

Proteomics data *repositories: PRIDE*

Dr. Juan Antonio VIZCAINO

PRIDE Group coordinator
PRIDE team, Proteomics Services Group
PANDA group
European Bioinformatics Institute
Hinxton, Cambridge
United Kingdom



EBI is an Outstation of the European Molecular Biology Laboratory.

Overview ...

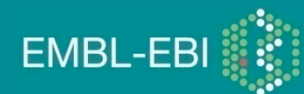
- Why sharing proteomics data?
- Introduction to existing proteomics repositories
- Proteomics data bottlenecks
- PRIDE in detail...
- ProteomeXchange consortium

A ONE-SLIDE INTRODUCTION TO MASS SPEC PROTEOMICS

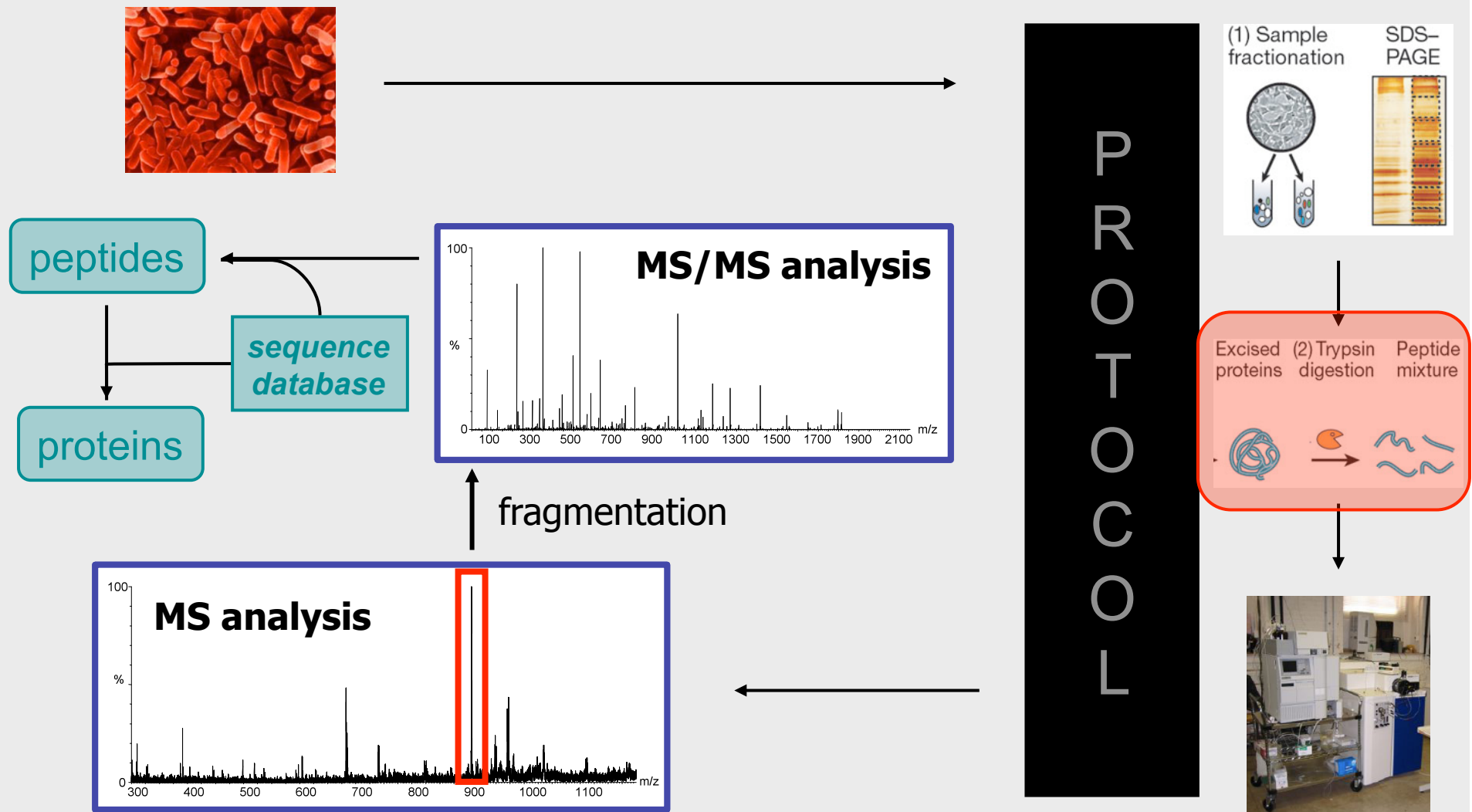
Juan A. Vizcaíno
juan@ebi.ac.uk



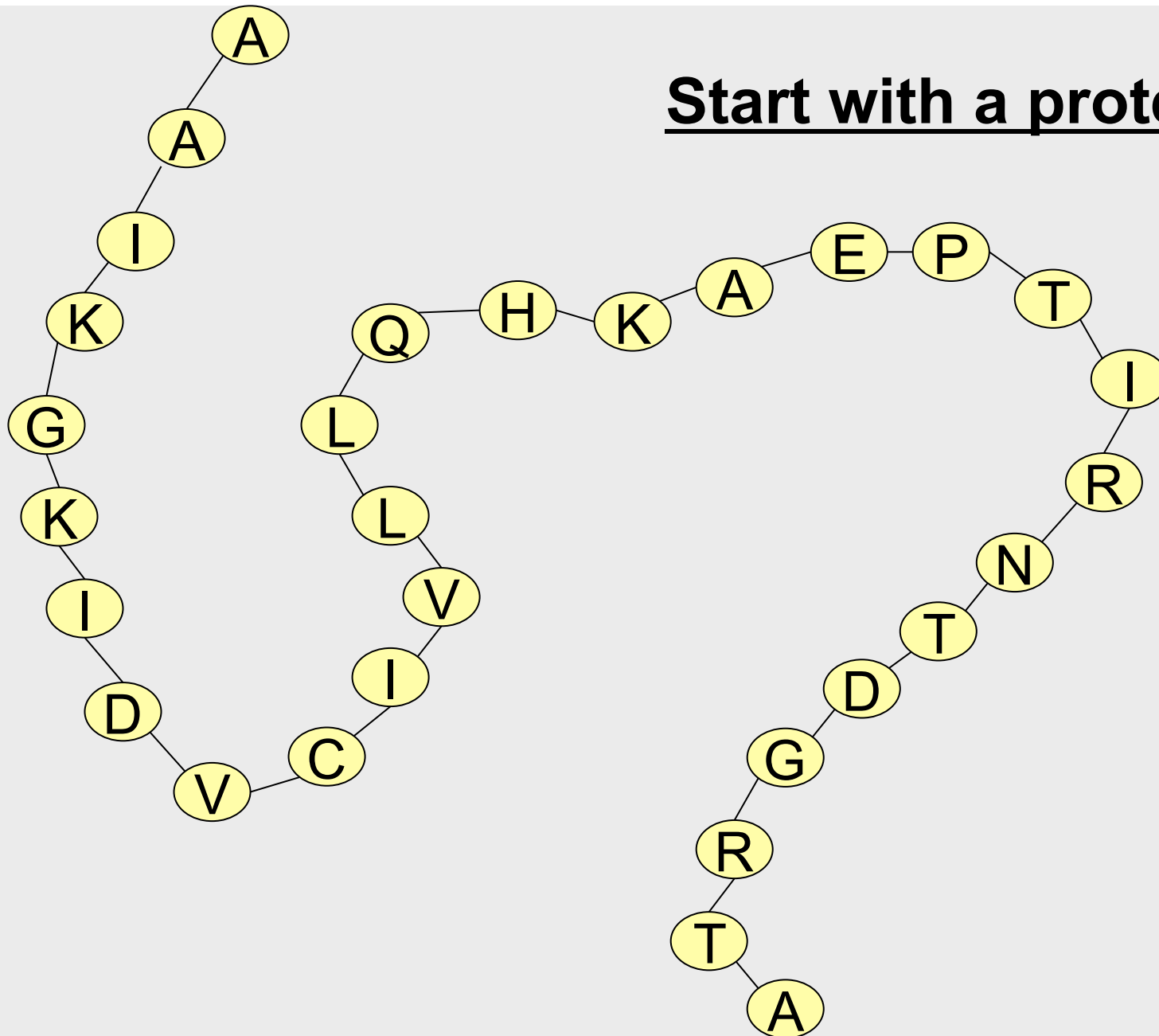
BSPR/EBI Educational Workshop
Hinxton, 16 July 2010



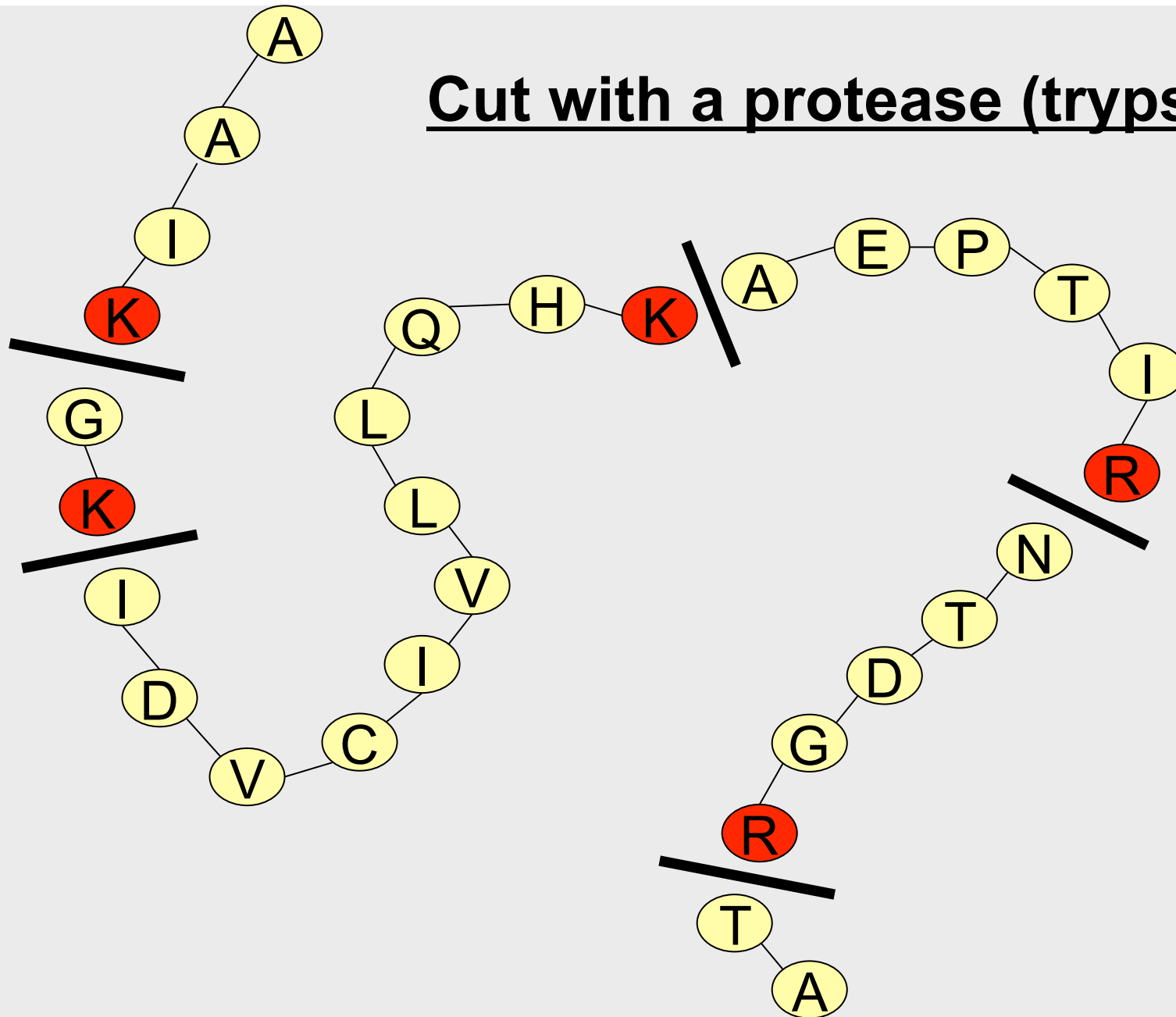
MS proteomics: overall workflow



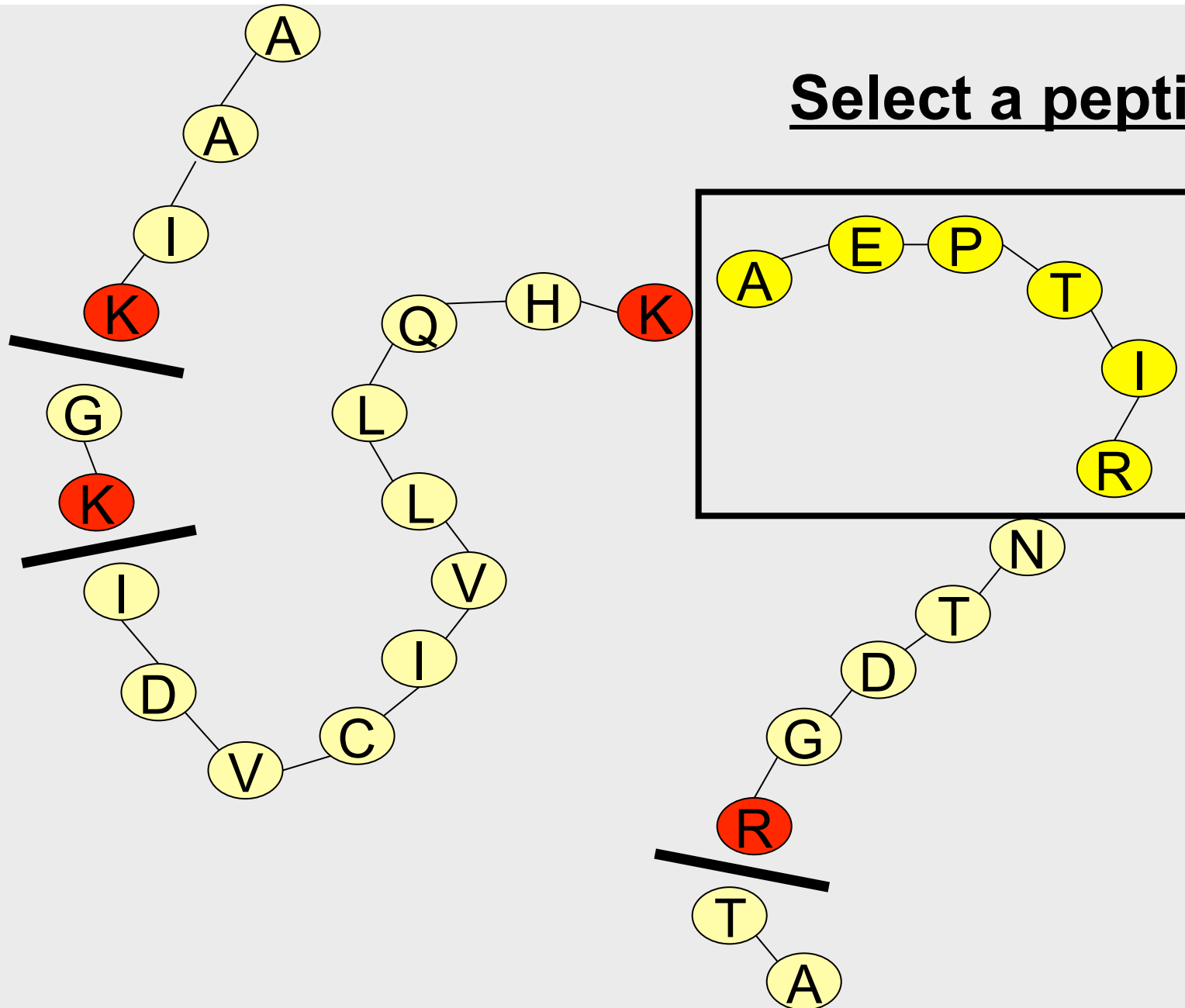
Start with a protein



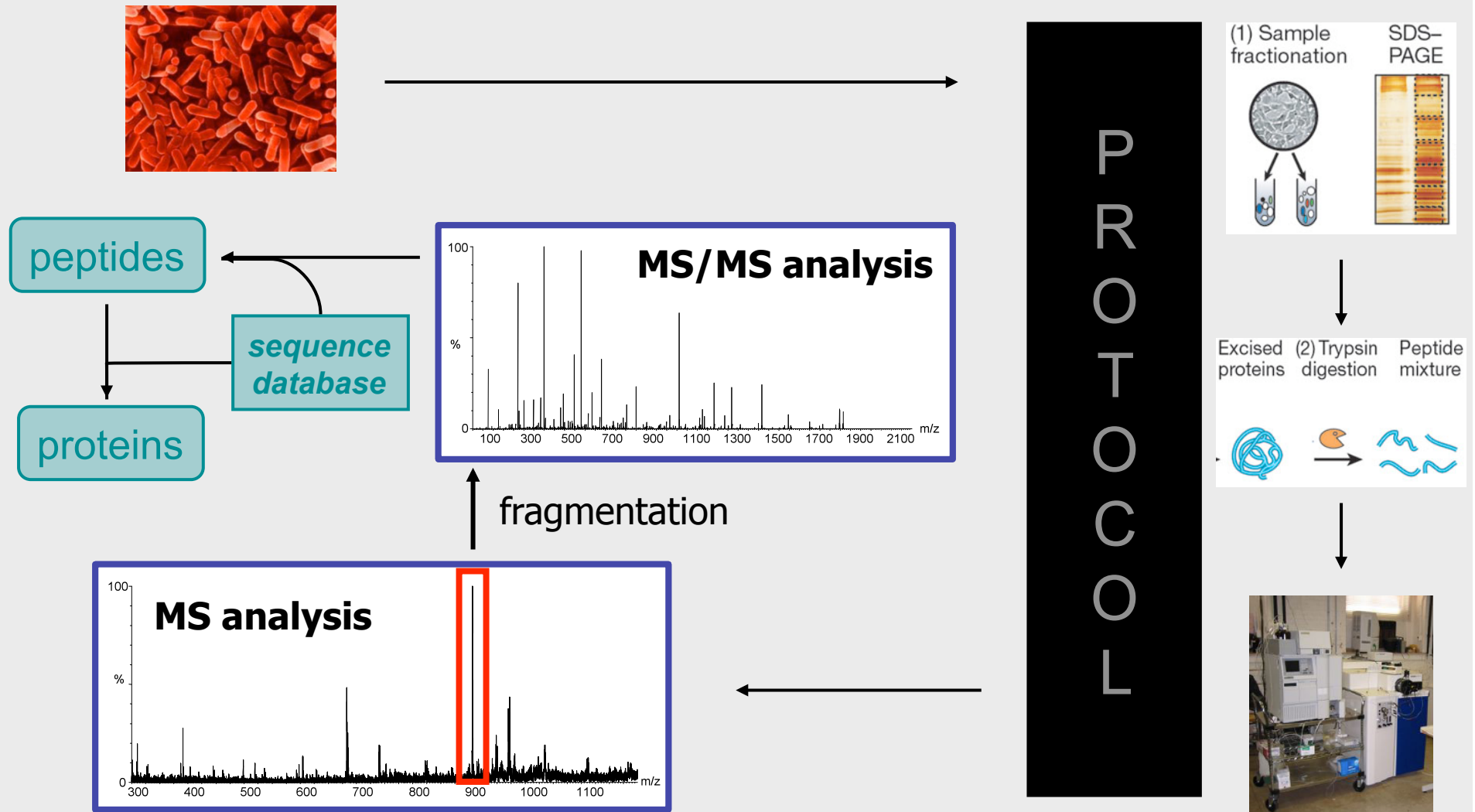
Cut with a protease (trypsin)



Select a peptide



MS proteomics: overall workflow



THE RATIONALE BEHIND SHARING PROTEOMICS DATA

Juan A. Vizcaíno
juan@ebi.ac.uk



BSPR/EBI Educational Workshop
Hinxton, 16 July 2010



Need of data sharing in the proteomics field

EDITORIAL

nature
biotechnology

Credit where credit is overdue

A universal tagging system that links data sets with the author(s) that generated them is essential to promote data sharing within the proteomics and other research communities.

Science progresses most rapidly when researchers provide access to their data. This is not only good scientific practice. It facilitates the confirmation of original results. It provides others with a starting point to explore new or related hypotheses. It speeds the identification of errors and discourages fraud. And it minimizes inefficient use of funding in duplicating experiments. And yet, full data disclosure in proteomics, and many other fields, remains a work in progress. If practicing scientists are to be truly incentivized to spend time and effort on sharing data, funders and publishers need to develop a universally recognized tagging system that would link investigators to their deposited data. In this way, publicly disclosed data sets would become part of a researcher's publication record, allowing such efforts to be recognized by employers and funders alike.

Next month marks the two-year anniversary of the publication of guidelines specifying the minimum reporting requirements for papers describing proteomics and molecular interaction experiments (*Nat. Biotechnol.* 25, 887–893, 894–898, 2007). Both sets of standards encourage deposition of data in public repositories, a practice that at the time was not universally adopted in proteomics.

We have carried out an informal survey of all manuscripts published in the year following publication of the two guidelines by the 68 authors of those two papers. The analysis reveals that a majority of the guideline authors published at least one manuscript last year for which no accompanying data were archived. If the proponents of data-reporting guidelines—most of whom are better resourced than other researchers in the field—are not depositing all of their data in a public repository, it is unlikely that the wider community is doing so either.

One issue that inhibits openness is the perception that full data disclosure may result in the loss of an edge over competing research groups. Occasionally, data are withheld while intellectual property is secured. More often, though, a failure to share simply reflects the considerable time and effort associated with formatting, documenting, annotating and releasing data. In this regard, the availability of new tools, such as an application (p. 598) to facilitate deposition of data in PRIDE (a public archive for mass spectrometry and protein identification data) should prove helpful.

For proteomics, the rapidly evolving technology and the complexity of the data itself pose particular challenges. Concerns about the quality of proteomics data generated by mass spectrometry have long plagued the field, raising the issue of whether peers have sufficient faith in other groups' work to not only value the data lodged in public repositories but also make the effort to deposit their own. Here too, though, progress is being made. A study reported in this issue (p. 633) demonstrates the high reproducibility of a targeted proteomic approach for biomarker discovery from plasma among several laboratories. Such a result would have been difficult to achieve using the technology and approaches of a few years ago.

But data quality is only part of the problem in overcoming the community's reticence about disclosure. For many researchers, the software provided by the public repositories for searching and analyzing proteomics data is not as efficient and user friendly as it could be. An analysis published last month by the Human Proteomics Organization cited the misassignment of peptides to ambiguously annotated proteins by database search engines as one of the major hindrances to researchers in the field (*Nat. Methods* 6, 423–430, 2009). What's more, despite the recent launch of yet another archive for mass spectrometry and protein identification data—the US National Center for Biotechnology Information's Peptidome repository (p. 600)—the various proteomics databases have yet to introduce a standardized data format that would allow the seamless exchange of data. Contrast this with the genome databanks, where the pooling of nucleotide sequence data in a common format has been integral to consistency, accessibility and, above all, utility of sequence data for reanalysis.

With all of these impediments, it's not surprising that proteomics researchers have been slow to embrace data disclosure. It is equally clear that disclosure edicts and recommendations from funding agencies and scientific journals have been insufficient to ensure widespread proteomics data release, despite evidence that the papers of researchers who share their data have an increased citation rate (*PLoS ONE* 2, e308, 2007). Clearly, other incentives are needed.

One option would be to provide researchers who release data to public repositories with a means of accreditation. This would take the form of a universally standardized tag for data that could be searched and recognized by both funding agencies and employers. An ability to search the literature for all online papers that used a particular data set would enable appropriate attribution for those who share. In essence, the tag would be a digital object identifier (DOI), currently best known for its use in unambiguously identifying papers online.

