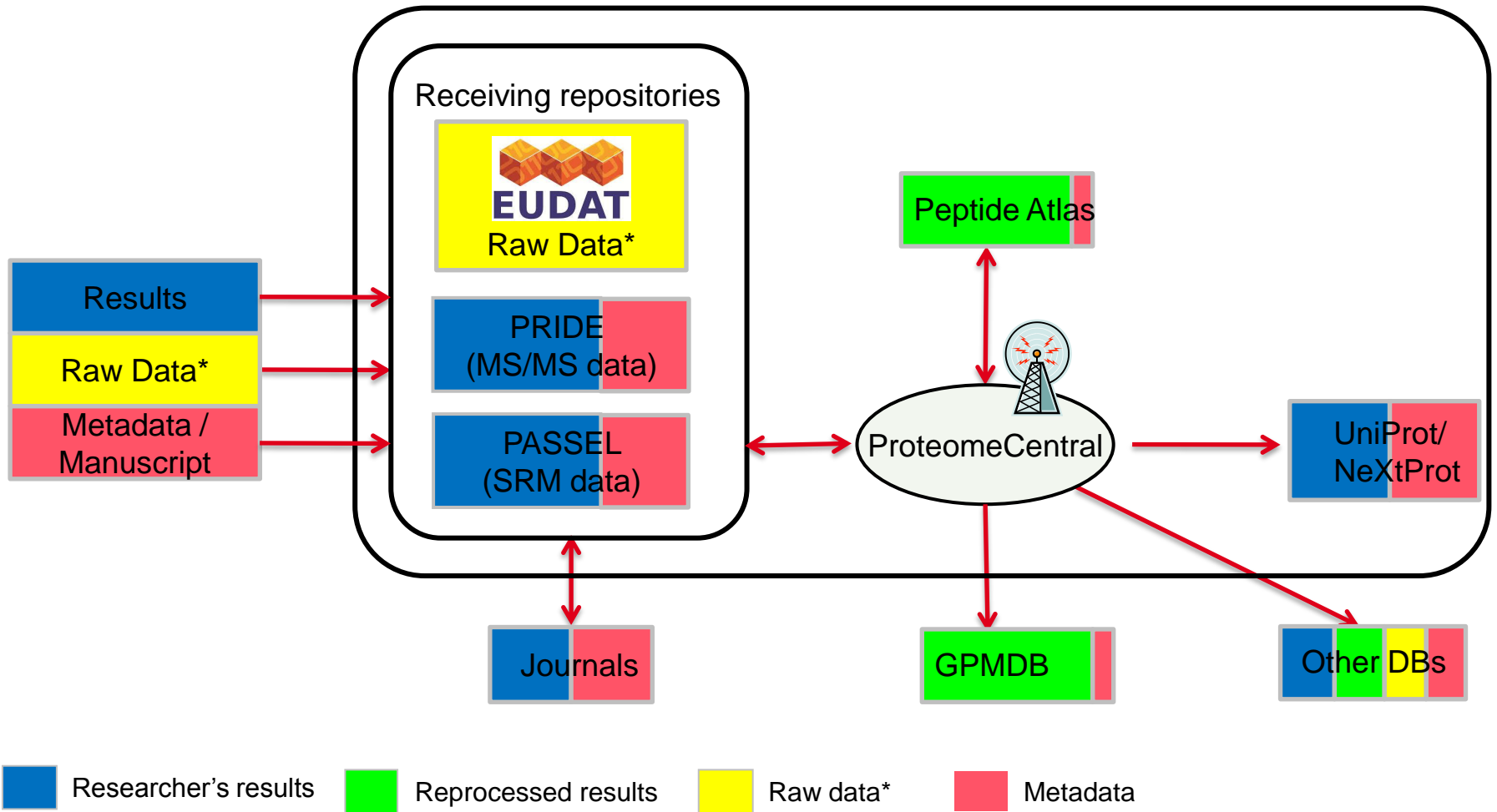
Replication of ELIXIR data in EUDAT data centers