Similar to citation information about publications, citation information about a researcher's data DOIs could be gathered by funders assessing future support and used by institutions in performance evaluation. Researchers who disclose data sets that subsequently prove particularly useful to the community would end up with highly cited data DOIs, and could thereby be rewarded accordingly.

Such a system would not solve all the problems slowing data disclosure in proteomics and elsewhere. But it would provide greater incentive than the present system of evaluation, which is skewed almost exclusively to publications in high-profile journals and citation metrics. Data DOIs would not only enhance a researcher's reputation but also establish priority of data generation. Most important of all, they would provide a way to acknowledge the time and effort individuals must invest in sharing data, which ultimately benefits the scientific community as a whole.

© 2009 Nature America, Inc. All rights reserved.



NATURE BIOTECHNOLOGY VOLUME 27 NUMBER 7 JULY 2009

579

Juan A. Vizcaíno
juan@ebi.ac.uk



BSPR/EBI Educational Workshop
Hinxton, 16 July 2010



Proteomics data sharing: why?

- 1) Data producers are not always the best data analysts
Sharing of data allows analysts access to real data, and in turn allows better analysis tools to be developed
- 2) Meta-analysis of data can recycle previous findings for new tasks
Putting findings in the context of other findings increases their scope
- 3) Sharing data allows independent review of the findings
When actual replication of an experiment is often impossible, a re-analysis or spot checks on the obtained data become vitally important
- 4) Direct benefit for the field: fragmentation models, spectral libraries, ...

Simply sharing data is not enough...

| Table 1. Identities of stress-induced proteins | | | Table 1. Identification of exosomal proteins based on MALDI-TOF peptide mass fingerprinting or MS/MS-derived sequences | | | | | | Phase | Stress | (kDa)/pI | Accession no. | Species | | | |
|--|---------|------------------|--|---|------------------------------------|-------------------------------|----------------------|-------------------|-------|--------|----------|---------------|---------|-----------------------|--------------|-------------|
| Spot ID | Synonym | Function | Band (Fig. 1) | Protein Name | Identification Method ^a | Accession Number ^b | Molecular Mass (kDa) | Matching Peptides | | | | | | Sequence Coverage (%) | Experimental | Theoretical |
| 1202 | SCO0525 | Hypothetical p | A | Mac-1 α -chain = CD11b | | Théry et al., 1999 (14) | | | | | | | | | | |
| 3307 | SCO2988 | UDP-glucose | 1 | Complement C3 ^c | MS/MS (7) | 4093220 | | | | | | | | | | |
| 3509 | SCO2180 | Putative dihydro | 1 | PK-120 ^d | MS/MS (3) | Not in databases | | | | | | | | | | |
| 6413 | SCO6027 | Probable acet | 1 | α 2-Macroglobulin ^e | MS/MS (2) | Not in databases | | | | | | | | | | |
| 6419 | SCO1494 | 3-Dehydroquini | 2 | Plasminogen ^f | MS | P06868 | 91 | 28 | 37 | | | | | | | |
| 6823 | SCO5477 | Putative oligo | 3 | Alix | MS | 6755002 ^d | 96 | 26 | 34 | | | | | | | |
| 118 | SCO1340 | Conserved hy | 3 | Mac-1 β -chain = CD18 | MS/MS (6) | P11835 | 85 | 27 | 38 | | | | | | | |
| 1104 | SCO2368 | Conserved hy | 4 | hsp90- β = hsp84 | MS/MS (1) | P11499 | 83 | 30 | 38 | | | | | | | |
| 1617 | SCO5373 | ATP synthase | 5 | Serum albumin ^g | MS | P02769 | 69 | 42 | 66 | | | | | | | |
| 2601 | SCO5373 | ATP synthase | B | | MS/MS (3) | | | | | | | | | | | |
| 3616 | SCO5371 | ATP synthase | B, C | hsc73 | | Théry et al., 1999 (14) | | | | | | | | | | |
| 5721 | SCO4814 | Bifunctional p | 6 | MFG-E8/lactadherin | MS | P05218 | 50 | 20 | 44 | | | | | | | |
| 1515 | SCO2180 | Putative dihydro | 7 | Tubulin β | MS | Q07076 | 50 | 6 | 13 | | | | | | | |
| 1616 | SCO3661 | Putative chap | 7 | Annexin VII = synexin | MS/MS (3) | | | | | | | | | | | |
| 2706 | SCO3671 | Heat shock p | 7 | Bovine coagulation factor X ^c | MS/MS (2) | P00743 | 54 | | | | | | | | | |
| 2906 | SCO5999 | Aconitase | 7 | Annexin I | MS/MS (3) | Q95121 | 46 | | | | | | | | | |
| 3504 | SCO1936 | Putative trans | 7 | Tumor susceptibility protein (tsf) 101 | MS/MS (2) | 3184260 ^d | 44 | | | | | | | | | |
| 5310 | SCO0506 | NH(3)-depend | 7 | Rab GDP dissociation inhibitor (GDI) 3 | MS | Q61598 ^e | 51 | 10 | 21 | | | | | | | |
| 7417 | SCO5477 | Putative oligo | 7 | Elongation factor (EF) 1- α -1 | MS/MS (2) | P10126 | 50 | | | | | | | | | |
| 505 | SCO1998 | 30S ribosoma | 7 | EIF-4A-II | MS/MS (2) | P10630 | 46 | 7 | 25 | | | | | | | |
| 1711 | SCO1352 | Xaa-pro amin | 8 | Annexin I | MS | P10107 | 39 | | | | | | | | | |
| 2618 | SCO0681 | Putative ferred | 8 | Reverse transcriptase/pol (murine leukemia virus) | MS/MS (2) | 61790 ^d | | | | | | | | | | |
| 2722 | SCO1998 | 30S ribosoma | D | γ -Actin | MS/MS (1) | | | | | | | | | | | |
| 4407 | SCO5113 | Oligopeptide | E | G protein G _{2e} subunit | MS/MS (2) | | | | | | | | | | | |
| 4509 | SCO2390 | Beta-ketoacyl | F | Annexin II | MS | P48036 | 36 | 16 | 54 | | | | | | | |
| 1803 | SCO2181 | 2 Oxoglutarat | F | Annexin V | MS | | | | | | | | | | | |
| 2113 | SCO4277 | Hypothetical p | 9 | Annexin IV | MS/MS (4) | P97429 | 36 | 20 | 63 | | | | | | | |
| 3101 | SCO3899 | Hypothetical p | 10 | Galactin-3 = Mac-2 | MS/MS (4) | P16110 | 27 | 11 | 37 | | | | | | | |
| 4309 | SCO1081 | Putative elect | 11 | Syntenin | MS/MS (6) | 2197106 ^d | 32 | 17 | 35 | | | | | | | |
| 4512 | SCO5212 | 3-Phosphosh | G | Gag polyprotein (murine leukemia virus) | MS/MS (6) | | | | | | | | | | | |
| 5514 | SCO3629 | Putative aden | G | MHC class II β -chain 14-3-3 protein η | MS | P11576 | 28 | 21 | 68 | | | | | | | |
| | | | 12 | 14-3-3 protein η | MS/MS (4) | P35215 | 28 | 20 | 63 | | | | | | | |
| | | | 12 | 14-3-3 protein γ | MS/MS (2) | 3065929 ^d | | | | | | | | | | |
| | | | 13 | Apolipoprotein A-I ^c | MS/MS (1) | P15497 | 30 | 25 | 67 | | | | | | | |
| | | | H | CD9 | MS | | | | | | | | | | | |
| | | | 14 | Thioredoxin peroxidase II | MS | P35700 | 22 | 8 | 43 | | | | | | | |
| | | | 14 | Rab 11 | MS/MS (6) | | | | | | | | | | | |
| | | | 14 | κ -Casein ^g | MS/MS (1) | P46638 | 24 | | | | | | | | | |
| | | | 14 | Rab-7 | MS/MS (2) | P02668 | 21 | | | | | | | | | |
| | | | 15 | Ferritin light chain ^f | MS | P51150 | 24 | 5 | 26 | | | | | | | |
| | | | 16 | Rap1B | MS/MS (3) | | | | | | | | | | | |
| | | | 16 | Cofilin | MS | O46415 | 20 | 15 | 73 | | | | | | | |
| | | | 17 | Histone H3 | MS | P09526 | 21 | 14 | 57 | | | | | | | |
| | | | 18 | Histone H2B | MS | P18760 | 19 | 10 | 50 | | | | | | | |
| | | | 19 | Histone H2A | MS | Z85979 ^f | 15 | 7 | 45 | | | | | | | |
| | | | 19 | Histone H4 | MS | P10853 | 14 | 13 | 82 | | | | | | | |
| | | | 20 | Profilin I | MS | P20670 | 14 | 12 | 67 | | | | | | | |
| | | | 20 | Hemoglobin γ -chain ^g | MS | 90626 ^d | 11 | 15 | 90 | | | | | | | |
| | | | 21 | Hemoglobin α -chain ^g | MS | P10924 | 15 | 11 | 60 | | | | | | | |
| | | | 21 | | MS | P02081 | 16 | 16 | 74 | | | | | | | |
| | | | 21 | | MS | P01966 | 15 | 9 | 66 | | | | | | | |

+ SOD1

A nuance: *available* data vs. *accessible* data

When data is only made available as arbitrarily formatted tables, it carries important limitations

- Source data are not made available
 - No peer review validation possible
 - Very little raw materials for testing innovative *in silico* techniques are available
 - Traceability of data is lost quickly in downstream results
- Automated (re-)processing of the results (e.g., identifications) is impossible
- Data producers do not actually feed their results and knowledge back to the community

Community standards for proteomics



The Human Proteome Organisation (HUPO)
Proteomics Standards Initiative (PSI)



<http://www.psidev.info>

- Creates minimal requirements, standard formats, and CV's and ontologies
- Composed of several workgroups

| | | |
|-------------------------------|--------------|----------------------------|
| <i>Molecular Interactions</i> | <i>(MI)</i> | <i>PSI-MI format v2.5</i> |
| <i>Mass Spectrometry</i> | <i>(MS)</i> | <i>mzData, mzML format</i> |
| <i>Protein Separation</i> | <i>(PS)</i> | <i>GelML format</i> |
| <i>Proteomics Informatics</i> | <i>(PI)</i> | <i>mzIdentML format</i> |
| <i>Protein Modifications</i> | <i>(Mod)</i> | <i>PSI-MOD ontology</i> |

How do we make this all happen?

- **Journal guidelines**

Journal guidelines heavily influence the decisions taken by authors; by first requesting and subsequently mandating data submission to established repositories, they provide an important stick.

- **Funder support and guidelines**

Funders contribute both sticks and carrots. The sticks lie in the grant application guidelines; they can require a plan for data management and dissemination. The carrot is in providing specific funding for this aspect of science.

- **Data repositories**

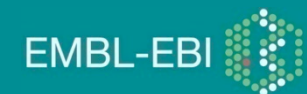
The availability of reliable, freely available repositories is key; submission thresholds should be kept low and **added value** needs to be provided. Furthermore, feedback loops need to be established in order to ensure that accumulated data flows back to the user community. Repositories thus provide mostly carrots.

PROTEOMICS DATA REPOSITORIES AVAILABLE TODAY

Juan A. Vizcaíno
juan@ebi.ac.uk



BSPR/EBI Educational Workshop
Hinxton, 16 July 2010



Existing proteomics repositories

- Main public repositories:

- PROteomics IDentifications database (PRIDE)
- Global Proteome Machine (GPMDB)
- Peptide Atlas
- Tranche
- NCBI Peptidome

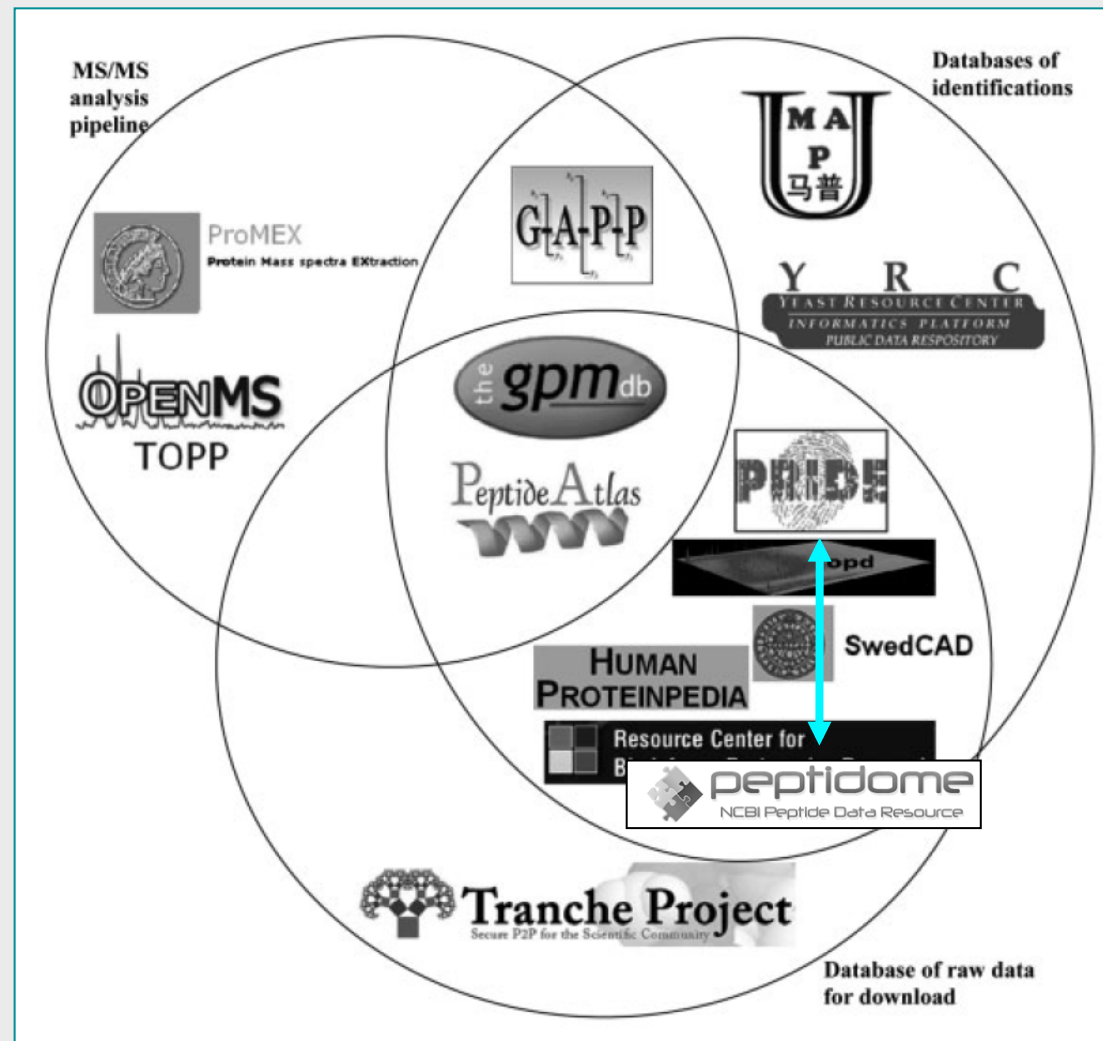


- Smaller scale repositories, more specialized:

Among others: Human Proteinpedia, Genome Annotation Proteomics Pipeline (GAPP), MAPU, SwedCAD, PepSeeker, Open Proteomics Database, ...

- Very diverse: different aims, functionalities, ...

A comprehensive view on existing systems

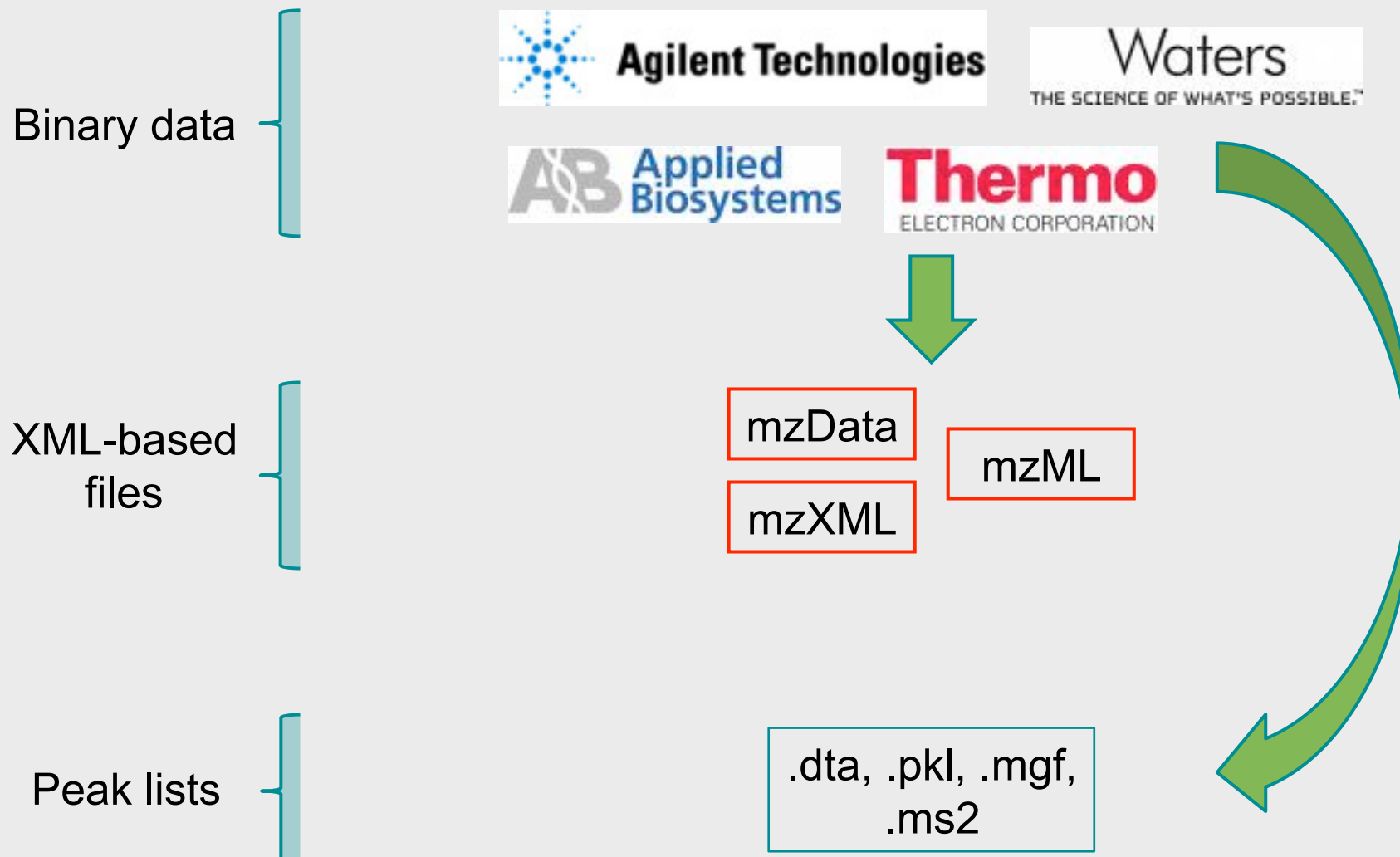


From: Mead *et al.*, *Proteomics*, 2009

Types of information stored

- 1) **Original experimental data** recorded by the mass spectrometer (primary data)

Primary data



Types of information stored

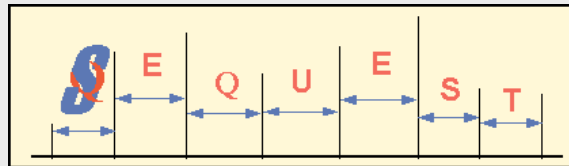
- 1) **Original experimental data** recorded by the mass spectrometer (primary data)
- 2) **Identification results** inferred from the original primary data

Peptide and Protein Identifications

GENEBIO
PHENYX

{MATRIX}
{SCIENCE}

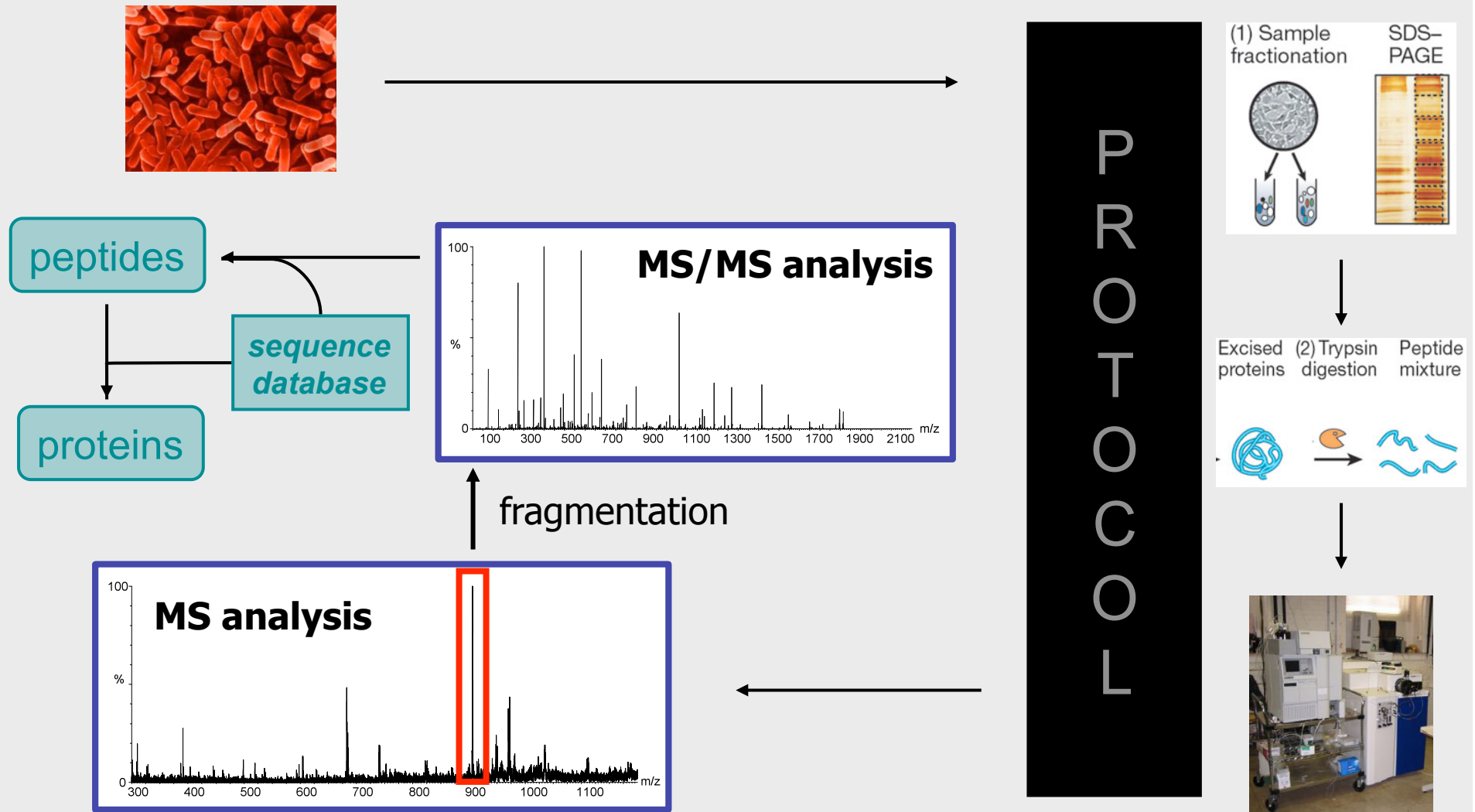
ProteinPilot[™]
Software 2.0



mzIdentML,
mascot .dat,
sequest .out,
SpectrumMill .spo
pep.xml, prot.xml

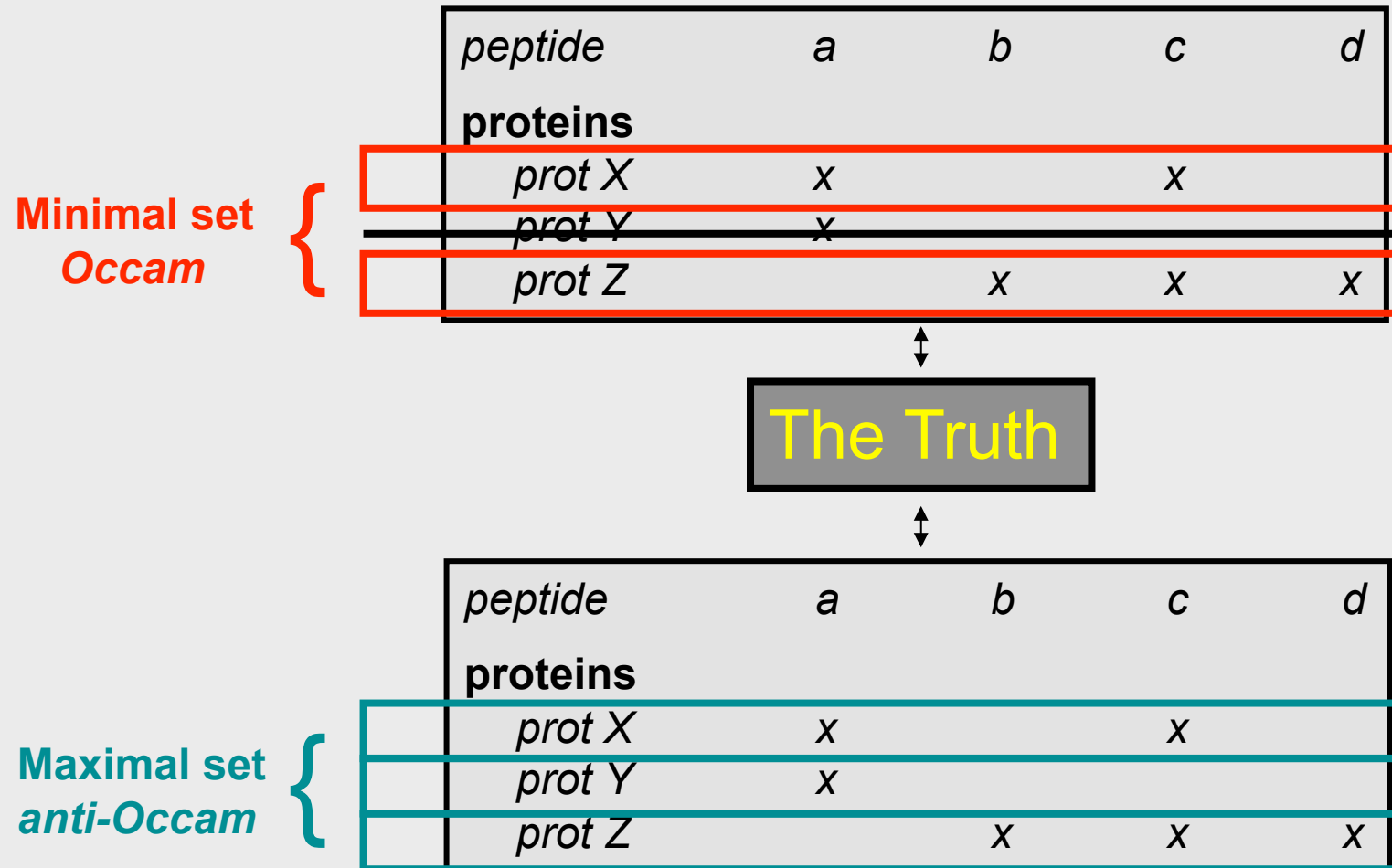
Only qualitative data!

MS proteomics: overall workflow

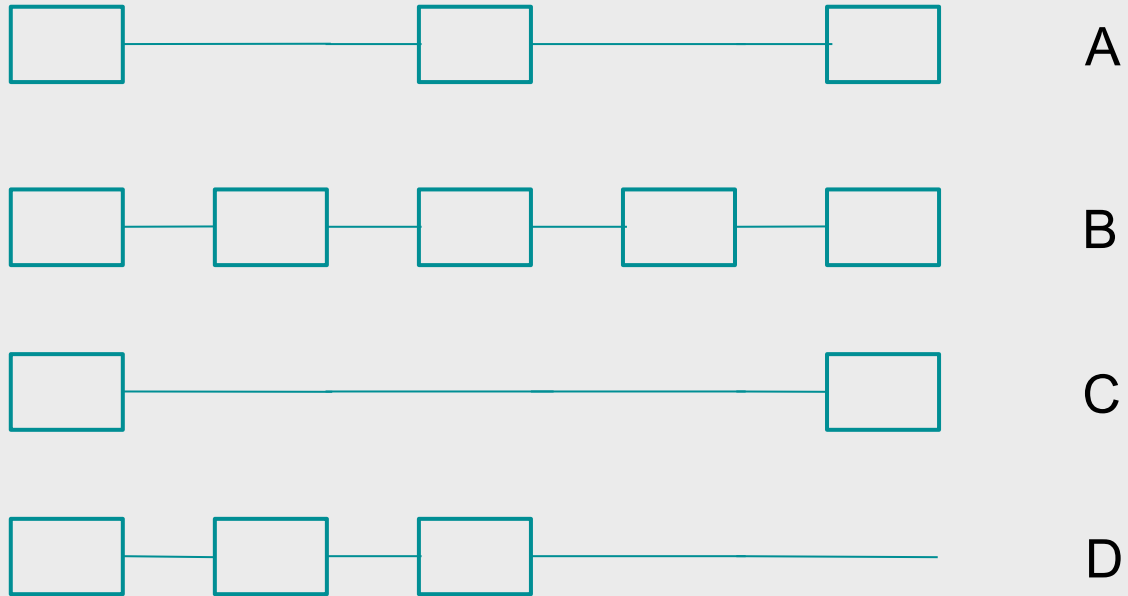


Intermezzo: Protein inference

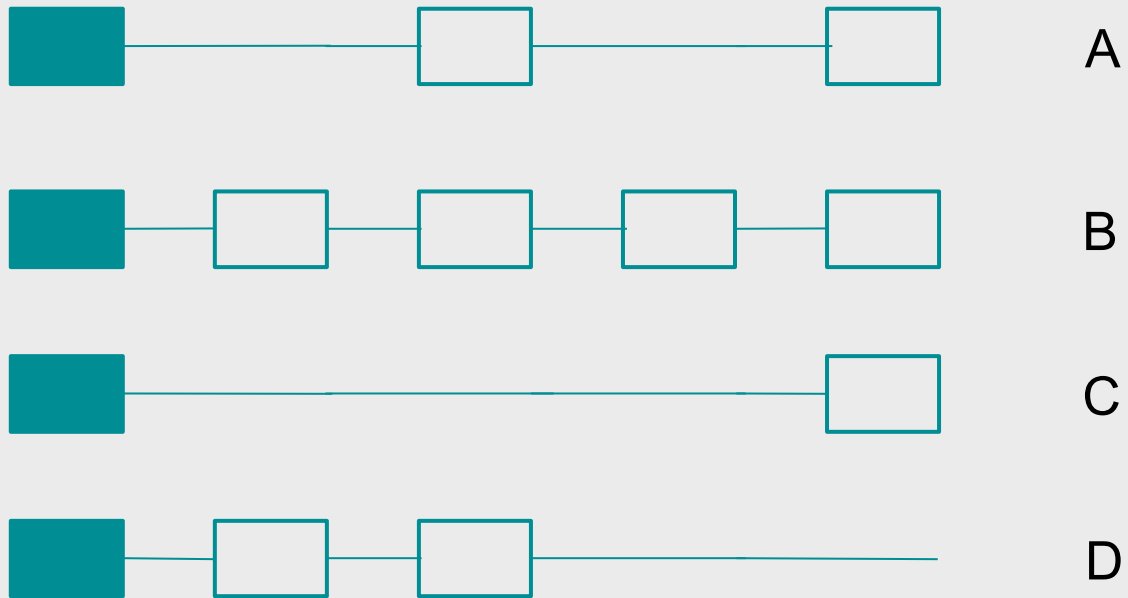
The *minimal* and *maximal* explanatory sets



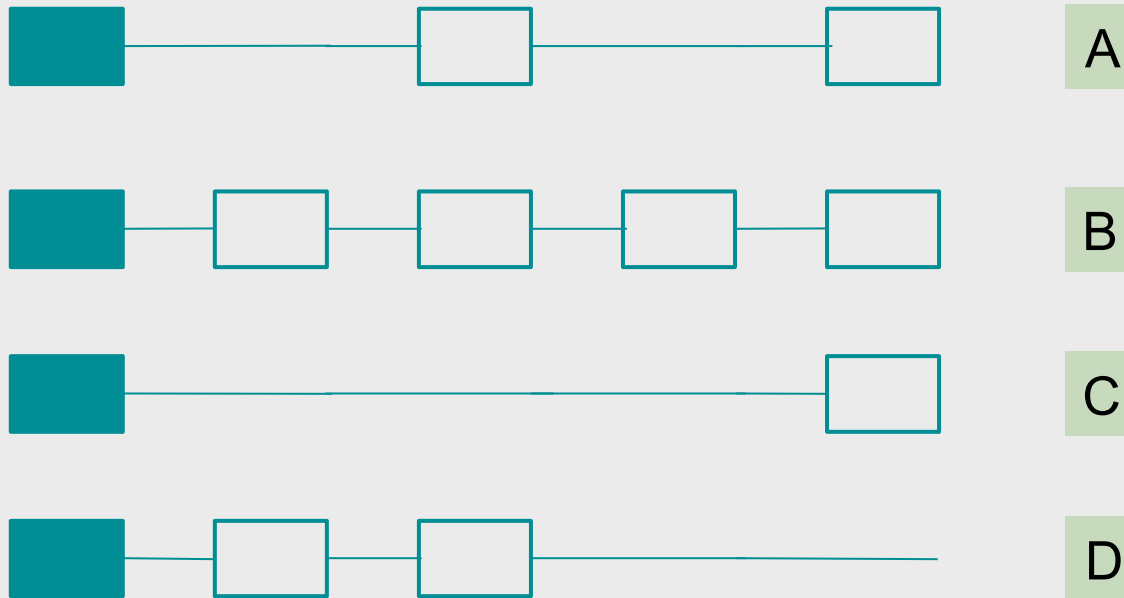
An additional layer of complexity...