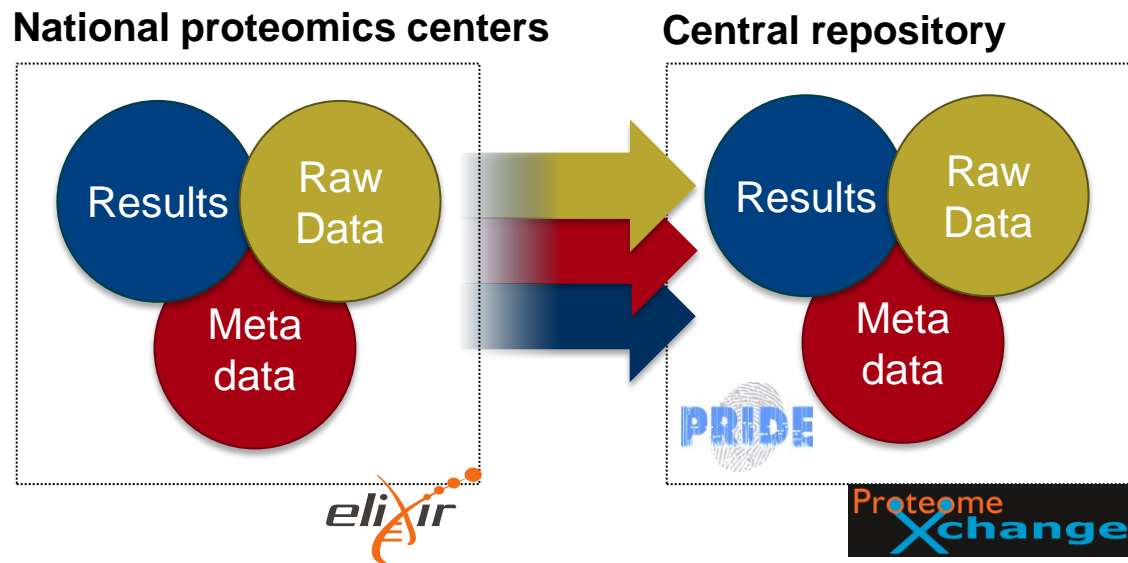
Delegation of ELIXIR data in EUDAT data centers



Delegation of ELIXIR data in EUDAT data centers



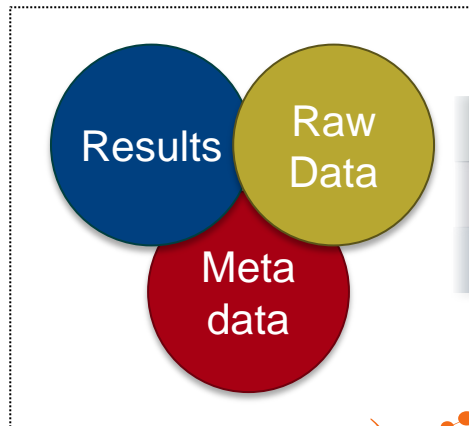
Replication of ELIXIR data in ELIXIR data centers using EUDAT technology