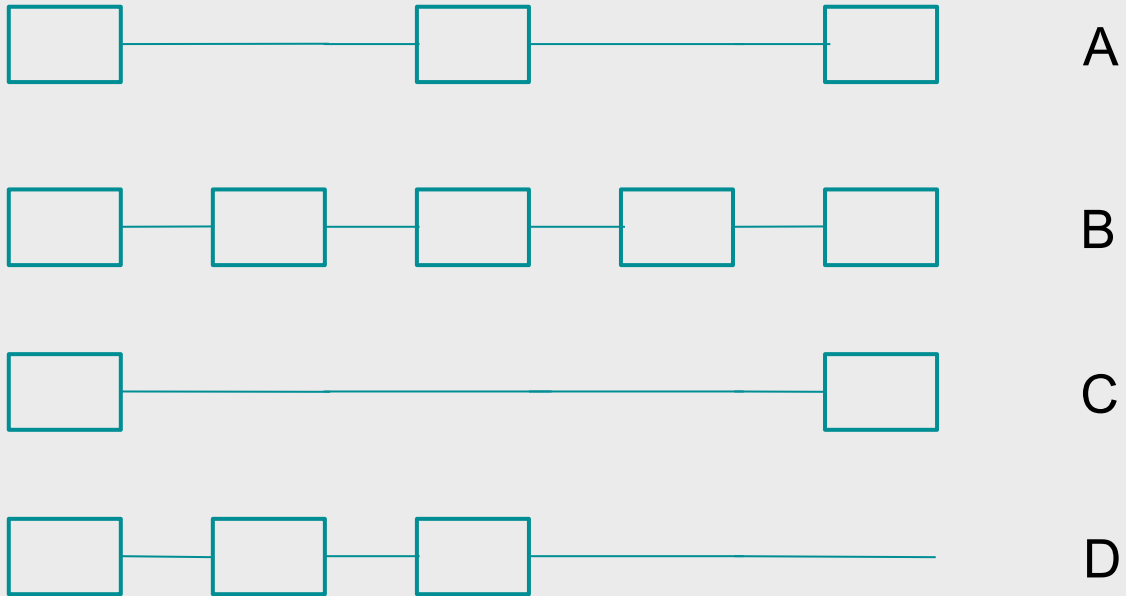
Protein inference



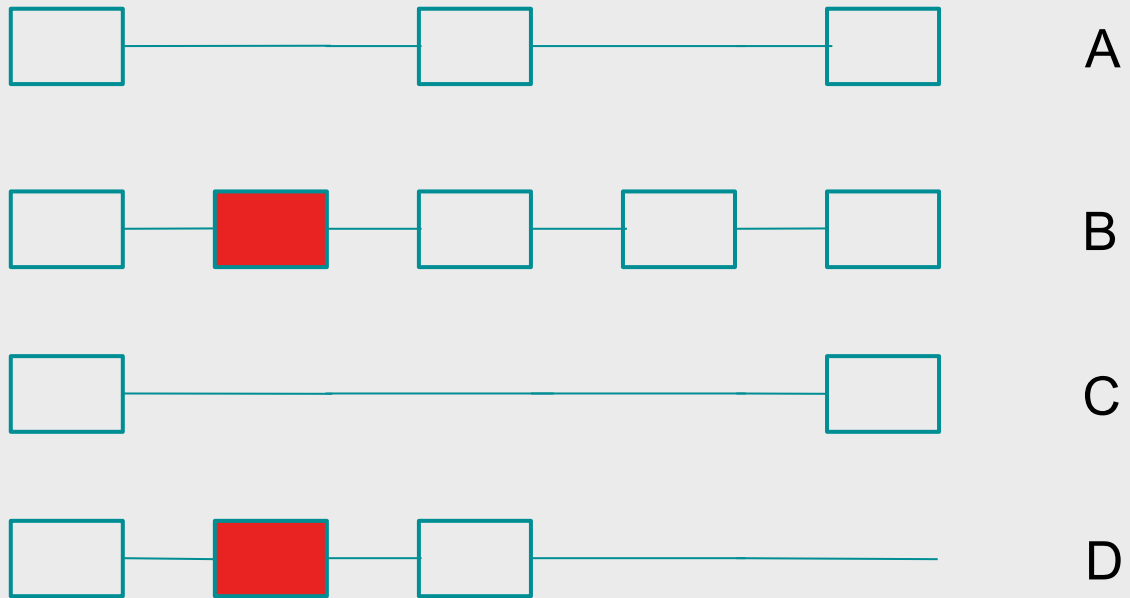
Protein inference



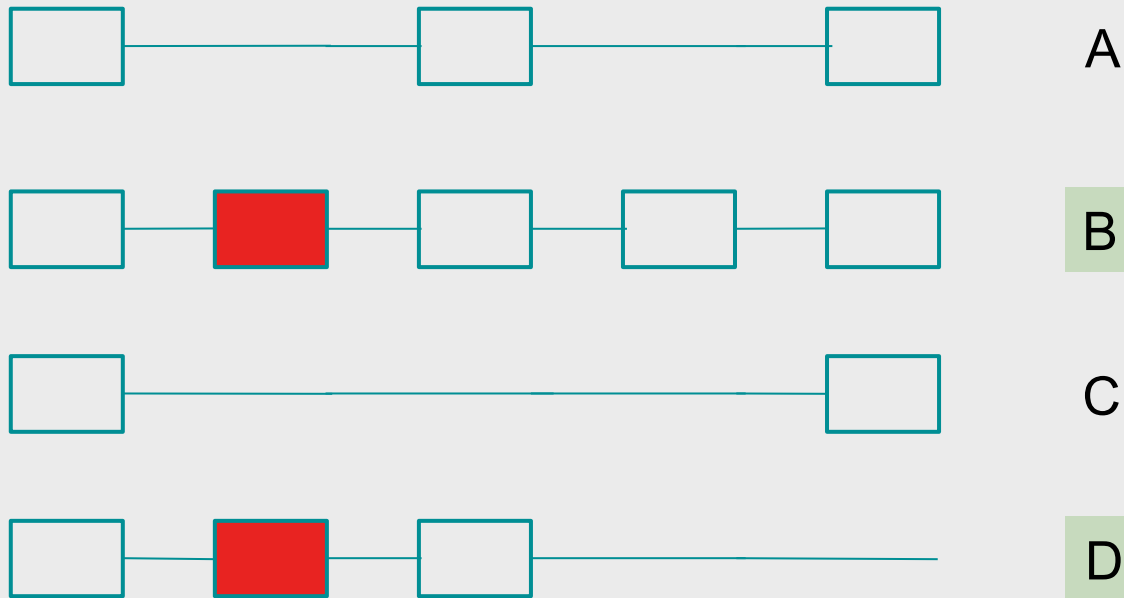
Protein inference



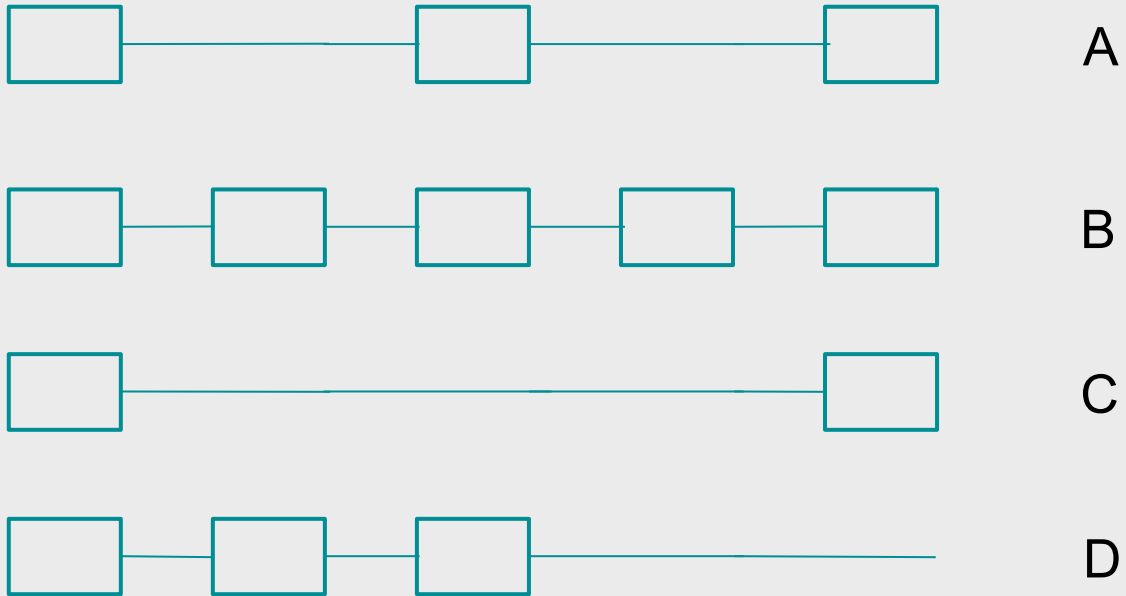
Protein inference



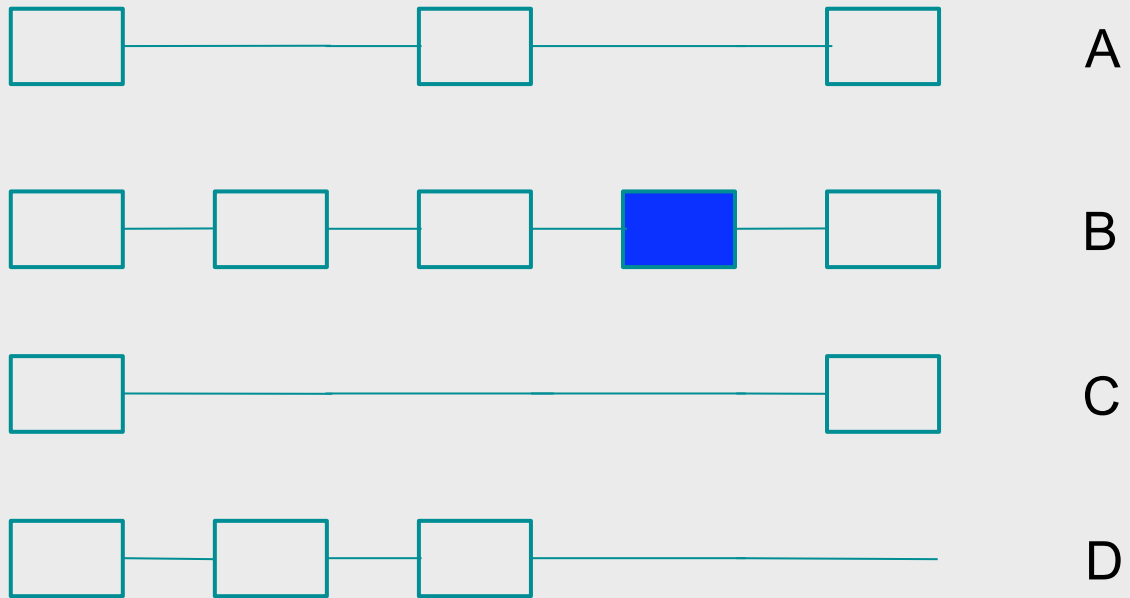
Protein inference



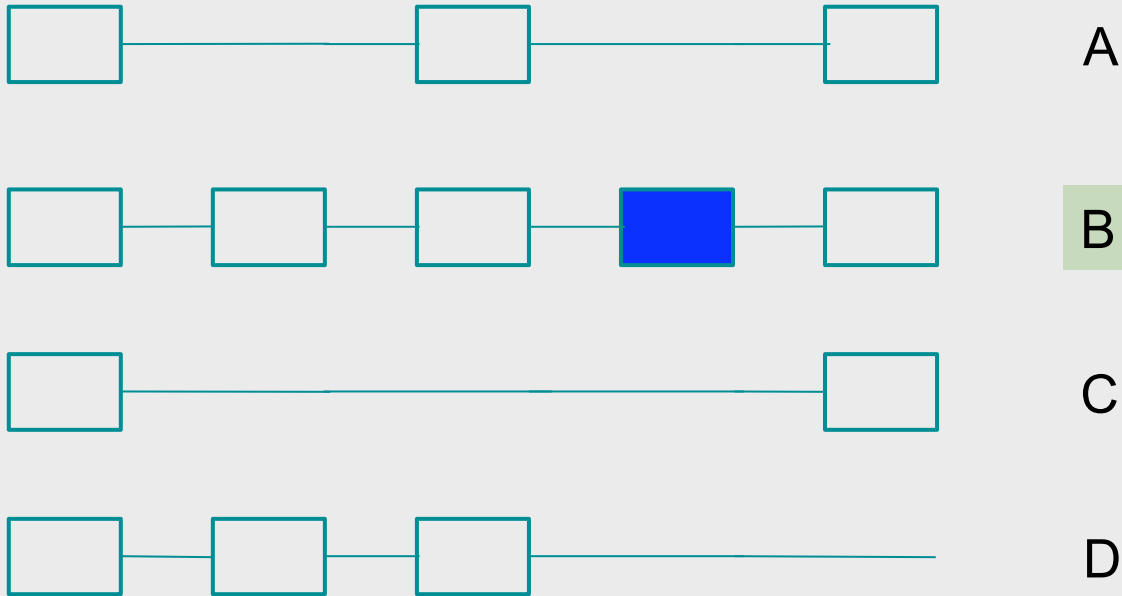
Protein inference



Protein inference



Protein inference



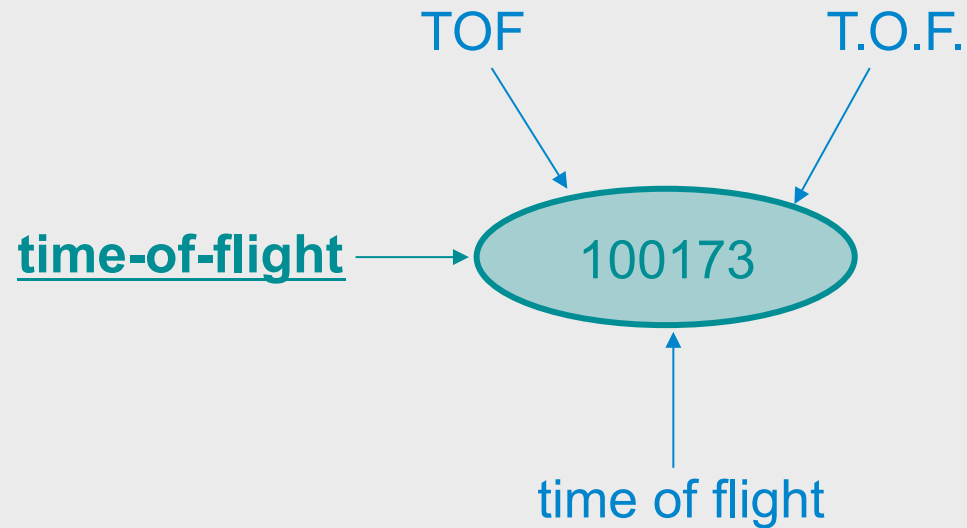
Unambiguous
peptide

Types of information stored

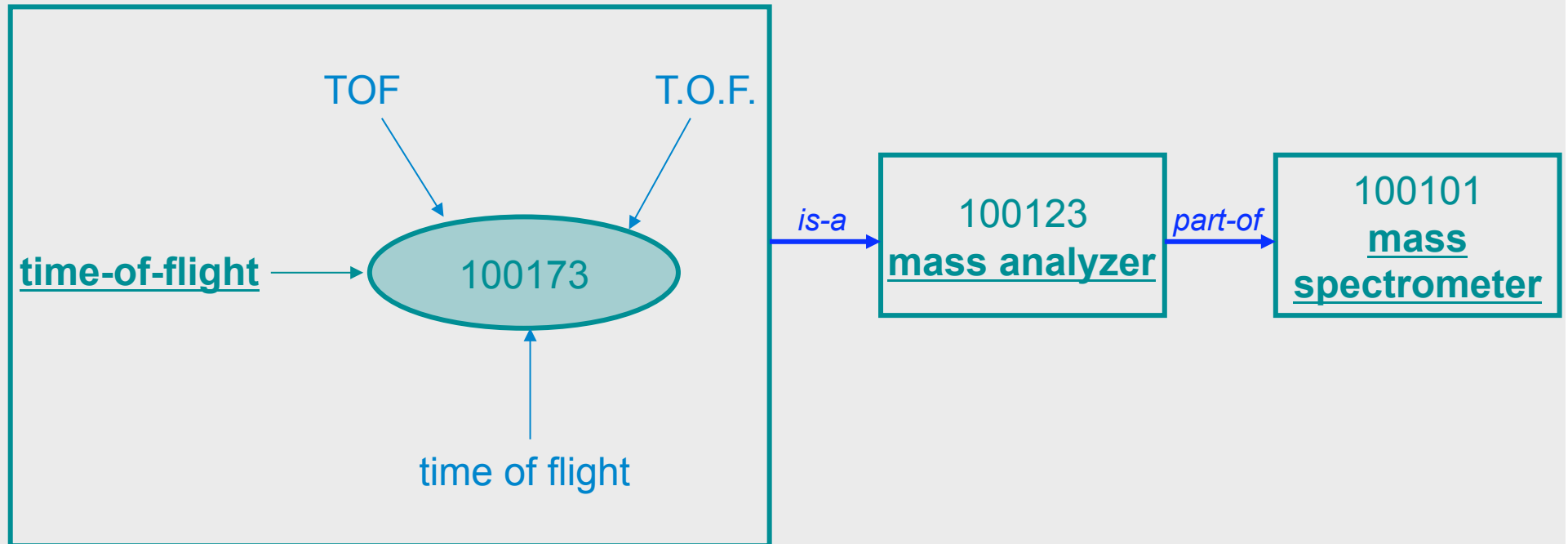
- 1) **Original experimental data** recorded by the mass spectrometer (primary data)
- 2) **Identification results** inferred from the original primary data
- 3) Experimental and technical **metadata**

Controlled Vocabularies (CVs)

Term
Synonyms



Relationships between CV terms



CVs, ontologies (here: PSI-MOD)

<http://www.ebi.ac.uk/ols>

OLS - Ontology Lookup Service

MOD Ontology Browser

protein modification

- uncategorized protein modification
- protein modification categorized by isobaric sets
- protein modification categorized by chemical process
 - alkylated residue
 - crosslinked residues
 - deamidated residue
 - cyclized residue
 - acylated residue
 - acylated residue
 - N-acetylated residue
 - N-acetyl-L-alanine
 - N-acetyl-L-aspartic acid
 - N-acetyl-L-cysteine
 - N-acetyl-L-glutamic acid
 - N-acetyl-L-glutamine
 - N-acetyl-glycine
 - N-acetyl-L-isoleucine
 - N-acetyl-L-methionine
 - N-acetyl-L-proline
 - N-acetyl-L-serine
 - N-acetyl-L-threonine
 - N-acetyl-L-tyrosine
 - N-acetyl-L-valine
 - N2-acetyl-L-arginine
 - N,O-diacetylated L-serine
 - N-acetylated L-lysine
 - N2-acetyl-L-lysine

Legend:

- is a
- develops from
- part of
- other

Help (hide)

Double-click a term to see its children. The ontology browser is populated dynamically. If the children for a given term, there may be a small delay while the browser fetches. Click to high see any information associated with it. Hover over a term to see its relation with its immediate terms will not display any relational information.

Relations

N-acetyl-L-serine is_a N-acetylated residue

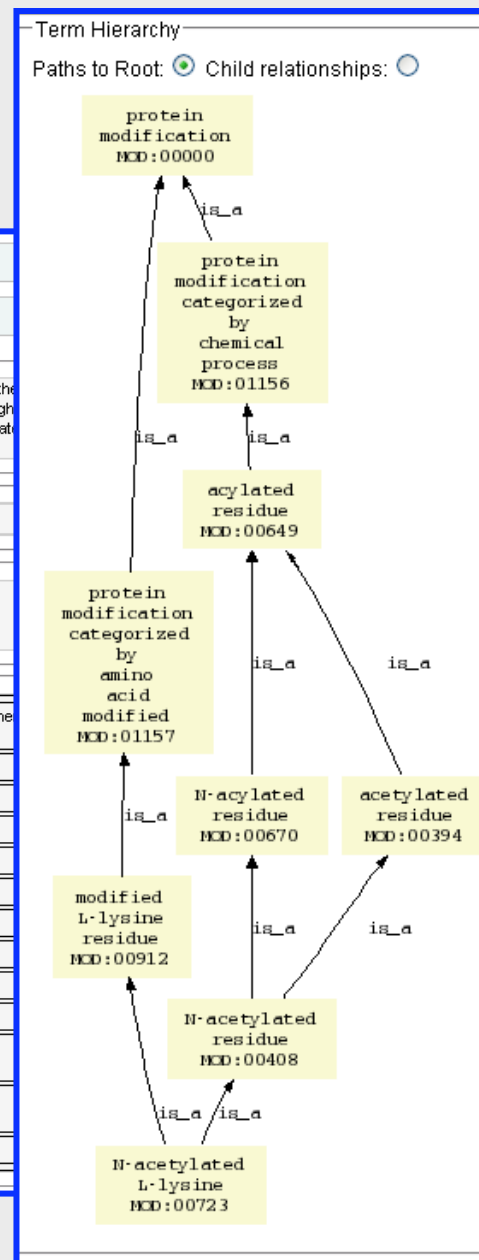
Term Information

ID: MOD:00723

Name: N-acetylated L-lysine

Associated information

| | |
|-----------------|---|
| definition | A protein modification that effectively converts an L-lysine residue to either L-lysine, or N6-acetyl-L-lysine. |
| DiffAvg | 42.04 |
| DiffFormula | C 2 H 2 N 0 O 1 |
| DiffMono | 42.010565 |
| Formula | C 8 H 15 N 2 O 2 |
| MassAvg | 171.22 |
| MassMono | 171.113353 |
| Origin | K |
| Source | Natural |
| TermSpec | none |
| preferred name | N-acetylated L-lysine |
| exact synonym | AcLys |
| xref_definition | PSI-MOD:ref |



Types of information stored

- 1) **Original experimental data** recorded by the mass spectrometer (primary data)
- 2) **Identification results** inferred from the original primary data
- 3) Experimental and technical **metadata**
- 4) **Quantitation** information

Wide variety of quantitative techniques...

Quantitation: Overview

Many different approaches to protein quantitation using mass spectrometry data have been described in the literature. For a short, recent review, see [Ong, S. E. and Mann, M., Mass spectrometry-based proteomics turns quantitative, Nature Chemical Biology 1 252-262 \(2005\)](#). In terms of the "mechanics" of their implementation, most of the popular approaches can be classified into a relatively small number of **protocols**:

- **Reporter**: Quantitation based on the relative intensities of fragment peaks at fixed m/z values within an MS/MS spectrum. For example, [ITRAQ](#) and [Tandem Mass Tags](#)
- **Precursor**: Quantitation based on the relative intensities of extracted ion chromatograms (XICs) for precursors within a single data set. This is by far the most widely used approach, which can be used with any chemistry that creates a precursor mass shift. For example, [¹⁸O](#), [AQUA](#), [ICAT](#), [ICPL](#), [Metabolic](#), [SILAC](#), etc., etc.
- **Multiplex**: Quantitation based on the relative intensities of sequence ion fragment peaks within an MS/MS spectrum. This is a [novel approach](#), which can be used with labels located at the peptide terminus, such as ¹⁸O or SILAC at K or R in combination with tryptic digestion.
- **Replicate**: Label free quantitation based on the relative intensities of extracted ion chromatograms (XICs) for precursors in multiple data sets aligned using mass and elution time.
- **emPAI**: Label free quantitation for the proteins in a mixture based on protein coverage by the peptide matches in a database search result.
- **Average**: Label free quantitation for the proteins in a mixture based on the application of a rule to the intensities of extracted ion chromatograms (XICs) for the peptide matches in a database search result.

Some protocols can be fully implemented within a Mascot result report because all the necessary information is present in the peak list. These protocols are [Reporter](#), [Multiplex](#), and [emPAI](#). In fact, emPAI is "always on", and will be reported whenever an MS/MS search contains at least 100 spectra.

The other three protocols require additional information from the raw data file, either because it is necessary to integrate the elution profile of each precursor peptide or because information is required for precursor peptides that were not used to trigger MS/MS scans, so are missing from the peak list. So, for [Precursor](#), [Replicate](#), and [Average](#), the quantitation report is generated in Mascot Distiller, which has access to both the Mascot search results and the raw data.

{*MATRIX*}
{*SCIENCE*}

Quantitation techniques



Label free

Gel-based quantitation
approaches



- Different philosophies
- Very heterogeneous data formats
- Techniques not very well established

Very problematic data for proteomics repositories

PRIDE

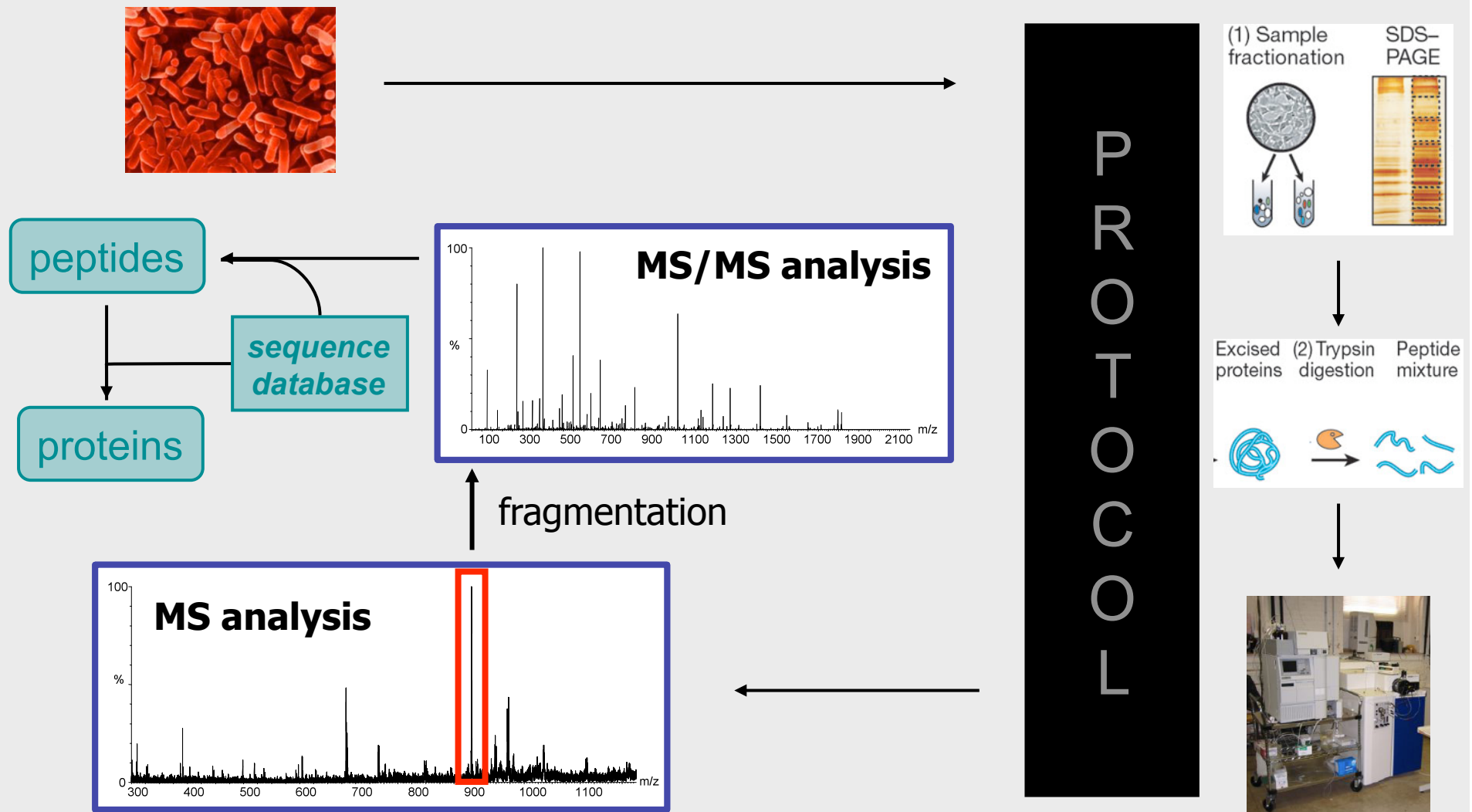
Juan A. Vizcaíno
juan@ebi.ac.uk



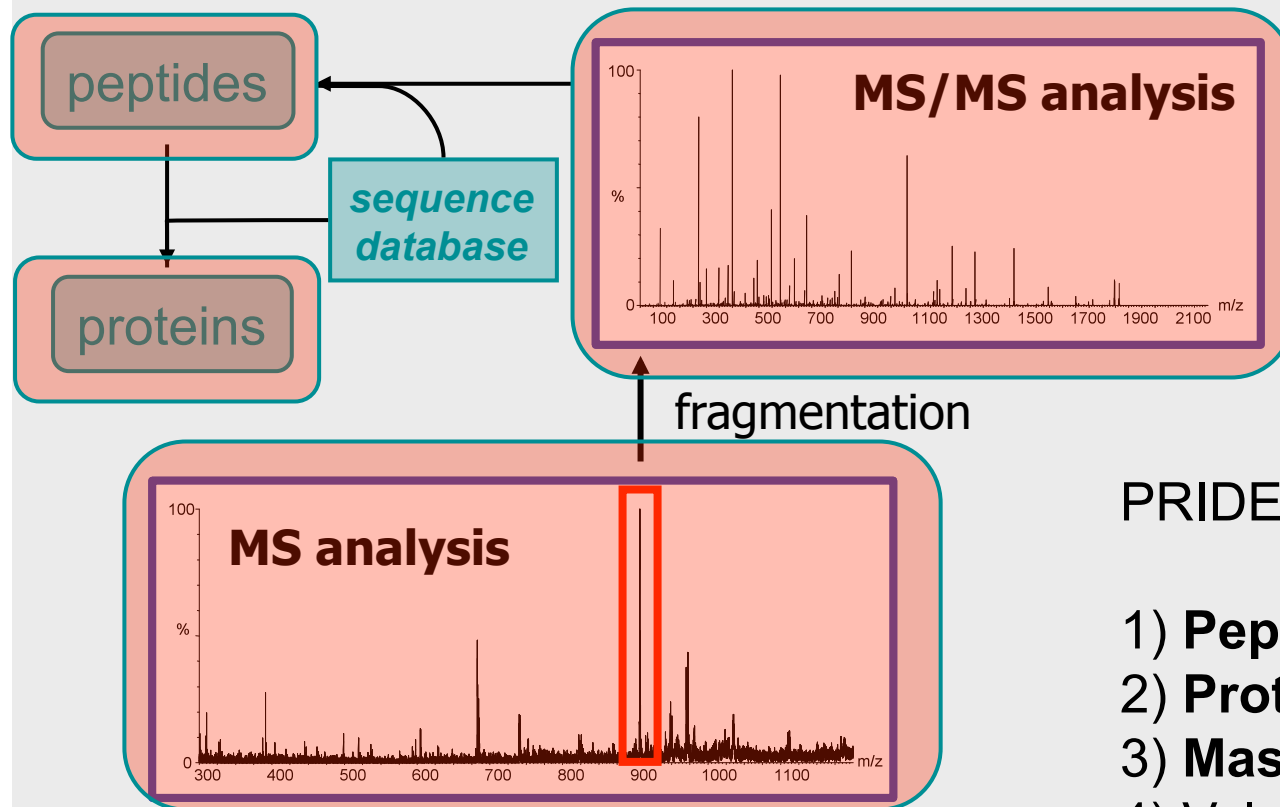
BSPR/EBI Educational Workshop
Hinxton, 16 July 2010



MS proteomics: overall workflow



PRIDE database (www.ebi.ac.uk/pride)



PRIDE stores:

- 1) **Peptide IDs**
- 2) **Protein IDs**
- 3) **Mass spectra** as peak lists
- 4) Valuable additional **metadata**

PRIDE: why is it there?

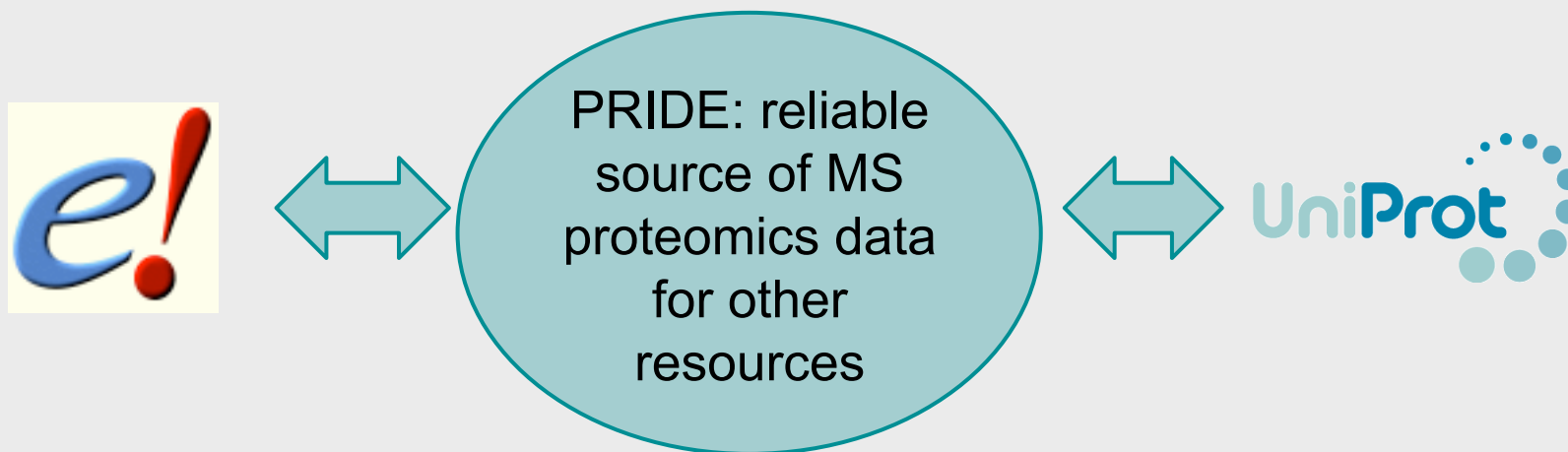


- Repository to support publications (proteomics MS derived data)

PRIDE: why is it there?



- Repository to support publications (proteomics MS derived data)
- Source of proteomics data for other data resources

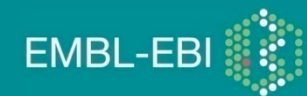


THE LOOK OF PRIDE

Juan A. Vizcaíno
juan@ebi.ac.uk



BSPR/EBI Educational Workshop
Hinxton, 16 July 2010



PRIDE web interface – overview

PRIDE Proteomics IDentifications database (PRIDE) PRIDE Basic Statistics

Search filtered on: Sample Parameters (Species, tissue, disease etc.) , with parameters Type: 'BTO', Value: 'BTO:0000142'

View Instructions

This Table Describes 88 Experiments.

To sort by any of the first seven columns, click the heading.
(Repeated clicking changes the direction of the sort.) Compare Experiments

| Accession | Title | Species | Tissue | Cell Type | GO Term | Disease | Protein Count | Peptide Count | Spectra Count | Retrieve Details (Output format set above.) | Compare Protein Identification Sets | Select Reference Experiment |
|-----------|--|----------------------|------------------------------|-----------|---------|-------------|---------------|---------------|---------------|---|-------------------------------------|-----------------------------|
| 1636 | Proteomics Mapping of Brain Plasma Membrane Proteins | Mus musculus (Mouse) | cerebral cortex, hippocampus | | | | 2356 | 9510 | 0 | Download | <input type="checkbox"/> | |
| 1637 | Characterization of the Mouse Brain Proteome Using Global Proteomic Analysis Complemented with Cysteinylyl-Peptide Enrichment (Protein ID Set 1) | Mus musculus (Mouse) | brain | | | DiseaseFree | 599 | 21218 | 0 | Download | <input type="checkbox"/> | |
| 1638 | Characterization of the Mouse Brain Proteome Using Global Proteomic Analysis Complemented with Cysteinylyl-Peptide Enrichment (Protein ID Set 2) | Mus musculus (Mouse) | brain | | | DiseaseFree | 800 | 11078 | 0 | Download | <input type="checkbox"/> | |
| 1639 | Characterization of the Mouse Brain Proteome Using Global Proteomic Analysis Complemented with Cysteinylyl-Peptide Enrichment (Protein ID Set 3) | Mus musculus (Mouse) | brain | | | DiseaseFree | 1000 | 8255 | 0 | Download | <input type="checkbox"/> | |
| 1640 | Characterization of the Mouse Brain Proteome Using Global Proteomic Analysis Complemented with Cysteinylyl-Peptide Enrichment (Protein ID Set 4) | Mus musculus (Mouse) | brain | | | DiseaseFree | 1600 | 7470 | 0 | Download | <input type="checkbox"/> | |

PART_OF

| | | | |
|-------|--|-------------|------------------------|
| 10116 | Rattus norvegicus (Rat) | BTO:0000383 | kidney tumor cell line |
| 4932 | Saccharomyces cerevisiae (Baker's yeast) | BTO:0000713 | test |
| | | BTO:0001629 | left ventricle |
| 602 | Salmonella typhimurium | BTO:0000759 | liver |

PRIDE web interface – experiment and protein

Experiment View
[the search page](#)

Human CSF analysis (LCQ)

| | |
|--------------|-------------|
| Accession: | I755 |
| Short Label: | LCQ C |
| Source: | Zhang human |
| PubM: | PubM |

Instrument:

Identification Detail View

Details for identification: IPI00400826.1

| Submitted Accession | IPI00400826.1 | | | | | | | | | | | | | | | | |
|--|--|---------------------|----------|---------------------------------|------------------------------|--|---------------------|-------------------------------|--|--|--------|---|-------------|--------------------------------|-------|--------------------------|-------|
| Search Database | ipi.HUMAN | | | | | | | | | | | | | | | | |
| Cross-References | <table border="1" style="width: 100%; border-collapse: collapse;"> <thead> <tr><th>Accession</th><th>Database</th></tr> </thead> <tbody> <tr><td>ENSP00000315130</td><td>ENSEMBL_HUMAN</td></tr> <tr><td>ENSP00000369812</td><td>ENSEMBL_HUMAN</td></tr> <tr><td>IPI00400826.1</td><td>IPI</td></tr> <tr><td>NP_001822.2</td><td>REFSEQ</td></tr> </tbody> </table> | Accession | Database | ENSP00000315130 | ENSEMBL_HUMAN | ENSP00000369812 | ENSEMBL_HUMAN | IPI00400826.1 | IPI | NP_001822.2 | REFSEQ | | | | | | |
| Accession | Database | | | | | | | | | | | | | | | | |
| ENSP00000315130 | ENSEMBL_HUMAN | | | | | | | | | | | | | | | | |
| ENSP00000369812 | ENSEMBL_HUMAN | | | | | | | | | | | | | | | | |
| IPI00400826.1 | IPI | | | | | | | | | | | | | | | | |
| NP_001822.2 | REFSEQ | | | | | | | | | | | | | | | | |
| Search Engine | proteinprophet | | | | | | | | | | | | | | | | |
| Score | 1.0 | | | | | | | | | | | | | | | | |
| Threshold | 0.9 | | | | | | | | | | | | | | | | |
| Sequence Coverage | 0.343 | | | | | | | | | | | | | | | | |
| Additional | <table border="1" style="width: 100%; border-collapse: collapse;"> <thead> <tr><th>Source</th><th>Name</th><th>Value</th></tr> </thead> <tbody> <tr><td>PRIDE</td><td>Protein description line</td><td>Clusterin isoform 1</td></tr> <tr><td>PRIDE</td><td>ProteinProphet probability score</td><td>1.0</td></tr> <tr><td>PRIDE</td><td>Indistinguishable alternative protein accession</td><td>IPI00291262</td></tr> </tbody> </table> | Source | Name | Value | PRIDE | Protein description line | Clusterin isoform 1 | PRIDE | ProteinProphet probability score | 1.0 | PRIDE | Indistinguishable alternative protein accession | IPI00291262 | | | | |
| Source | Name | Value | | | | | | | | | | | | | | | |
| PRIDE | Protein description line | Clusterin isoform 1 | | | | | | | | | | | | | | | |
| PRIDE | ProteinProphet probability score | 1.0 | | | | | | | | | | | | | | | |
| PRIDE | Indistinguishable alternative protein accession | IPI00291262 | | | | | | | | | | | | | | | |
| Mapped Protein sequence | <pre> 0001 MQVCSQPQRG CVREQSAINI APPSAHNAAS PGGARGHRVP LTEACKDSRI GGMKTKLLLF VGLLLTWESG QVLGDTQVSD 0080 0081 NELQEMSNQG SKYVNRKIQN AVNGVKIKIT LIEKTNEERK TLLSNLEEAK KKKEDALNET RESETKLKEI PGVCNETMMA 0160 0161 LWEECKPCLK QTCMKFYARV CRSGSLVGR QLEEFNLQSS PFYFMNGDR IDSLENDRO QTHMLDVMDD HFSRASSIID 0240 0241 ELFQDRFFTR EPQDTHYLP FSLPHRRPHF FFPKSRIVRS LMPFSPYEPL NFHAMFPFPL EMIHEAQQAM DIHFHSPAFO 0320 0321 HPPTEFIREG DDDRTVCREI RHNSTGCLRM KDQCDKREI LSVDCSTNMP SQAKLRRELD ESLQVAERLT RKMELLSKY 0400 0401 QWKMLNTSSL LEQLNEQFNW VSRLAMLTQG EDQYYLRVIT VASHTSDSDV PSGVTEVVVK LFDSDPITVT VPVEVSRKNP 0480 0481 KFMETVAEKA LQEYRKKHRE E </pre> | | | | | | | | | | | | | | | | |
| Submitted Protein Sequence | The protein sequence that was submitted with this identification is identical to the current mapped protein sequence. | | | | | | | | | | | | | | | | |
| Spectrum | <table border="1" style="width: 100%; border-collapse: collapse;"> <thead> <tr><th>Sequence</th><th>Start</th><th>End</th></tr> </thead> <tbody> <tr><td>RELDKSLQVAER</td><td>377</td><td>388</td></tr> </tbody> </table> <p style="text-align: center;">View Spectrum Information</p> <table border="1" style="width: 100%; border-collapse: collapse;"> <thead> <tr><th>Source Name</th><th>Value</th></tr> </thead> <tbody> <tr><td>PRIDE PeptideProphet probability score</td><td>0.8438</td></tr> <tr><td>PRIDE X correlation</td><td>2.650</td></tr> <tr><td>PRIDE Delta Cn</td><td>0.189</td></tr> <tr><td>PRIDE Sp</td><td>997.9</td></tr> </tbody> </table> | Sequence | Start | End | RELDKSLQVAER | 377 | 388 | Source Name | Value | PRIDE PeptideProphet probability score | 0.8438 | PRIDE X correlation | 2.650 | PRIDE Delta Cn | 0.189 | PRIDE Sp | 997.9 |
| Sequence | Start | End | | | | | | | | | | | | | | | |
| RELDKSLQVAER | 377 | 388 | | | | | | | | | | | | | | | |
| Source Name | Value | | | | | | | | | | | | | | | | |
| PRIDE PeptideProphet probability score | 0.8438 | | | | | | | | | | | | | | | | |
| PRIDE X correlation | 2.650 | | | | | | | | | | | | | | | | |
| PRIDE Delta Cn | 0.189 | | | | | | | | | | | | | | | | |
| PRIDE Sp | 997.9 | | | | | | | | | | | | | | | | |

changes in

*SkyPainter is a tool to determine mapped UniProt accessions

PRIDE web interface – mass spectra

Details for spectrum: 245 [Back to the search page](#)

| | |
|------------------|---|
| Spectrum ID: 245 | Peptides Identified From This Mass Spectrum |
|------------------|---|

Peak list Drag the arrows Mass Error 0.5 Daltons For De novo Sequencing: De novo Ion Series End 388

0.0001 Log 10 Include ions where z > 1

mzData Accession

PRIDE Experiment: 8126 Mass Spectrum ID: 4184

Intensity vs m/z plot showing various peaks labeled with peptide sequences and charge states (e.g., Y11, Y5, Y6, Y7, Y8, Y9, Y10, Y11, Co-elute).

PRIDE BioMart

The screenshot displays the PRIDE BioMart interface. At the top, there is a search bar with 'All Databases' selected and a search input field containing 'Enter Text Here'. Navigation tabs include 'Databases', 'Tools', 'Groups', 'Training', 'Industry', 'About Us', and 'Help'. The breadcrumb trail reads 'EBI > Databases > Proteomics > PRIDE > PRIDE BioMart'. Below the search bar, there are buttons for 'New', 'XML', 'Help', 'Count', and 'Results'. The main content area is divided into two columns. The left column contains filter options: 'Dataset: PRIDE', 'Attributes: Submitted Protein Accession', and 'Filters: Filter by Experiment Accession : [ID-list specified]'. The right column shows 'Display maximum: 10 rows as HTML' and 'Export all results to: File'. A table titled 'Submitted Protein Accession' lists the following values: 15079369, Q15526, Q9Y6A2, 21389381, 17149828, O00273, 21359982, P04901, P43251, and Q9H5N1. The footer of the interface indicates 'biomart version 0.5'.



The spectacular bit: across-BioMart queries!

Question: “Which proteins, identified in PRIDE, in blood plasma, → PRIDE are transcribed from genes located in chromosome 11” → Ensembl

The screenshot shows the Ensembl BioMart interface. The top navigation bar includes links for HOME, MARTVIEW, MARTSERVICE, DOCS, CONTACT, NEWS, and CREDITS. Below this, there are buttons for New, Count, and Results, along with options for URL, XML, Perl, and Help. The main interface is divided into two sections: 'Dataset 1895 / 37435 Genes Homo sapiens genes (NCBI36)' and 'Dataset 498 / 8173 Experiments PRIDE'. The 'Dataset 1895' section has filters for 'Chromosome: 11' and 'Attributes' including Ensembl Gene ID, Transcript ID, Start/End coordinates, Chromosome Name, and Gene Name. The 'Dataset 498' section has a filter for 'Filter by Tissue : blood plasma' and 'Attributes' including PRIDE Experiment Accession, Title, Submitted Protein Accession, and Uniprot Accession. The main table displays the results of the query, with columns for Ensembl Gene ID, Ensembl Transcript ID, Gene Start/End (bp), Chromosome Name, Associated Gene Name, Ensembl Protein ID, PRIDE Experiment Accession, Experiment Title, Submitted Protein Accession, and Uniprot Accession. The table lists 13 rows of data, all from chromosome 11, showing genes like HBB, APOA1, F2, SERPING1, APOC3, and HPX, and their corresponding proteins and PRIDE experiments.

| Ensembl Gene ID | Ensembl Transcript ID | Gene Start (bp) | Gene End (bp) | Chromosome Name | Associated Gene Name | Ensembl Protein ID | PRIDE Experiment Accession | Experiment Title | Submitted Protein Accession | Uniprot Accession |
|---------------------------------|---------------------------------|---------------------------|---------------------------|-----------------|----------------------|---------------------------------|----------------------------|---|-----------------------------|------------------------|
| ENSG00000221842 | ENST00000335295 | 5203272 | 5207201 | 11 | HBB | ENSP00000333994 | 13 | HUPO Plasma Proteome Project, Lab # 2 Expt # 17 | IPI00218816 | P68871 |
| ENSG00000118137 | ENST00000236850 | 116211677 | 116213571 | 11 | APOA1 | ENSP00000236850 | 13 | HUPO Plasma Proteome Project, Lab # 2 Expt # 17 | IPI00021841 | P02647 |
| ENSG00000118137 | ENST00000375320 | 116211677 | 116213571 | 11 | APOA1 | ENSP00000364469 | 13 | HUPO Plasma Proteome Project, Lab # 2 Expt # 17 | IPI00021841 | P02647 |
| ENSG00000118137 | ENST00000375323 | 116211677 | 116213571 | 11 | APOA1 | ENSP00000364472 | 13 | HUPO Plasma Proteome Project, Lab # 2 Expt # 17 | IPI00021841 | P02647 |
| ENSG00000118137 | ENST00000359492 | 116211677 | 116213571 | 11 | APOA1 | ENSP00000352471 | 13 | HUPO Plasma Proteome Project, Lab # 2 Expt # 17 | IPI00021841 | P02647 |
| ENSG00000180210 | ENST00000311907 | 46697331 | 46717631 | 11 | F2 | ENSP00000308541 | 13 | HUPO Plasma Proteome Project, Lab # 2 Expt # 17 | IPI00019568 | P00734 |
| ENSG00000149131 | ENST00000278407 | 57121436 | 57138902 | 11 | SERPING1 | ENSP00000278407 | 13 | HUPO Plasma Proteome Project, Lab # 2 Expt # 17 | IPI00291866 | P05155 |
| ENSG00000110245 | ENST00000375345 | 116205834 | 116208998 | 11 | APOC3 | ENSP00000364494 | 13 | HUPO Plasma Proteome Project, Lab # 2 Expt # 17 | IPI00021857 | P02656 |
| ENSG00000110245 | ENST00000227667 | 116205834 | 116208998 | 11 | APOC3 | ENSP00000227667 | 13 | HUPO Plasma Proteome Project, Lab # 2 Expt # 17 | IPI00021857 | P02656 |
| ENSG00000110169 | ENST00000265983 | 6408858 | 6418830 | 11 | HPX | ENSP00000265983 | 13 | HUPO Plasma Proteome Project, Lab # 2 Expt # 17 | IPI00022488 | P02790 |

www.biomart.org

DATA SUBMISSION TO PRIDE

Juan A. Vizcaíno
juan@ebi.ac.uk



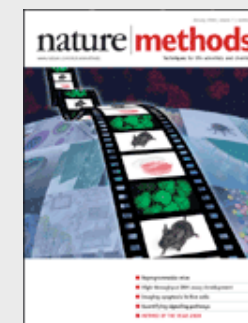
BSPR/EBI Educational Workshop
Hinxton, 16 July 2010



Journals recommend PRIDE as submission point

- Journal guidelines recommend now submission to proteomics repositories:

- *Proteomics*
- *Nature Biotechnology*
- *Nature Methods*
- *Molecular and Cellular Proteomics*



- Closer collaboration between *Proteomics* and PRIDE:
 - “Deposition of supporting data in a public, open access database like PRIDE or World-2DPAGE is strongly recommended, and **mandatory** for Dataset Briefs”

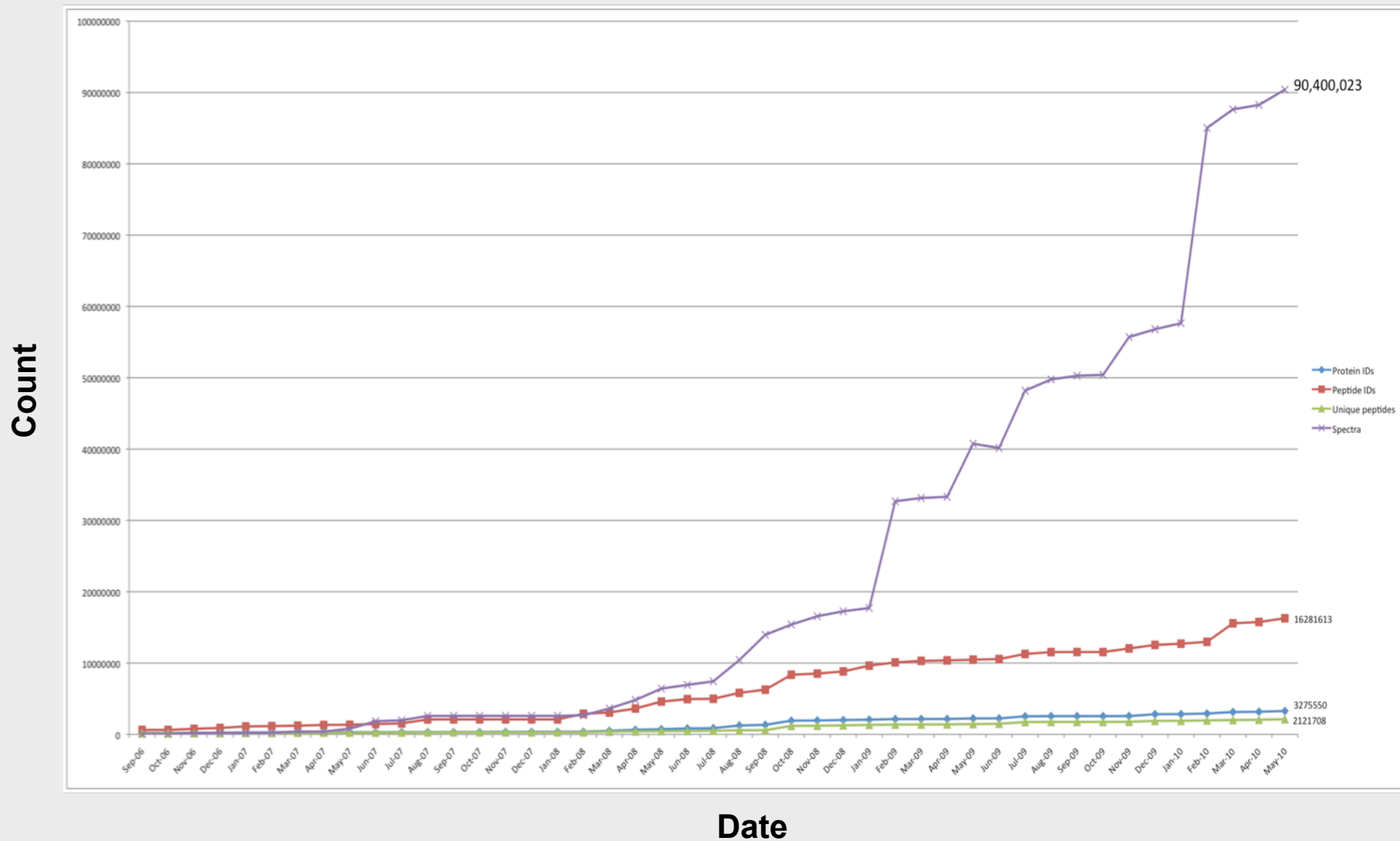
MCP new guidelines

New guidelines from **MCP** for **data deposition**:

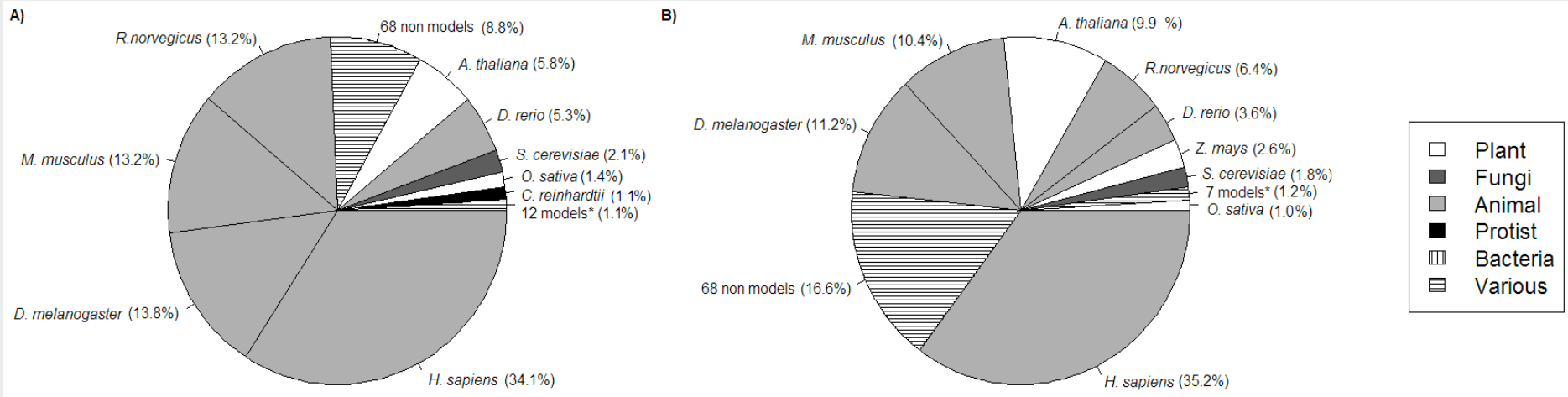
For all proteins identified on the basis of **ONE OR TWO unique peptide spectra**, the ability to view **annotated spectra** for these identifications must be made available. This can be achieved in one of three ways:

- 1) Submission of spectra and search results to a **public results repository** that is **equipped with a spectral viewer** (e.g. **PRIDE**, Peptidome etc). This information will appear as a **hyperlink** in the published article...
- 2) Submission (with the manuscript) of spectra and search results in a file format that allows visualization of the spectra using a **freely-available viewer**.
- 3) Submission (with the manuscript) of annotated spectra in an 'office' or PDF format.

PRIDE growth



PRIDE data content

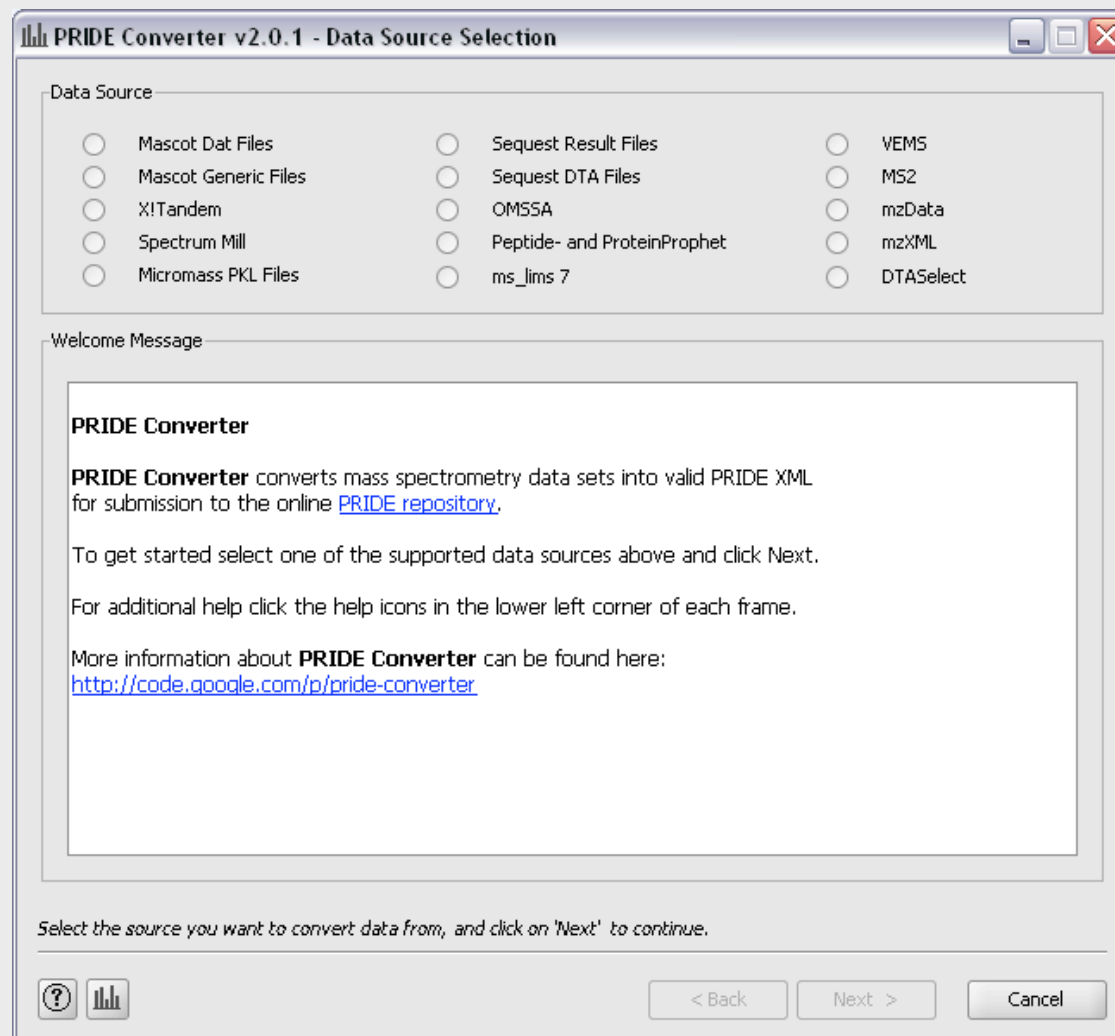


Protein IDs

Peptide IDs

Why? Submission made easier: PRIDE Converter

<http://code.google.com/p/pride-converter>



Barsnes *et al.*, 2009

PRIDE Converter – interface details

The screenshot displays three overlapping windows from the PRIDE Converter v1.16.2 software:

- Spectra Selection - Step 2 of 8:** Contains sections for 'Simple Spectra Selection' (with 'Select All Spectra' selected), 'Advanced Spectra Selection', and 'Manual Spectra Selection' (with a table for PID and Filename).
- Protocol Properties - Step 5 of 8:** Shows a table for 'Protocol' with columns for Name and Protocol S.
- Instrument Properties - Step 6 of 8:** Contains fields for 'Instrument Name' (Bruker Ultraflex), 'Source' (Matrix-assisted Laser Desorption Ionization [PSI:1000075]), and 'Detector' (Electron Multiplier Tube [PSI:1000111]). It also features an 'Analyzers' table and a 'Processing' section with a table for 'Processing Methods'.