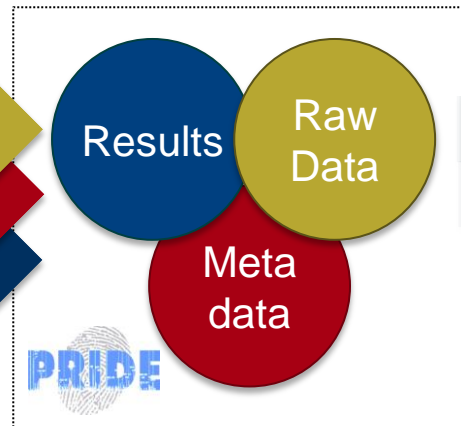
Plans

1.- ELXIR replication

National proteomics centers



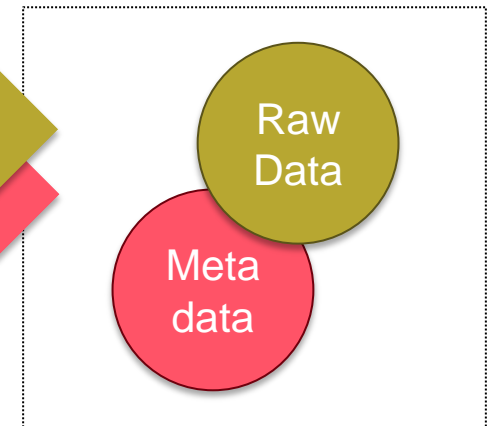
Central repository



PRIDE

ProteomeXchange

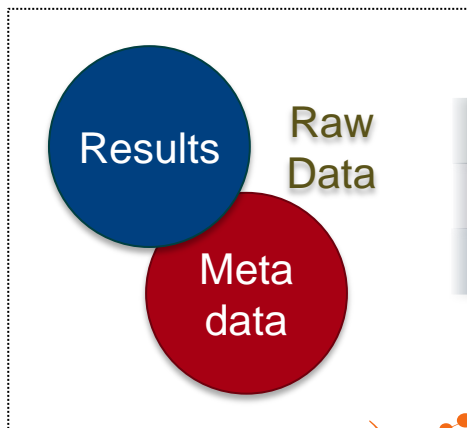
Data storage centers



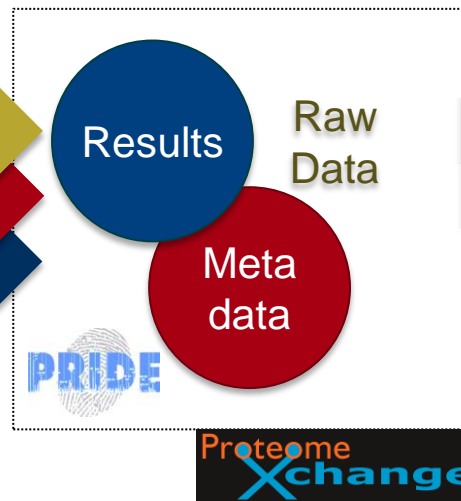
2.- EUDAT replication

Plans

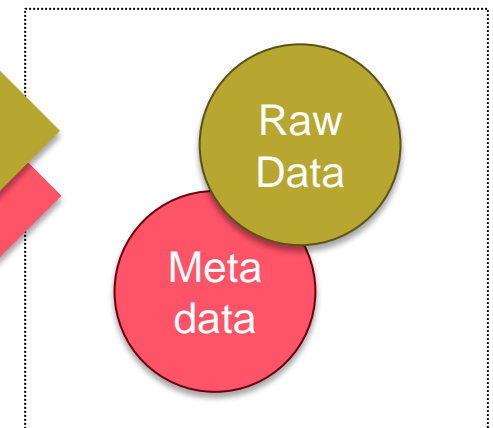
National proteomics centers



Central repository



Data storage centers



3.- delegation

ELIXIR Pilot action

ELIXIR BILS ProteomeXchange_27062014.pdf (1 page)