Processing Methods Table:

| CV Terms | Value |
|-------------------------------------|----------------------|
| 1 Deisotoping [PSI:1000033] | false |
| 2 ChargeDeconvolution [PSI:1000034] | false |
| 3 PeakProcessing [PSI:1000035] | CentroidMassSpectrum |

From PRIDE Converter to PRIDE FTP

PRIDE Converter v2.3.4 - Output Properties - Step 8 of 8

Output Folder: /Users/javizca/

Resubmission

Resubmission * Original Accession Number: []

* When resubmitting a PRIDE XML file please provide the original accession

Format Specific Parameters

Round Score and Threshold Down Before Comparison

Use Comma As [] Transforming Spectra. Please

OMSSA Folder: [] F007448.dat (1/1)

PRIDE Submission

File Created: []

If the data is part of a paper, please include a reference to PRIDE Converter: []

Select an output folder and click on 'Convert!' to create the PRIDE XML file.

PRIDE Login

PRIDE Registration

Request PRIDE FTP Access?

Submission Tips

The size of your PRIDE XML file is larger than the maximum file size for using the 'Direct Submission' via the PRIDE web page.

We therefore recommend using the PRIDE FTP server. To get access to the FTP server, please contact the PRIDE team at pride-support@ebi.ac.uk.

Close

PRIDE Converter v2.3.4 - Output Properties - Step 8 of 8

Output Folder: /Users/javizca/

Resubmission

Resubmission * Original Accession Number: []

* When resubmitting a PRIDE XML file please provide the original accession

Format Specific Parameters

Round Score and Threshold Down Before Comparison

Use Comma As [] Transforming Spectra. Please

OMSSA Folder: []

PRIDE Submission

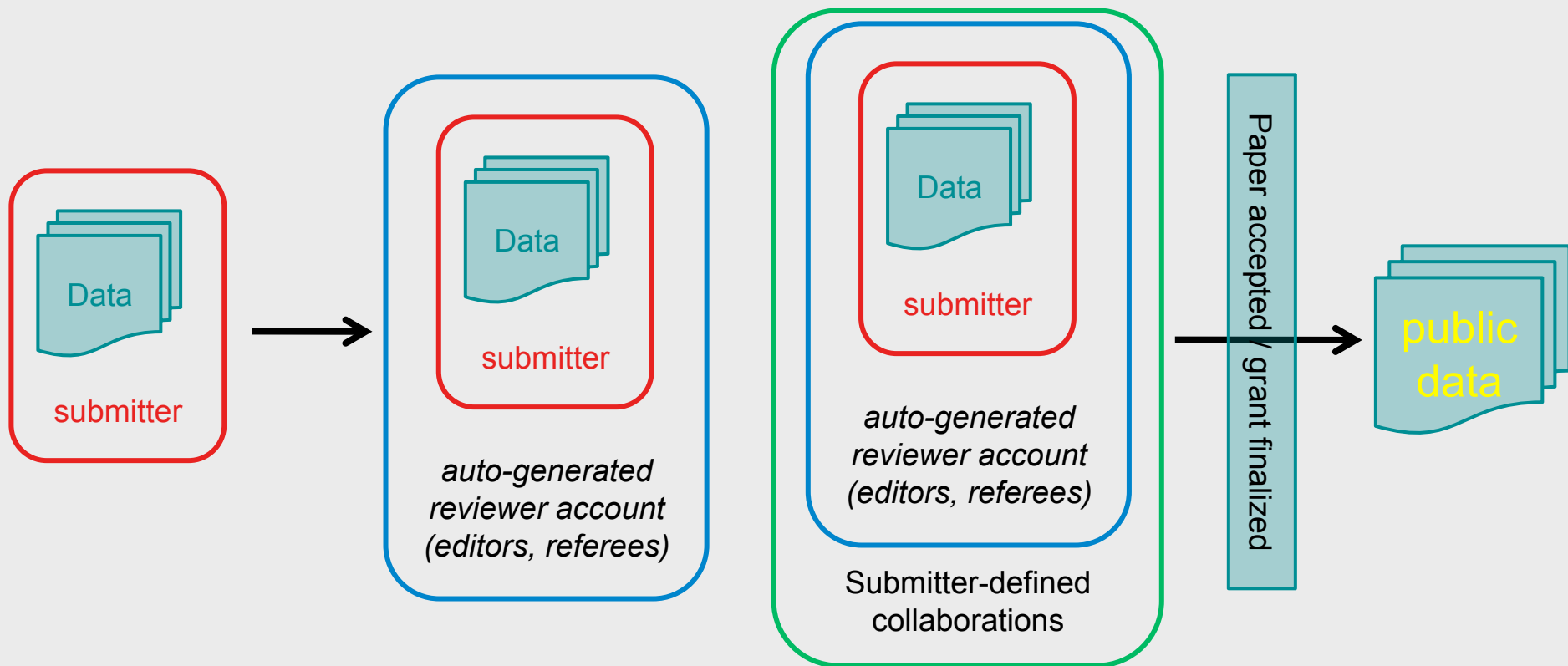
File Created: []

If the data is part of a paper, please include a reference to PRIDE Converter: [PMID:19587657](https://pubmed.ncbi.nlm.nih.gov/19587657/)

Select an output folder and click on 'Convert!' to create the PRIDE XML file.

< Back Convert! Exit

Data access privileges in PRIDE



PRIDE relies on a simple but very powerful group-based access system that can accommodate even more complex data release schemes than pictured here

OTHER PROTEOMICS REPOSITORIES

Existing proteomics repositories

- Main public repositories:

- PROteomics IDentifications database (PRIDE)
- Global Proteome Machine (GPMDB)
- Peptide Atlas
- Tranche
- NCBI Peptidome








- Smaller scale repositories, more specialized:

Among others: Human Proteinpedia, Genome Annotation Proteomics Pipeline (GAPP), MAPU, SwedCAD, PepSeeker, Open Proteomics Database, ...

- Very diverse: different aims, functionalities, ...

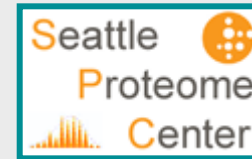
Other MS proteomics repositories

|  |  |  |  |  Tranche |
|---|---|---|---|--|
| <i>Reprocesses data</i> | <i>Reprocesses data</i> | No reprocessing | No reprocessing | No reprocessing |
| <i>Editorial control</i> | <i>Editorial control</i> | No editorial control | No editorial control | No editorial control |
| <i>Limited annotation</i> | <i>Limited annotation</i> | Detailed annotation | Detailed annotation | <i>Limited annotation</i> |
| ?? | 170 million peptides | 96 million spectra | 3.8 million spectra | ?? |
| ?? | 22.3 million protein IDs | 3.7 million protein IDs | 60,000 protein IDs | ?? |

PeptideAtlas

- Peptide identifications from MS/MS
- Data are reprocessed using the popular *Trans Proteomic Pipeline (TPP)*
- Uses *PeptideProphet* to derive a probability for the correct identification for all contained peptides

<http://www.peptideatlas.org>



ISB Home

PeptideAtlas

PEPTIDEATLAS HOME

Seattle Proteome Center

PEPTIDEATLAS:
Overview
Contacts
Data Contributors
Publications
Software
Database Schema
Feedback
FAQ

ATLAS DATA:
Data Repository
HPPP Data Central
PeptideAtlas Builds
Search Database

Contribute Data
Genome Browser
Setup

RELATED:
MRM Atlas
Phosphopep
Unipep
mspecLINE

SPECTRAL LIBS:
Libraries + Info
SpectraST Search

GLOSSARY/TERMS:
Atlas nomenclature
SGD nomenclature
Protein ID terms

LOGIN

PeptideAtlas

Search PeptideAtlas: GO

PeptideAtlas is a multi-organism, publicly accessible compendium of peptides identified in a large set of tandem mass spectrometry proteomics experiments. Mass spectrometer output files are collected for human, mouse, yeast, and several other organisms, and searched using the latest search engines and protein sequences. All results of sequence and spectral library searching are subsequently processed through the *Trans Proteomic Pipeline* to derive a probability of correct identification for all results in a uniform manner to insure a high quality database, along with false discovery rates at the whole atlas level. Results may be queried and browsed at the PeptideAtlas web site. The raw data, search results, and full builds can also be downloaded for other uses.

Related Resources

HO
O=C
HO

Phosphopep MRAA Atlas Unipep

Atlas News

News 2010-03: Members of the PeptideAtlas team have recently published *mspecLINE*, a [web application](#) that allows researchers to explore relationships between human diseases and the observed proteome.

News 2010-02: For the first time, a Mouse PeptideAtlas build, based on 64 samples from a variety of tissues and subcellular compartments, is available to [search](#).

News 2009-08: A new build of the Drosophila PeptideAtlas is now available to [search](#). The data is described in [this publication by Erich Brunner et al.](#)

News 2009-06-30: A new build of the Human PeptideAtlas is now available to [search](#). In this build, we used much more stringent criteria, spectra FDR 0.00001, to report the peptides identified, so the number of distinct peptides identified is much lower than the previous build.

News 2009-06: New version 3.0 spectrum libraries from NIST now available at the [PeptideAtlas Spectrum Library Central Page](#).

PeptideAtlas

- All peptides mapped to *Ensembl* using *ProteinProphet* (for human)
- Built by the Aebersold lab to help them find proteotypic peptides
- Provides proteotypic peptide predictions
- Limited metadata
- Great support for targeted proteomics approaches (SRM/MRM)

<http://www.peptideatlas.org>



ENSP00000222219

Protein Name: ENSP00000222219
 Gene Name: ENSG00000105612
 Description: pep.known-ccds chromosome:NCBI36:19:12847023:12853335:-1
 gene:ENSG00000105612 transcript:ENST00000222219 CCDS12284.1
 Ensembl Protein: [ENSP00000222219](#)
 Entrez Gene Symbol: [DNASE2](#)
 Entrez GeneID: [1777](#)
 Full Name: DNASE2, DNASE2A, DNL2: Deoxyribonuclease-2-alpha precursor
 IPI: [IPI00010348](#)
 RefSeq: [NP_001366](#)
 RefSeq: [REVIEWED:NP_001366](#)
 UniGene: Hs.118243
 UniProt: [O00115](#)
 UniProt Symbol: [DNS2A_HUMAN](#)
 Total peptides: 11

Sequence Motifs

Sequence Position: 0 100 200 300

Observed Peptides: PAp00005361, PAp00041670, PAp00041671, PAp00040423

Anchor Sequence: [Anchor sequence predicted by Signal P]

Extracellular Domain: [Extracellular domain predicted by TMHMM]

Sequence Coverage: [Protein coverage by observed peptides]

Sequence Position: 0 100 200 300

Sequence: MIPILLALL CVPAGALTCY **CDSGQPVDFW** VYIKLPALRG SGEAAQRGLQ YKYLDESSGG WRDGRALINS PEGAVRSILQ PLYRSNTSQL AFLLYNDQPD QPSKAQDSM RCHTKGVLLL DHGGFWLVN SVFNFPFPPAS SAAYSWPHSA CTYQTLILCV SFPFAQFSKM GQQLTYTYPN VYNYQLEGIF AQEFPDLENV VKGHVSGEP WNSSITLTSQ AGAVEQSFQK FSKFGDDLYS GWLAAALGYN LQQVFWHKTIV GILPSNCSDI WQVLNVNQIA FPGPAGPSFN STEDHSKWCV SPKGPWTCVG DMRNRQGEQ **RGGGFLCAQL PALMKAFQPL** VKNYQPCNMG ARKPSRAYKI

Protein Coverage = 12.5%

Observed Peptides

| Peptide Accession | Peptide Sequence | Best Prob | Best Adjusted Prob | N Obs | Empirical Proteotypic Score | SSRCalc Relative Hydrophob | N Protein Mappings | N Genome Locations | Sample IDs |
|-------------------|--------------------------|-----------|--------------------|-------|-----------------------------|----------------------------|--------------------|--------------------|------------|
| PAp00005361 | LTCYGDSSGQPVDFWVWYK | 0.999 | 0.970 | 4 | 0.75 | 38.98 | 1 | 1 | 27,7,8 |
| PAp00040423 | GGGTLCAGLPALWKAQPLVK | 1.000 | 1.000 | 4 | 0.25 | 45.12 | 1 | 1 | 54 |
| PAp00041670 | LTCYGDSSGQPVDFWVYKLPALR | 1.000 | 1.000 | 1 | 0.25 | 47.81 | 1 | 1 | 54 |
| PAp00041671 | LTCYGDSSGQPVDFWVYKLPALRG | 1.000 | 1.000 | 2 | 0.25 | 46.88 | 1 | 1 | 54 |

GPMDDB



- End point of the *GPM proteomics pipeline*, to aid in the process of validating peptide MS/MS spectra and protein coverage patterns.

the gpmdb

Search by: [accession](#) [gpm #](#) [sequence](#) [keyword](#) [ontology](#)
Information: [home](#) [statistics](#) [species](#) [thegpm](#) [about](#)

what is the gpm powered by GPMDDB send us email

Eukaryote proteomes
1 2 3 4 5 6 7

Boutique proteomes
human mouse frog
cow bacteria plant
fish rat

Algorithms
X! P3 X! Hunter

Information
gpmDB wiki
review lists

Images:

the gpm

companies
1. Proteome Software
2. Beavis Informatics

data
1. Tranche
2. PeptideAtlas
3. PRIDE
4. PNNL
5. Peptidome

information
1. Proteome Commons
2. Unimod
3. NCTA

organizations
1. HUPO
2. CNPN
3. US HUPO
4. EUPA

expression
1. Allen Atlas
2. The HPR

pathways
1. KEGG

site reference
Craig, et al. (2004) JPR, 3:1234-42.

What is GPMDDB

Search by protein description keywords
Keywords: [View matches](#)
Data source: [Homo sapiens - ensembi](#)
Examples: [ABLI \(human\) \(more ...\)](#)

Or search by data set keywords:
Keywords: [View matches](#)
Examples: [mitochondria \(more ...\)](#)

Global Proteome Machine News:

1. [Data set of the week: 27 June 2010](#)
mTAL Phosphoproteome Data.
2. [Data set of the week: 20 June 2010](#)
Proteomic analysis of mouse brain microsomes.
3. [Data set of the week: 13 June 2010](#)
The minor salivary gland proteome in Sjogren's syndrome.
4. [Data set of the week: 6 June 2010](#)
Identification of ricin and concanavalin A-binding Trypanosoma brucei glycoproteins.
5. [Mouse protein phosphorylation sites](#)
List of the 10,266 observed mouse phosphorylation sites in GPMDDB.

If you do not see a red dot below, you will need Adobe's [SVG plugin](#).

Hosted by: Manitoba Centre for Proteomics and Systems Biology

<http://gpmdb.thegpm.org/>

GPMDDB



- End point of the *GPM proteomics pipeline*, to aid in the process of validating peptide MS/MS spectra and protein coverage patterns.
- Data are reprocessed using the popular *X!Tandem* or *X!Hunter* spectral searching algorithm
- Also provides proteotypic peptides

1. Data set of the week: 27 June 2010
mTAL Phosphoproteome Data.

2. Data set of the week: 20 June 2010
Proteomic analysis of mouse brain microsomes.

advanced page
view saved xml data

Lookup model:
GPM go

what is the gpm powered by tandem send us email

Eukaryote proteomes
1 2 3 4 5 6 7

Boutique proteomes
human mouse frog
cow bacteria plant
fish rat

Algorithms
X! P3 X! Hunter

Information
gpmDB wiki
review lists

the gpm

GPM Cyclone, simple search form

1. **spectra**
common, mzXML, mzData, DTA, PKL or MGF only
 Browse...

2. **taxon**
Select one or more.
 Eukaryotes Prokaryotes Viruses

GRCh 37 (ENSEMBL)
Human (SwissProt)
H. sapiens (NCBI Unigene)
H. sapiens (NCBI RefSeq)

cRAP artifacts
none

none
H. sapiens microbiome
Acaryochloris marina MBIC11017
Acetobacter pasteurianus IFO 3283 01
Acholeplasma laidlawii PG 8A
Acidaminococcus fermentans DSM 20731 uid43471
Acidimicrobium ferrooxidans DSM 10331
Acidiphilium cryptum JF-5

1. Include reversed sequences: none | mixed | only |
2. all ¹⁵N amino acids

Find proteins with peptide log(e) < -1 and protein log(e) < -1

3. **measurement errors**
1. Fragment mass error: 0.4 Da

4. **residue modifications**
1. Complete modifications 1:
Carbamidomethyl (C)
57.021464@C specify your own
2. Complete modifications 2:
No further mods specify your own
3. Potential modifications:
none
Oxidation (M) specify your own 15.994915@M
Oxidation (W)
Deamidation (N)

4. Use sequence annotations yes no

<http://gpmdb.thegpm.org/>

GPMDB



- Powerful visualization features
- Provides very limited annotation with GO, BTO
- Some support to targeted approaches is available

the gpmdb

Search by: [accession](#) | [gpm #](#) | [sequence](#) | [keyword](#) | [ontology](#)
Information: [home](#) | [statistics](#) | [species](#) | [thegpm](#) | [about](#)

Accession number: [View matches](#)

285 matches for *ENSP00000249364*
| [ensembl](#) | [ncbi](#) | [omim](#) | [unigene](#) | [hapmap](#) | [snps](#) | [geo](#) | [human protein atlas](#) |
| [kegg](#) | [hmdb](#) | [grid](#) | [hgnc](#) | [uniprot](#) | [peptideatlas](#) | [pride](#) | [gpmDB](#) |

Calumenin precursor (Crocabin) (IEF SSP 9302). Source: Uniprot/SWISSPROT O43852

Annotated domains:
IPR013623 NADPH oxidase Respiratory burst
IPR002048 Calcium-binding EF-hand
IPR013684 Miro-like
IPR013567 EF hand associated, type-2

Best models for *ENSP00000249364* [Show all](#)

| # | log(e) | model | coverage |
|-----|--------|-----------|----------|
| 1. | -122.2 | G P X | |
| 2. | -122 | G P X | |
| 3. | -121.8 | G P X | |
| 4. | -121.7 | G P X | |
| 5. | -121.4 | G P X | |
| 6. | -118.2 | G P X | |
| 7. | -118 | G P X | |
| 8. | -116.4 | G P X | |
| 9. | -116.3 | G P X | |
| 10. | -116 | G P X | |
| 11. | -115.9 | G P X | |

<http://gpmdb.thegpm.org/>

NCBI Peptidome



- No reprocessing
- Detailed annotation (no CVs)
- Review system
- Sibling resource to PRIDE

NCBI Peptidome Home Browse Data Search Data Submit Data Submission Guidelines Contact Us

NCBI » Peptidome » Browse Data » PSM1002

Sample PSM1002






| Name | Accession | Organism | Gene | Length | Mass | Peptides | Spectra | Define |
|----------------------|-------------|--------------------------|------|--------|----------|----------|---------|---|
| GENEFINDER0000007354 | - | - | - | - | - | 1 | 1 | pep:GeneFinder chromosome:SGD1.01.VII:500685:504574:1 transcript:GE |
| GENEFINDER0000007497 | - | - | - | - | - | 1 | 1 | pep:GeneFinder chromosome:SGD1.01.VII:649530:650925:-1 transcript:GE |
| K1C1_HUMAN | - | - | - | - | - | 7 | 37 | UPSP:K1C1_HUMAN P3527 homo sapiens (human), keratin, type I cytoske |
| K1CJ_HUMAN | - | - | - | - | - | 13 | 100 | UPSP:K1CJ_HUMAN P13645 homo sapiens (human), keratin, type I cytoski |
| K1CM_HUMAN | - | - | - | - | - | 2 | 4 | UPSP:K1CM_HUMAN P13646 homo sapiens (human), keratin, type I cytosk |
| K1CN_HUMAN | - | - | - | - | - | 1 | 1 | UPSP:K1CN_HUMAN P02533 homo sapiens (human), keratin, type I cytosk |
| K22E_HUMAN | P35908.1 | Homo sapiens | KRT2 | 645 | 65865.31 | 10 | 60 | RecName: Full=Keratin, type II cytoskeletal 2 epidermal; AltName: Full=Cytc |
| K2C1_HUMAN | P04264.5 | Homo sapiens | KRT1 | 644 | 66017.66 | 13 | 89 | RecName: Full=Keratin, type II cytoskeletal 1; AltName: Full=Cytokeratin-1; |
| K2C4_HUMAN | P19013.4 | Homo sapiens | KRT4 | 534 | 57285.24 | 1 | 1 | RecName: Full=Keratin, type II cytoskeletal 4; AltName: Full=Cytokeratin-4; |
| K2CA_HUMAN | - | Saccharomyces cerevisiae | - | - | - | 2 | 2 | UPSP:K2CA_HUMAN P02538 homo sapiens (human), keratin, type II cytosk |
| Q0045 | P00401.2 | Saccharomyces cerevisiae | COX1 | 534 | 58798.13 | 1 | 11 | RecName: Full=Cytochrome c oxidase subunit 1; AltName: Full=Cytochrome |
| Q0085 | NP_009313.1 | Saccharomyces cerevisiae | ATP6 | 259 | 29099.04 | 2 | 8 | Atp6p [Saccharomyces cerevisiae] |
| Q0105 | NP_009315.1 | Saccharomyces cerevisiae | COB | 385 | 43655.96 | 1 | 1 | Cobp [Saccharomyces cerevisiae] |
| Q0250 | P00410.1 | Saccharomyces cerevisiae | COX2 | 251 | 28567.26 | 1 | 2 | RecName: Full=Cytochrome c oxidase subunit 2; AltName: Full=Cytochrom |
| TRYP_PIG | P00761.1 | Sus scrofa | - | 231 | 24409.42 | 15 | 345 | RecName: Full=Trypsin; Flags: Precursor |
| UPSP:HSP71_YEAST | - | - | - | - | - | 1 | 1 | P10591 saccharomyces cerevisiae (baker's yeast), heat shock protein ssa1 |
| YAL003W | P32471.4 | Saccharomyces cerevisiae | EFB1 | 206 | 22627.12 | 8 | 42 | RecName: Full=Elongation factor 1-beta; Short=EF-1-beta; AltName: Full=Ti |
| YAL005C | P10591.4 | Saccharomyces cerevisiae | SSA1 | 642 | 69657.23 | 34 | 268 | RecName: Full=Heat shock protein SSA1; AltName: Full=Heat shock protein |
| YAL012W | P31373.2 | Saccharomyces cerevisiae | CYS3 | 394 | 42542.04 | 10 | 28 | RecName: Full=Cystathionine gamma-lyase; AltName: Full=Gamma-cystath |
| YAL023C | P31382.2 | Saccharomyces cerevisiae | PMT2 | 759 | 86869.83 | 2 | 5 | RecName: Full=Dolichyl-phosphate-mannose-protein mannosyltransferase |