ELIXIR Pilot Action Proposal

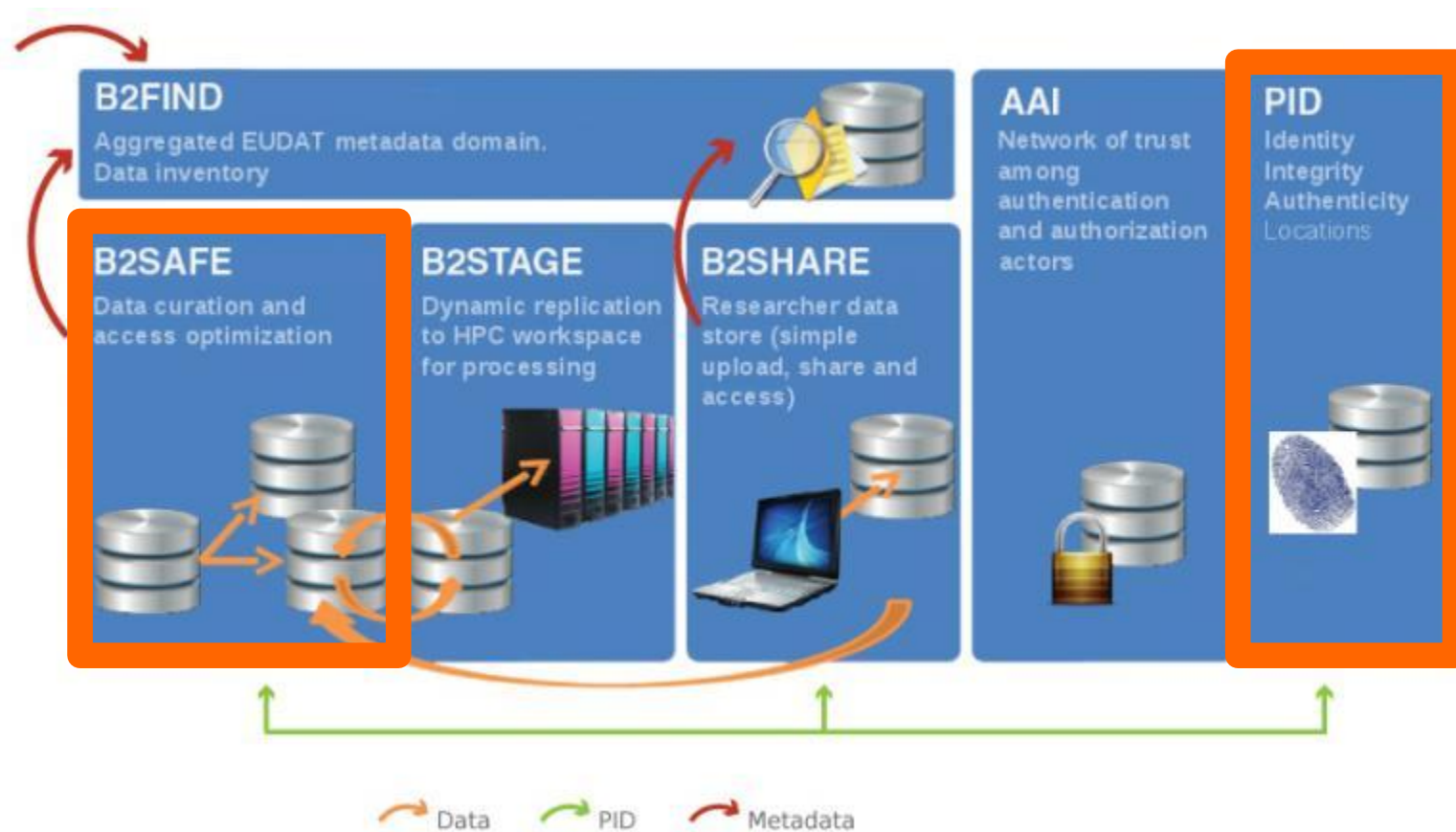
Date of submission: 2014-06-27

Pilot name: BILS-ProteomeXchange integration using EUDAT resources

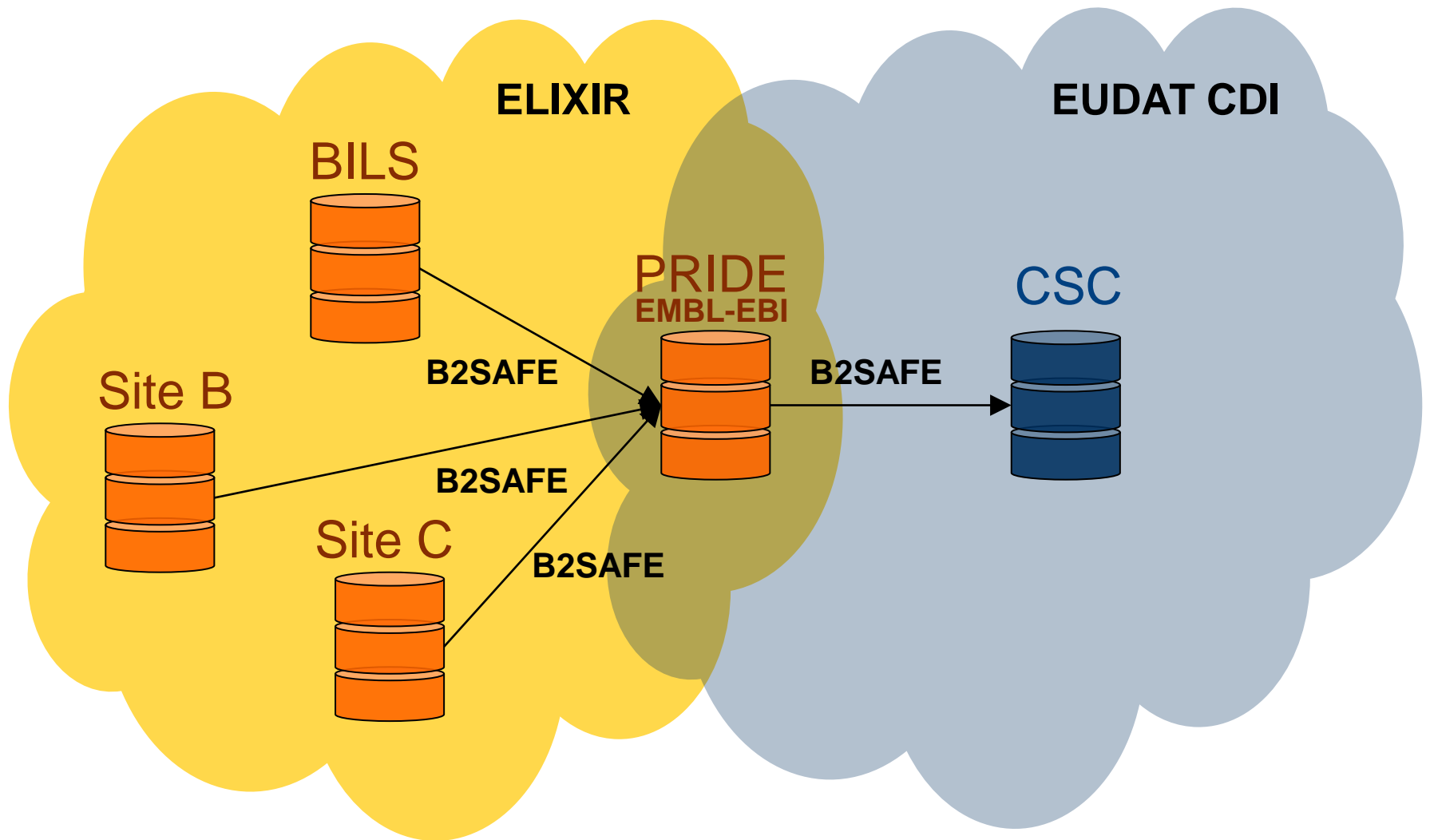
Motivation

While the current data deluge creates a need for distributed data storage and replication, it is essential to enable data access through a single access interface. In the proposed pilot action the aim is to integrate the raw data repositories for mass spectrometry (MS) proteomics data run by BILS (Sweden) and ProteomeXchange¹ (via the PRIDE database, EMBL-EBI, UK), using the European infrastructure EUDAT (<http://www.eudat.eu/>). The ProteomeXchange consortium has been recently set up to facilitate and standardise submission and dissemination practices for proteomics data resources. The proposed pilot is within the "Data resources" and "Tools" programmes. This pilot will report back to the ELIXIR community and will serve as an example to connect national data storage services and international repositories through ELIXIR. This pilot will also show the potential of collaboration among research infrastructures and e-infrastructures to better manage the data deluge. It will help to evaluate the requirements of such federated systems.