Page 1 of 38

Displaying proteins 1 - 20 of 744

<http://www.ncbi.nlm.nih.gov/peptidome/>

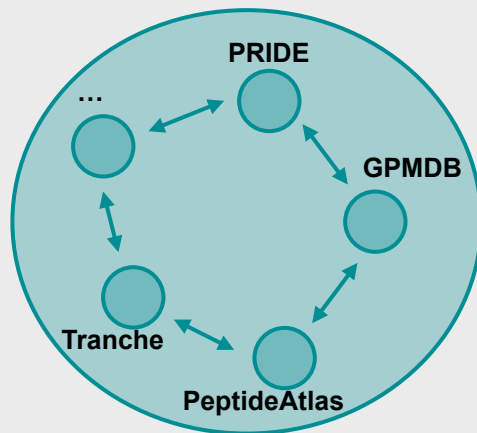
Other MS proteomics repositories

|  |  |  |  |  |
|---|---|---|---|---|
| <i>Reprocesses data</i> | <i>Reprocesses data</i> | No reprocessing | No reprocessing | No reprocessing |
| <i>Editorial control</i> | <i>Editorial control</i> | No editorial control | No editorial control | No editorial control |
| <i>Limited annotation</i> | <i>Limited annotation</i> | Detailed annotation | Detailed annotation | <i>Limited annotation</i> |
| ?? | 162 million peptides | 92 million spectra | 3.8 million spectra | ?? |
| ?? | 21.5 million protein IDs | 3.5 million protein IDs | 60,000 protein IDs | ?? |

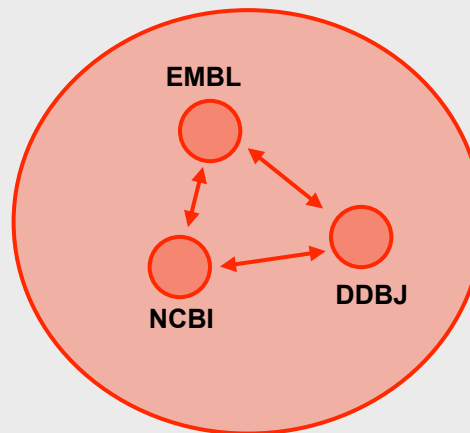
PRIDE AND OTHER REPOSITORIES: ProteomeXchange

For sharing, superstructures must be built

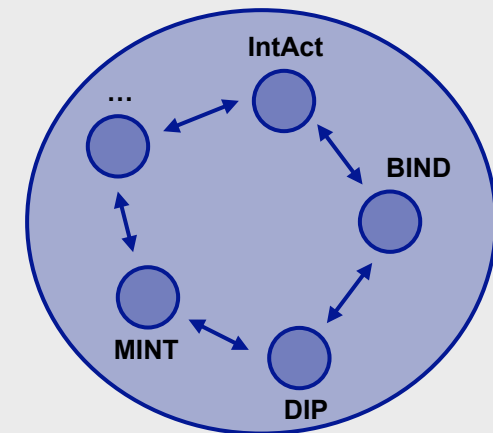
Often, multiple repositories will emerge more or less simultaneously in a particular field. By exchanging data, and by collaborating on data acquisition an increase in coverage as well as a more comprehensive dataset is obtained by each individual resource. Such superstructures do require additional infrastructure, however.



mass spec
ProteomeXchange



sequence databases
(INSDC)



interactions
IMEx

ProteomeXchange consortium



- Sharing proteomics data between existing proteomics repositories
- Includes PeptideAtlas, GPMDB, NCBI Peptidome and PRIDE, with data sharing infrastructure provided by Tranche
- Submission guidelines document finalized, it was proven on three different datasets
- ProteomeXchange is primarily **user-oriented**: the idea is to provide a **single point of submission**, but **multiple points of data visualization and analysis**

Proteomics data submission strategy for ProteomeXchange

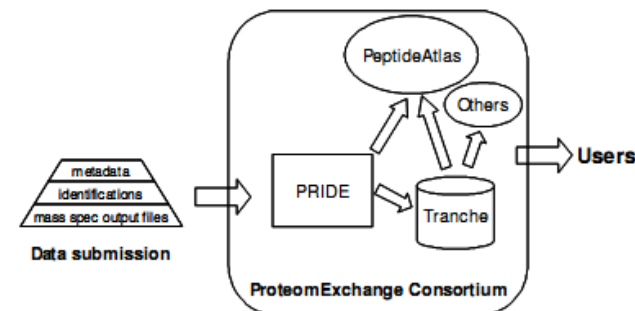
1. Summary

This document provides detailed guidelines for the submission of mass spectrometry-derived proteomics data to the ProteomeXchange consortium¹ databases PRIDE^{2,3}, PeptideAtlas^{6,7}, and Tranche^{7,8}. First the policy is summarized in this section; then in subsequent sections, definitions of terms, descriptions of the relevant resources, details on the submission path, and policies regarding data ownership and data privacy are provided. This policy has been adopted by the HUPO Plasma Proteome Project^{9,10} for the collection of its Phase II data; it is hoped that widespread adoption will follow.

Each submission shall consist of three major components: mass spectrometer output files, study metadata, and peptide/protein identifications (further details in section 4; definitions provided in section 2). All submissions will include all three components and will be made to the PRIDE repository using data sufficiency guidelines established by PRIDE as described below.

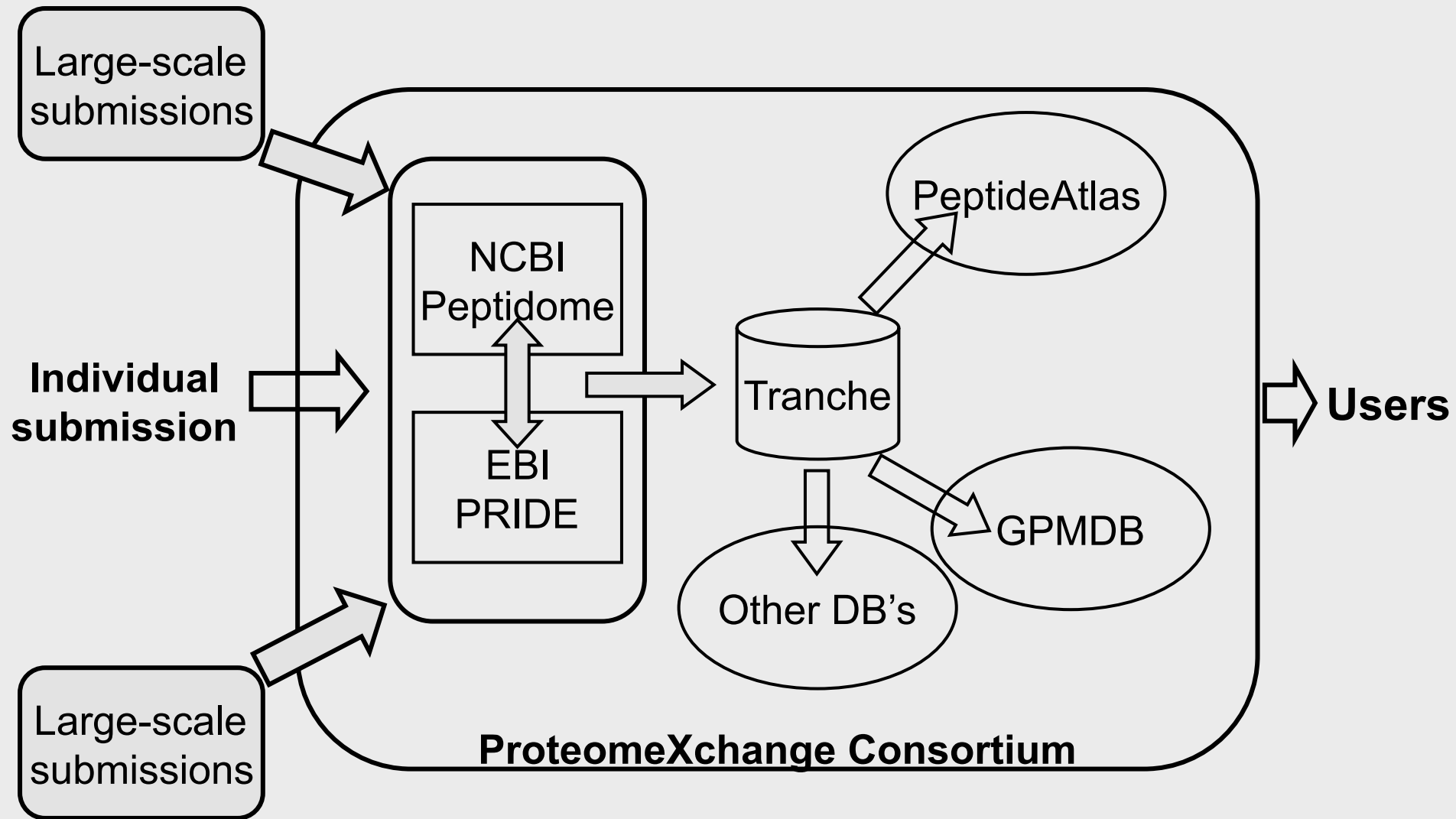
At the time when the submitted data are declared publicly available by the submitter, all mass spectrometer output files will be deposited in the Tranche repository. Hash keys required to download this information from Tranche and study metadata will be displayed in PRIDE and actively transmitted to PeptideAtlas and any other participating ProteomeXchange repositories (see section 3 for information about the individual repositories) for further processing.

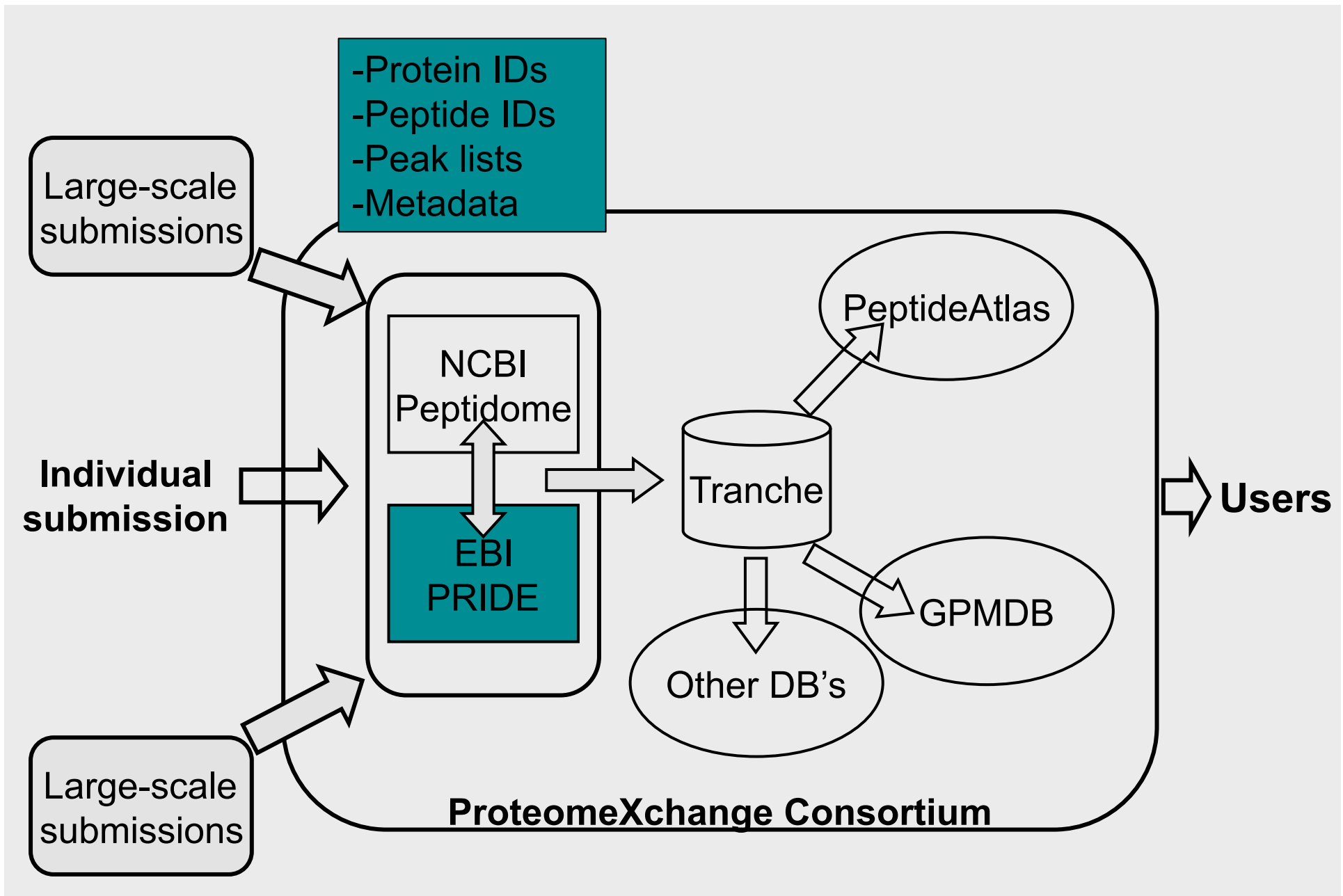
This insures that a simple one-time submission from a contributor is automatically distributed to all ProteomeXchange repositories with sufficient information.

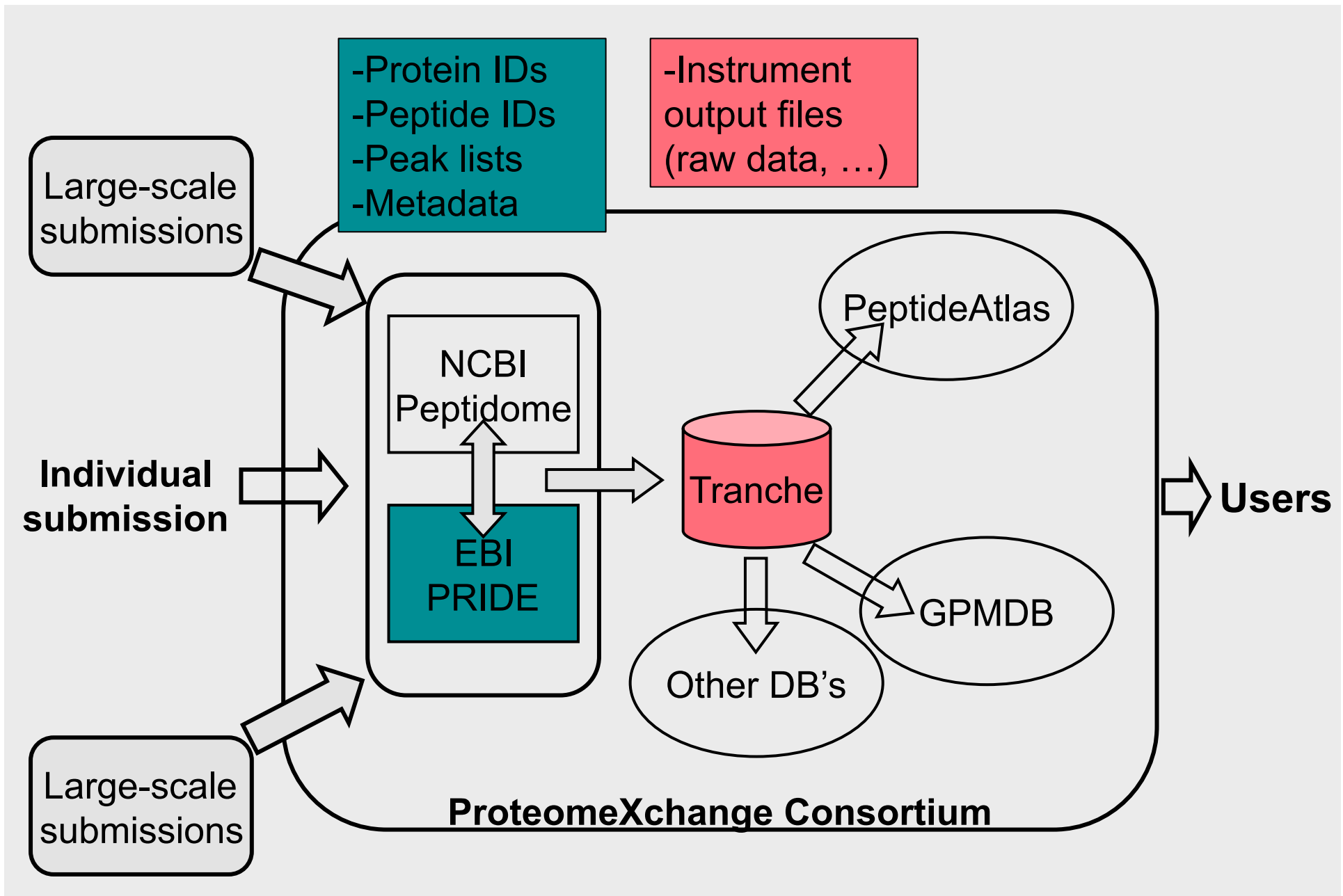


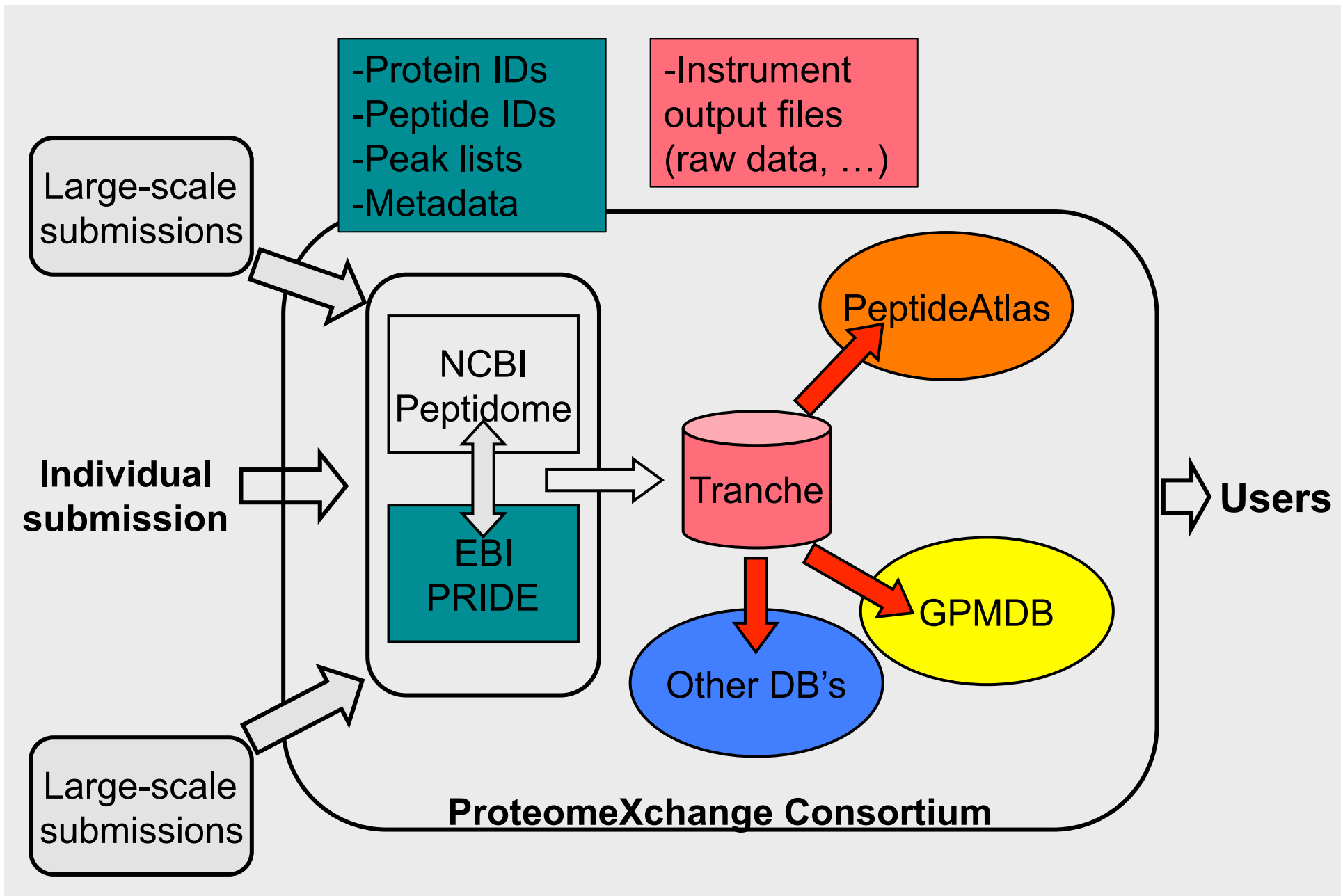
Summary Figure: data submissions are sent to the ProteomeXchange Consortium via PRIDE. The ProteomeXchange partners then ensure data are distributed internally, ultimately giving users the ability to access the data from any participating database.

www.proteomexchange.org





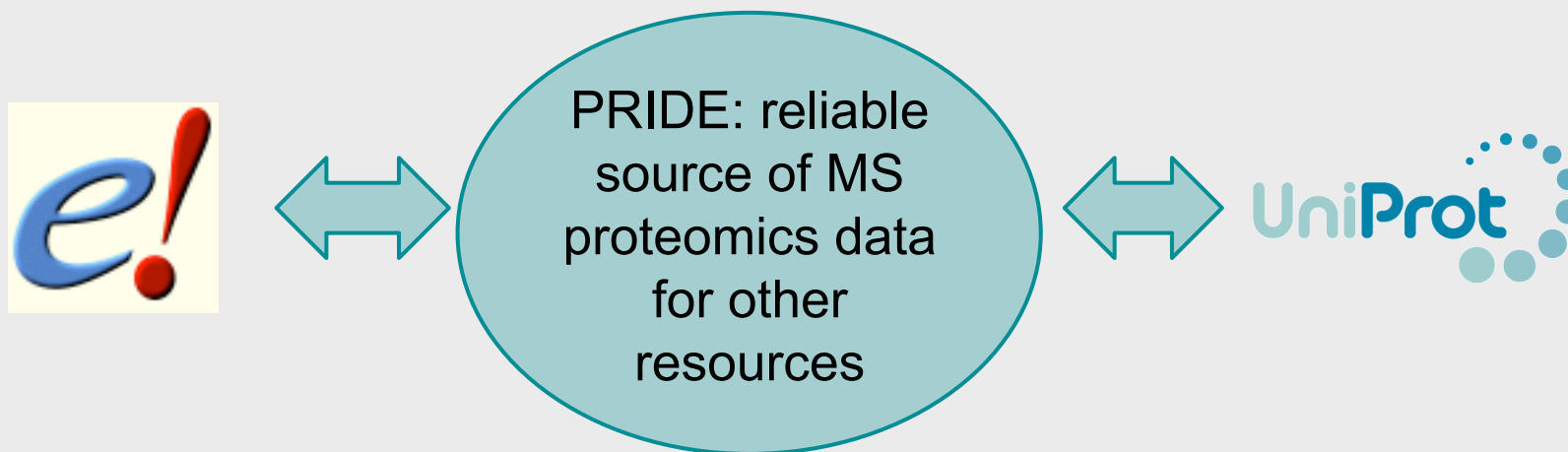




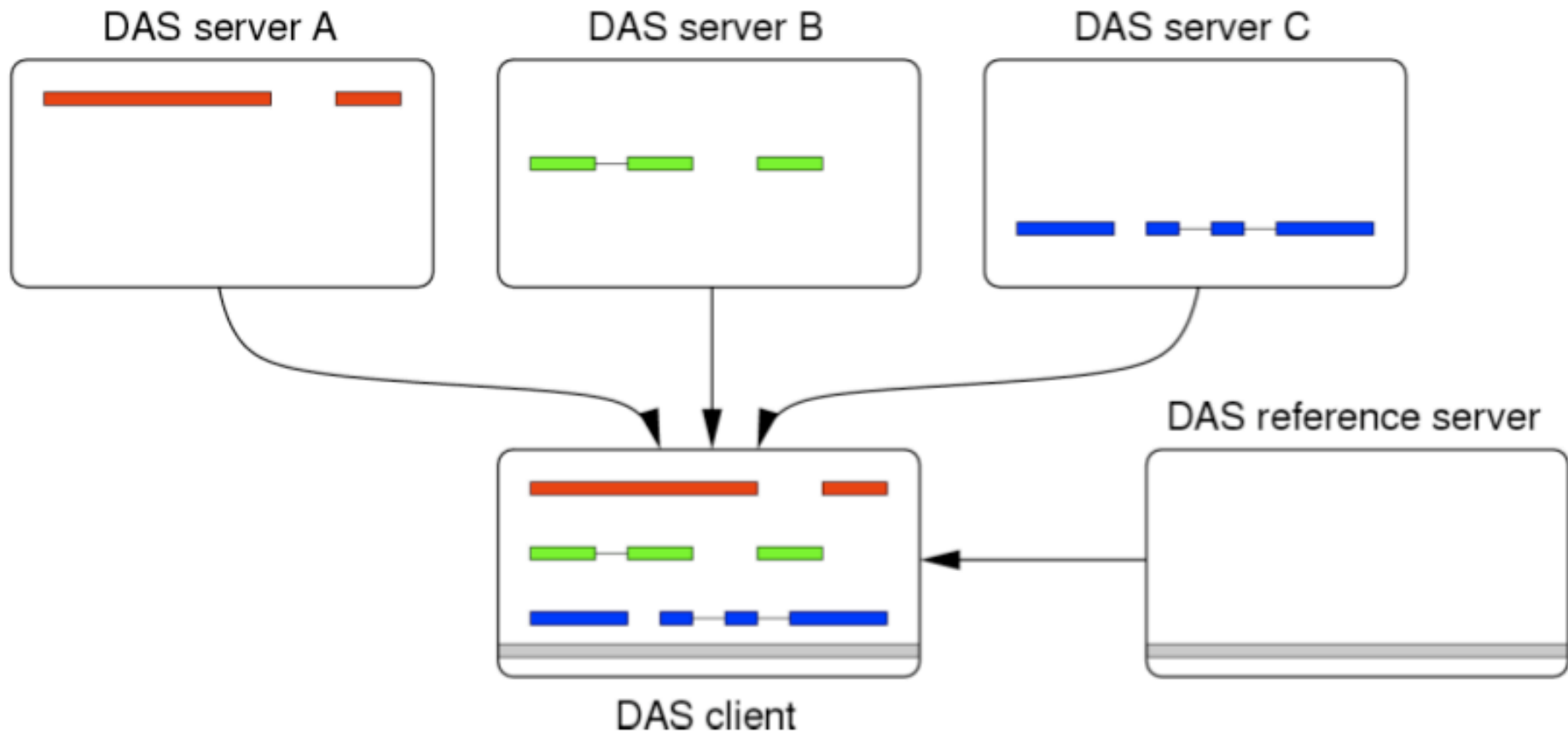
PRIDE: why is it there?