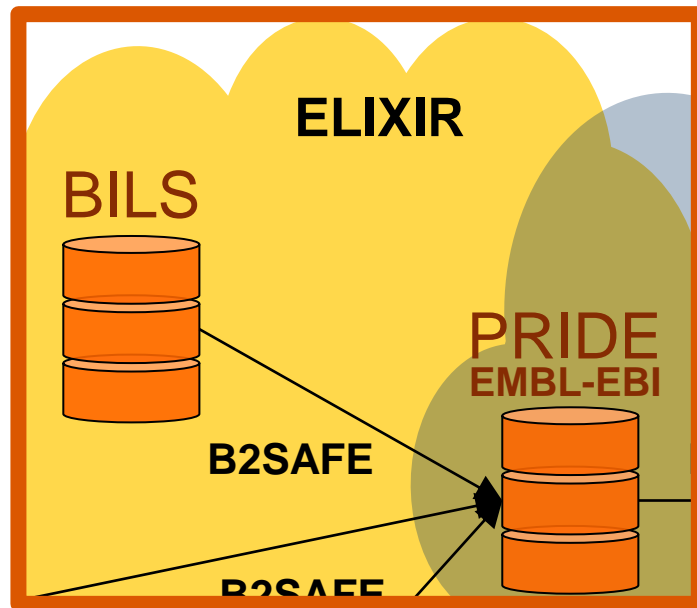
EUDAT services



File sharing model



Pilot – EUDAT adoption: ELIXIR replication

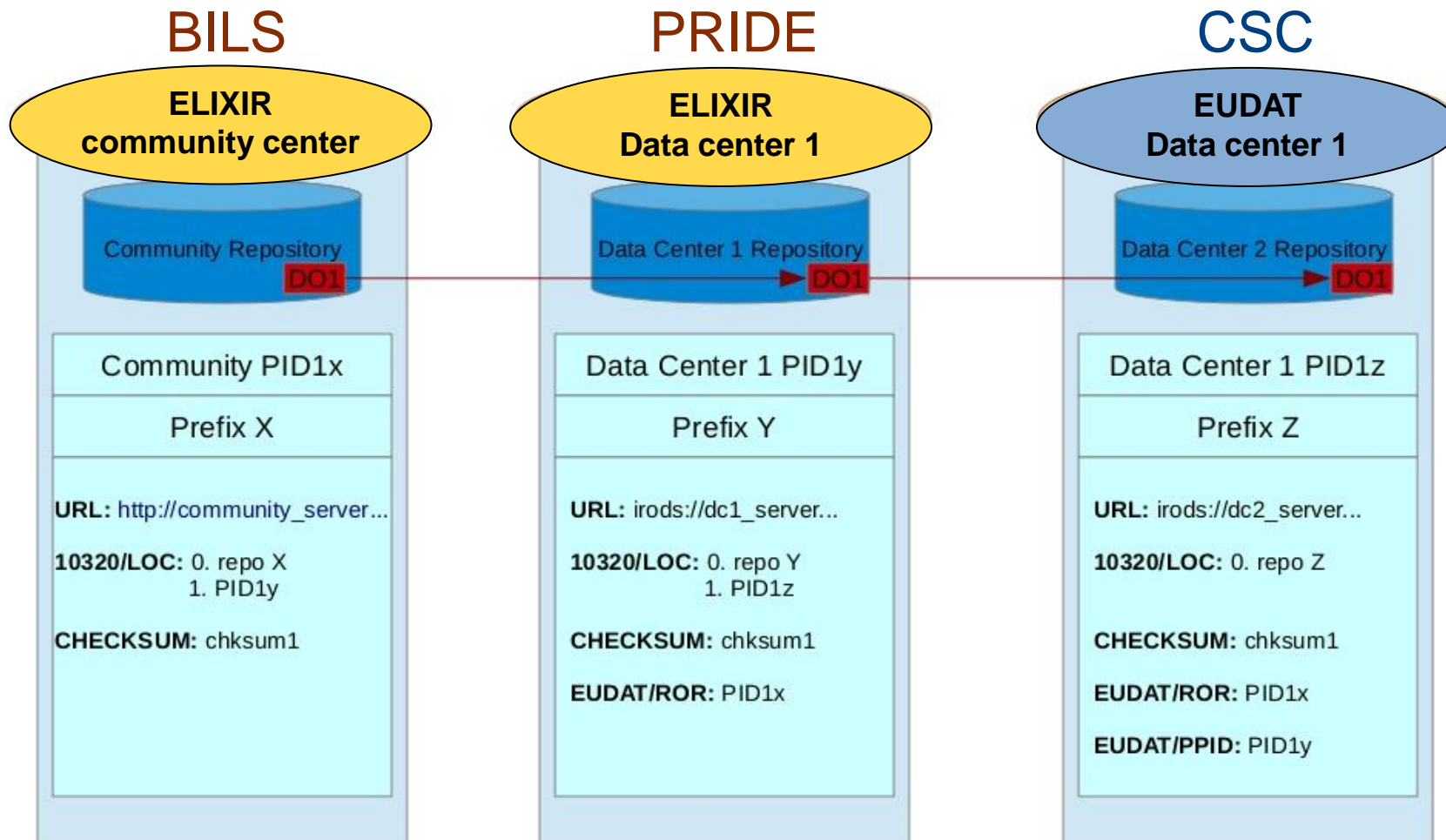


National proteomics centers

Central repository



PIDs



Status

- BILS
 - Migrating from existing Swestore dCache to iRODS
 - Testing compatibility with B2SAFE
 - Latest iRDOS not compatible with B2SAFE?
- PRIDE
 - iRODS service installed
 - B2SAFE module have been deployed at EMBL-EBI (PRIDE)
 - Test B2SAFE replication PRIDE -> CSC
 - DOI for datasets
 - PID for dataset files
 - Web service to associate datasets to dataset files

Status

In progress

- Handle System Registration
- Test requests of EPIC/EUDAT identifiers

Open questions

- BILS local PIDs?
- Sync back from PRIDE to BILS for modifications/additions at PRIDE?
- Data push or pull model?
- Replication of process data requires previous validation

Participants

EUDAT/CSC

- Jani Heikkinen
- Damien Lecarpentier

EMBL-EBI/systems

- Andy Jenkinson
- Steven Newhouse

EMBL-EBI/PRIDE

- Juan Antonio Vizcaíno
- Henning Hermjakob

BILS

- Mikael Borg
- Fredrik Levander
- Bengt Persson

ELIXIR Hub

- Rafael C Jimenez



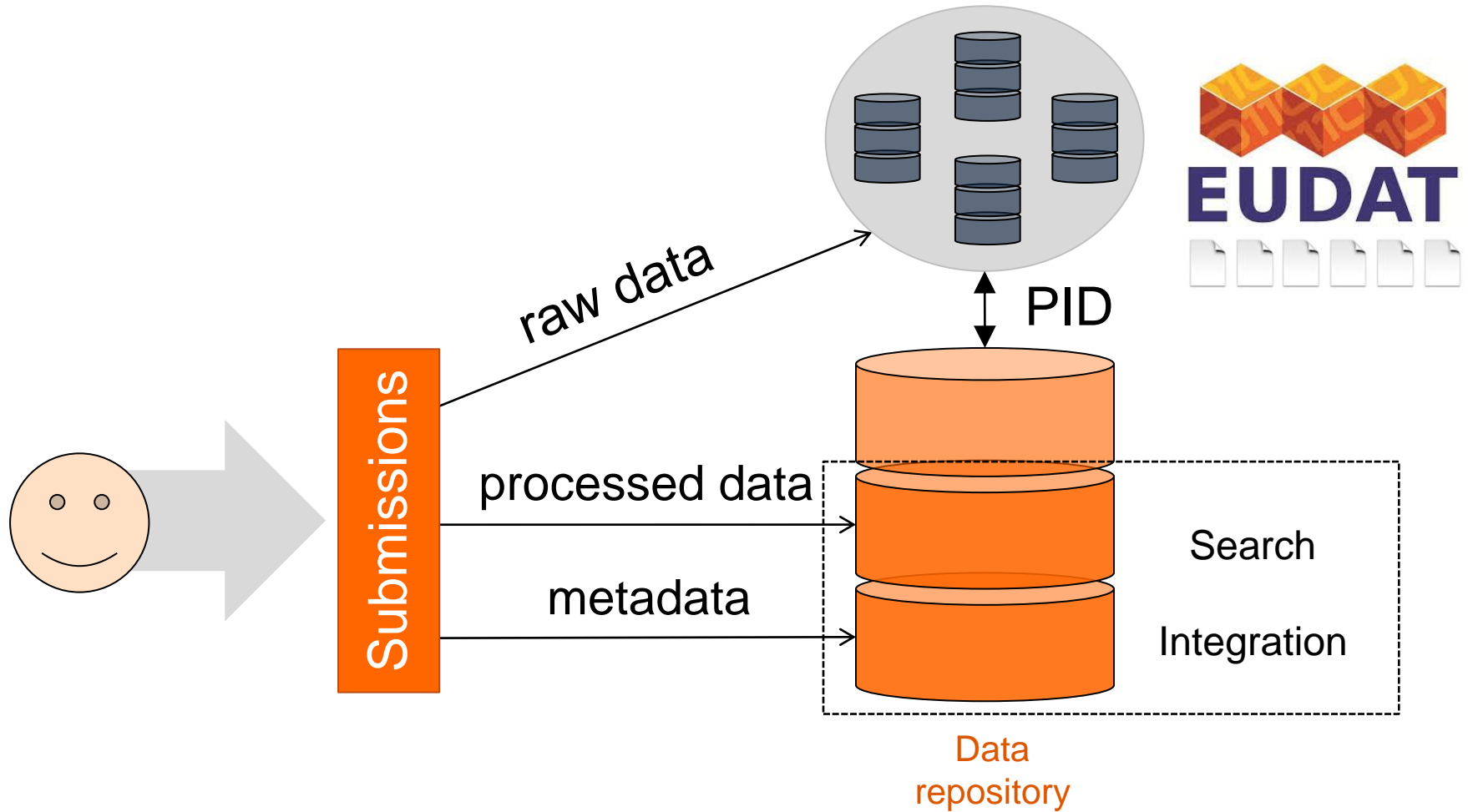
Thank you for your attention

European Life Sciences Infrastructure for Biological Information

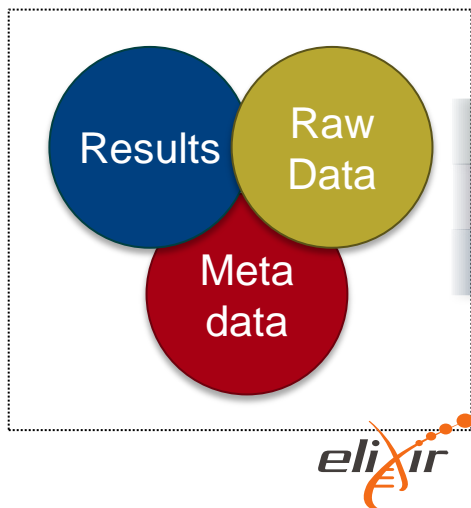
www.elixir-europe.org



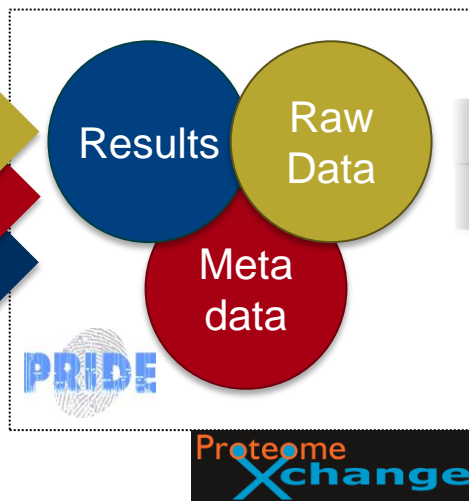
Delegation of raw data



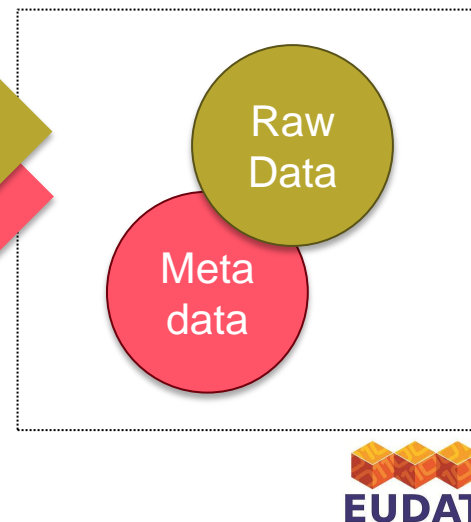
National proteomics centers



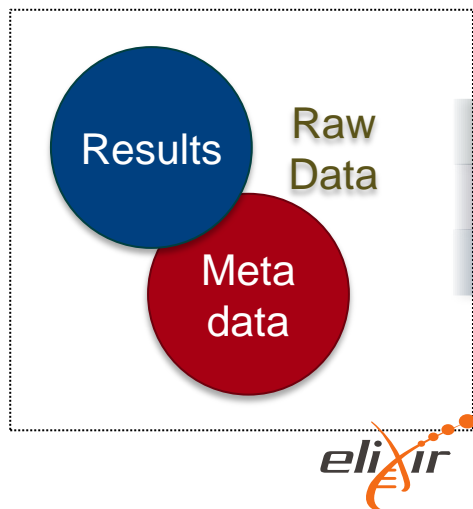
Central repository



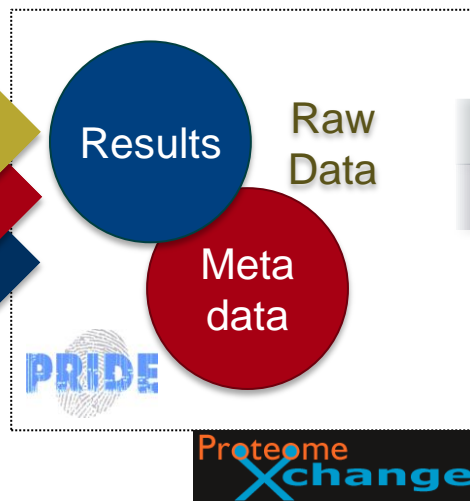
Data storage centers



National proteomics centers



Central repository



Data storage centers