- Repository to support publications (proteomics MS derived data)
- **Source of proteomics data for other data resources**



Distributed Annotation System



(<http://www.ebi.ac.uk/dasty/>)

PRIDE DAS server: Dasty example (1)

SEARCH

Protein ID: Registry label:

"UniProt" protein sequence [coordinate system](#)

Examples: P05067, P03973, P13569, MDM2_MOUSE, BRCA1_HUMAN, ...

External links: [UniProt](#), [Spice](#), [Strap](#)

CHECKING

Annotation servers loaded: 100%

System information:

... Dasty2 finished to sort the graphic by type

FILTERING BY

MANIPULATION OPTIONS (Positional features)

POSITIONAL FEATURES

PROTEIN STRUCTURE

View in a pop-up window

Structure [1xmi]

PDB Region: 429 To:671
Uniprot Region: 429 To:671

| FEATURE TYPE | LABELS | FEATURE ANNOTATIONS | SERVER NAME | EVIDENCE (Category) |
|---------------------------|-----------------|---------------------|---------------|--------------------------|
| ABC transporter | SSF90123, SS... | | ⚠ SUPERFAMILY | miscellaneous |
| family annotation | Cyclic AMP-d... | | ✅ InterPro | inferred from InterPro |
| family annotation | CFTR protein... | | ✅ InterPro | inferred from sequence |
| family annotation | CYSFIBREGLTR... | | ✅ InterPro | inferred from sequence |
| family annotation | cAMP cl chan... | | ✅ InterPro | inferred from sequence |
| family annotation | Cystic fibro... | | ✅ InterPro- | inferred from InterPro |
| Component:Protein | P13569 | | ⚠ SUPERFAMILY | structural |
| O-phosphorylated L-serine | PHOSPHORYLAT... | | ✅ netphos | inferred from electronic |
| O-phosphorylated L-serine | PHOSPHORYLAT... | | ✅ cbs total | inferred from electronic |
| O-phosphorylated L- | PHOSPHORYLAT... | | ✅ netphos | inferred from electronic |
| O-phosphorylated L- | PHOSPHORYLAT... | | ✅ cbs total | inferred from electronic |
| O4'-phosphorylated L- | PHOSPHORYLAT... | | ✅ netphos | inferred from electronic |
| O4'-phosphorylated L- | PHOSPHORYLAT... | | ✅ cbs total | inferred from electronic |
| glycosylated residue | UNIPROTKB P1... | | ✅ uniprot | inferred by curator |
| β-loop containing | SSF52540, SS... | | ⚠ SUPERFAMILY | miscellaneous |
| Peptide | MLHSVLOAPMST... | | ✅ PRIDE | Peptide |
| region | G3DSA:3.40.5... | | ✅ InterPro | inferred from sequence |
| region | PTHR19242 | | ✅ InterPro | inferred from sequence |
| region | SSF52540, SS... | | ✅ InterPro | inferred from sequence |
| polypeptide domain | ABC transmem... | | ✅ uniprot | inferred by curator |
| polypeptide domain | ABC TM 1 | | ✅ InterPro | inferred from sequence |
| polypeptide domain | NACHT:452-47... | | ✅ Prosite | inferred from reviewed |
| polypeptide domain | GUANYLATE KI... | | ✅ Prosite | inferred from reviewed |
| polypeptide domain | ABC TRANSPOR... | | ✅ InterPro | inferred from sequence |

Data sharing requires proper infrastructure

- **Community supported, standardized data formats**
Necessary to allow efficient access to the data
- **Controlled vocabularies (CVs) and ontologies**
To provide unambiguous context and metadata to the actual data, as well as enabling powerful queries to be performed on the data
- **Minimal reporting requirements for specific data types**
Ensures the presence of certain bits of information without which interpretation is ambiguous, hampered or impossible
- **Publicly available, online repositories**
Bioinformatics grew up along side the internet, and this is reflected in the successful online data sharing mechanism already in place in the life sciences. *The repositories should implement the standards, use the CV's and ontologies, and adhere to the minimal requirements.*

Coming soon... support for quantitative data

Details for identification: IPI00068506.1

| Submitted Accession | IPI00068506.1 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
|---------------------------------|--|-------------|----------------------------|---------------------------------|-------------|---------------------------------|---------------|-------------------------------|-------------|-------------------------------|--------------|-------------------------------|---------------|-------------------------------|--------------|--|--------|---|-------|---------|---|---------|-----|---------|------|---------|------|---------|------|---------|------|
| Search Database | IPI_HUMAN | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| Cross-References | <table border="1"> <thead> <tr> <th>Accession</th> <th>Database</th> </tr> </thead> <tbody> <tr> <td>ENST00000222388</td> <td>ENSEMBL</td> </tr> <tr> <td>ENSP00000222388</td> <td>ENSEMBL_HUMAN</td> </tr> <tr> <td>IPI00068506.1</td> <td>IPI</td> </tr> <tr> <td>NP_005683.2</td> <td>REFSEQ</td> </tr> <tr> <td>Q75MJ1.1</td> <td>TREMBL</td> </tr> <tr> <td>Q75MJ1_HUMAN</td> <td>TREMBL</td> </tr> <tr> <td>Q9UG63.1</td> <td>TREMBL</td> </tr> </tbody> </table> | Accession | Database | ENST00000222388 | ENSEMBL | ENSP00000222388 | ENSEMBL_HUMAN | IPI00068506.1 | IPI | NP_005683.2 | REFSEQ | Q75MJ1.1 | TREMBL | Q75MJ1_HUMAN | TREMBL | Q9UG63.1 | TREMBL | <p>These mappings have been obtained using the Protein Identifier Cross Reference (PICR) Service at the EBI. They are based on 100% sequence identity and, as a further requirement where applicable, all submitted peptide sequences must match. Mappings shown in light grey are historical and correspond to inactive entries in the source databases.</p> | | | | | | | | | | | | | |
| Accession | Database | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| ENST00000222388 | ENSEMBL | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| ENSP00000222388 | ENSEMBL_HUMAN | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| IPI00068506.1 | IPI | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| NP_005683.2 | REFSEQ | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| Q75MJ1.1 | TREMBL | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| Q75MJ1_HUMAN | TREMBL | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| Q9UG63.1 | TREMBL | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| Search Engine | Mascot | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| Score | 63.89 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| Threshold | 34.48 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| % Sequence Coverage | - | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| Quantitative Data | <table border="1"> <thead> <tr> <th>Source Name</th> <th>Value</th> </tr> </thead> <tbody> <tr> <td>PRIDE TMT 126</td> <td>6192618.999</td> </tr> <tr> <td>PRIDE TMT 127</td> <td>6454344.492</td> </tr> <tr> <td>PRIDE TMT 128</td> <td>7620348.236</td> </tr> <tr> <td>PRIDE TMT 129</td> <td>1.98224667E7</td> </tr> <tr> <td>PRIDE TMT 130</td> <td>1.659637806E7</td> </tr> <tr> <td>PRIDE TMT 131</td> <td>1.18205045E7</td> </tr> </tbody> </table> | Source Name | Value | PRIDE TMT 126 | 6192618.999 | PRIDE TMT 127 | 6454344.492 | PRIDE TMT 128 | 7620348.236 | PRIDE TMT 129 | 1.98224667E7 | PRIDE TMT 130 | 1.659637806E7 | PRIDE TMT 131 | 1.18205045E7 | <table border="1"> <thead> <tr> <th>Source Name</th> <th>Value</th> </tr> </thead> <tbody> <tr> <td>TMT 126</td> <td>9</td> </tr> <tr> <td>TMT 127</td> <td>9.4</td> </tr> <tr> <td>TMT 128</td> <td>11.1</td> </tr> <tr> <td>TMT 129</td> <td>28.9</td> </tr> <tr> <td>TMT 130</td> <td>24.2</td> </tr> <tr> <td>TMT 131</td> <td>17.3</td> </tr> </tbody> </table> | | Source Name | Value | TMT 126 | 9 | TMT 127 | 9.4 | TMT 128 | 11.1 | TMT 129 | 28.9 | TMT 130 | 24.2 | TMT 131 | 17.3 |
| Source Name | Value | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| PRIDE TMT 126 | 6192618.999 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| PRIDE TMT 127 | 6454344.492 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| PRIDE TMT 128 | 7620348.236 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| PRIDE TMT 129 | 1.98224667E7 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| PRIDE TMT 130 | 1.659637806E7 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| PRIDE TMT 131 | 1.18205045E7 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| Source Name | Value | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| TMT 126 | 9 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| TMT 127 | 9.4 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| TMT 128 | 11.1 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| TMT 129 | 28.9 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| TMT 130 | 24.2 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| TMT 131 | 17.3 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| Additional | <table border="1"> <thead> <tr> <th>Source User</th> <th>Name MascotConfidenceLevel</th> <th>Value</th> </tr> </thead> <tbody> <tr> <td></td> <td></td> <td>95.0</td> </tr> </tbody> </table> | Source User | Name MascotConfidenceLevel | Value | | | 95.0 | | | | | | | | | | | | | | | | | | | | | | | | |
| Source User | Name MascotConfidenceLevel | Value | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| | | 95.0 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| Mapped Protein sequence | <pre> 0001 MPSDLAKKKA AKKKEAAKAR QRPKRGHEEN GDVVTEPQVA EKNEANGRET TEVDLLTKEL EDFEMKKA^{AA} RAVTGV^{LASH} 0080 0081 PNSTDVHIIIN LSLTFHGQEL LSDTKLELNS GRRYGLI^{GLN} GIGK^{SMLLSA} IGKREVP^{IPE} HIDIYHL^{TRE} MP^{PSDKT}PLH 0160 0161 CVMEVD^{TERA} MLEKEAERLA HEDAEC^{EKLM} ELYERLEELD ADKAEMRAS I^{LHGLG}TPA MQRK^{LKDFS} GGWRMR^{VALA} 0240 0241 RALFIRPFML LLDEPTNHL^{LD} LDACVWLEEE LKTFKRILVL VSHSQDFL^{NG} VCTNI^{IHMHN} KKLKY^{YTGNY} DQYV^{KTRLEL} 0320 0321 EENQMKRFHW EQDQIAHMKN YIARF^{HGSA} K^{LARQAQ}SKE KTLQKMMAS^C LTERV^{VSDKT} LSFY^{FPPCGK} IPP^{PVIMVQN} 0400 0401 VSFKYTKDGF CIYNNLEFGI DL^{DTR}VAL^{VG} PNGAGK^{STLL} K^{LLTG}ELLPT DGM^{IRK}SHSV KIGRY^{HQLQ} EQ^{LDL}LSPL 0480 0481 EYMKCYPEI KEKEEMRKII GR^{YGLT}CKQ^Q VSP^{IRN}LSDG QKCRV^{CLAWL} AWQ^{NPHM}LFL DEPT^{NHL}DIE TIDALAD^{AIN} 0560 0561 EFEGGM^{LVS} HDFRLIQ^{VVA} QE^{IWV}CEK^{QT} I^{TKW}PGDILA YKEHL^{KSLV} DEEP^{QLTKRT} HNV^{CTL}TLAS L^{PRP} </pre> | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |

Do you want to know a bit more...?

D736-D742 Nucleic Acids Research, 2010, Vol. 38, Database issue
doi:10.1093/nar/gkp964

Published online 11 November 2009

The Proteomics Identifications database: 2010 update

Juan Antonio Vizcaino¹, Richard Côté¹, Florian Reisinger¹, Harald Barsnes², Joseph M. Foster¹, Jonathan Rameseder^{1,2}, Henning Hermjakob¹ and Lennart Martens^{1,3*}

¹EMBL Outstation, European Bioinformatics Institute, Wellcome Trust Genome Campus, Hinxton, Cambridge, UK, ²Department of Informatics, University of Bergen, Norway and ³Computational and Systems Biology Initiative, Massachusetts Institute of Technology, Cambridge, MA 02139, USA

Received September 6, 2009; Revised October 6, 2009; Accepted October 13, 2009

ABSTRACT

The Proteomics Identifications database (PRIDE, <http://www.ebi.ac.uk/pride/>) at the European Bioinformatics Institute has become one of the main repositories of mass spectrometry-derived proteomics data. For the last 2 years, PRIDE data holdings have grown substantially, comprising 60 different species, more than 2.5 million protein identifications, 11.5 million peptides and over 50 million spectra by September 2009. We here describe several new and improved features in PRIDE, including the revised submission process, which now includes direct submission of fragmentation annotations. Correspondingly, it is now possible to visualize spectrum fragmentation annotations on tandem mass spectra, a key feature for compliance with journal data submission requirements. We also describe recent developments in the PRIDE BioMart interface, which now allows integrative queries that can join PRIDE data to a growing number of biological resources such as Reactome, Ensembl, InterPro and UniProt. This ability to perform extremely powerful across-domain queries will certainly be a cornerstone of future bioinformatics analyses. Finally, we highlight the importance of data sharing in the proteomics field, and the corresponding integration of PRIDE with other databases in the ProteomeXchange consortium.

INTRODUCTION

Mass spectrometry (MS) is currently the most commonly used technology for the identification and quantification of proteins. Like in any other 'omics' field, the amount of data generated by MS-based proteomics has increased exponentially in the last few years, which has prompted the development of several data repositories. The Proteomics Identifications database (PRIDE)

(<http://www.ebi.ac.uk/pride/>) was developed at the European Bioinformatics Institute (EBI), as a repository for the results of MS-based proteomics experiments, allowing data from a vast range of approaches, instruments and analysis platforms to be stored and disseminated in a common structured and queryable format. Originally established as a production service in 2005, PRIDE has previously been described (1–3) along with guidelines on using the database and its associated tools (4,6).

PRIDE does not stand alone in this field, however, as several other proteomics databases have been established over the past few years. GPMDB (7), PeptideAtlas (8) and Proteopedia (9) are among the most important representatives of these (10). Additionally, the Tranche system (<http://tranche.proteomcommons.org/>) provides a data transfer layer relying on peer-to-peer Internet protocol technology. Finally, the most recently launched proteomics repository is the NCBI PeptideMine (11), a centralized, public proteomics data repository not dissimilar from PRIDE in its objectives. For an up-to-date review covering the capabilities of a comprehensive selection of proteomics MS repositories see Mead *et al.* (12).

PRIDE stores three different kinds of information: MS and MS/MS mass spectra as peak lists, the derived peptide and protein identifications (IDs) and any associated metadata. Indeed, one of the advantages PRIDE offers over other proteomics databases lies in the amount of structured metadata it contains, which is a key requirement to put the stored data in biological and/or technical context. Furthermore, together with the newly released NCBI PeptideMine, the established PRIDE database constitutes an actual structured data repository, and does not assume any editorial control over submitted data.

Another important feature of PRIDE is that it allows data to remain private while anonymously sharing it with journal editors and reviewers through so-called 'reviewer log-in accounts'. As a result, PRIDE is now the recommended submission point for proteomics data for several journals such as *Nature Biotechnology* (13), *Nature*

CORRESPONDENCE

PRIDE Converter: making proteomics data-sharing easy

To the Editor:

Your editorial on 'Democratizing Proteomics Data' correctly addressed the increasing importance of making proteomics data publicly available so that it can be audited, reanalyzed or reused. To make global data-sharing in the field work, however, it is important to minimize the burden of uploading data into publicly available databases, such as PRIDE¹. To this end, we have written a freely available, open source tool called PRIDE Converter that makes it straightforward to submit proteomics data to PRIDE from most common data formats.

Public availability of data is the standard *modus operandi* for most of the life sciences, ranging from genome sequences, over microarray data, to protein information. Some of the best known examples in the field of proteomics include protein sequences in UniProt (<http://www.uniprot.org/>), protein structures in the Protein Databank (<http://www.rcsb.org/>) and protein modifications in UniMod and RESID (<http://www.unimod.org/> and <http://www.ebi.ac.uk/RESID/>). As highlighted in your 2007 editorial¹, making data publicly available in a standardized and structured way enables other researchers to access and reanalyze the data, and to use the collected results in novel ways.

Indeed, much of the progress over the past years in emerging fields, such as mass spectrometry (MS)-based proteomics, is directly related to the public availability of data obtained in earlier efforts², specifically the genome sequencing projects. Not surprisingly, the need for data-sharing in the field of proteomics itself was quickly pointed out³. Several proteomics MS data repositories have since been established, with GPMDB, PRIDE, PeptideAtlas and Proteopedia among the most prominent⁴. With this infrastructure in place, journals have followed suit by starting to request deposition of MS-related data in these databases^{5,6}.

The PRIDE repository at the European Bioinformatics Institute (<http://www.ebi.ac.uk/pride/>) occupies a special place in the list of proteomics databases, in that it constitutes an actual data repository and

does not assume editorial control over submitted data¹. Additionally, it provides a simple yet powerful infrastructure to support anonymous peer review of submitted data while maintaining the submission as private in the system⁷. The PRIDE database has so far accumulated data on more than 9,500 experiments, collectively containing more than 40 million mass spectra, identifying well over 1.4 million unique peptide sequences, which in turn infer more than 100,000 unique Ensembl proteins across all species. Submitting an MS-based proteomics data set to a structured repository, such as PRIDE, has many advantages over alternative ways of making peptide and protein identifications publicly available, such as uploading raw data files on a web page⁸ or submitting text or PDF files as supplementary information to a journal⁹. Furthermore, centralized repositories can also offer additional services and tools to the scientific community, based on uploaded data. PRIDE for instance includes tools for (i) visualizing protein coverage, peptide modifications and spectrum annotations, (ii) automatic

mapping of protein accession numbers to identifiers from all other commonly used proteomics databases using the PICR service¹⁰ and (iii) comprehensive protein list comparisons (through Venn diagrams)¹¹.

Submitting data to PRIDE could be challenging for some users, however. PRIDE relies on an XML-based data format for submissions, which is built around the Proteomics Standards Initiative mzData standard for mass spectrometry (<http://www.ebi.ac.uk/pride/schema/mzDataDocumentation.do>)¹². And although the PRIDE XML format is well documented, converting proteomics data to PRIDE XML could present difficulties, especially for wet-lab scientists without a strong bioinformatics background or informatics support. To alleviate this problem, two tools for converting data into PRIDE XML have already been developed: the ProteomeHarvest PRIDE Submission Spreadsheet, which is a Microsoft Excel-based (<http://www.ebi.ac.uk/pride/protomeharvest/>), and the PRIDE Wizard for Macosx result files (<http://www.ebi.ac.uk/pride/>).

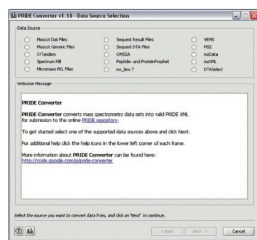


Figure 1 PRIDE Converter opening screen, showing the available input formats.

Proteomics 2009, 9, 1–8

DOI 10.1002/pmic.200900402

1

STANDARDISATION & GUIDELINES

A guide to the Proteomics Identifications Database proteomics data repository

Juan Antonio Vizcaino¹, Richard Côté¹, Florian Reisinger¹, Joseph M. Foster¹, Michael Mueller¹, Jonathan Rameseder^{1,2}, Henning Hermjakob^{1,2} and Lennart Martens¹

¹EMBL Outstation, European Bioinformatics Institute, Wellcome Trust Genome Campus, Hinxton, Cambridge, UK, ²Computational and Systems Biology Initiative, Massachusetts Institute of Technology, Cambridge, MA, USA

The Proteomics Identifications Database (PRIDE, www.ebi.ac.uk/pride/) is one of the main repositories of MS derived proteomics data. Here, we point out the main functionalities of PRIDE both as a submission proteomics data and as a source for proteomics data. We describe the main features for data retrieval and visualization available through the PRIDE web and BioMart interfaces. We also highlight the mechanism by which tailored queries in the BioMart can join PRIDE to other resources such as Reactome, Ensembl or UniProt to execute extremely powerful across-domain queries. We then present the latest improvements in the PRIDE submission process, using the new easy-to-use, platform-independent graphical user interface submission tool PRIDE Converter. Finally, we speak about future plans and the role of PRIDE in the ProteomeXchange consortium.

Received: June 9, 2009
Revised: June 24, 2009
Accepted: June 25, 2009

Keywords:
Bioinformatics / Data repository / Mass spectrometry

1 Introduction

Bioinformatics tools and data repositories provide one of the main pillars of biology in the 21st century. Indeed, public availability of biological data via the Internet has changed the way biologists plan, execute and interpret their studies. Some of the best known protein-related resources include UniProt [1] for protein sequences and annotation, the Protein Databank [2] and other members of the wwPDB consortium [3] for protein structures, Interact [4] and other components of the IMEx consortium [5] for protein interactions, InterPro [6] for protein domains, and UniMod [7] and RESID [8] for protein modifications.

Like in any other 'omics' field, the amount of data generated by MS based proteomics has increased exponentially in the last few years, which prompted the development of several data repositories. At the same time, proteomics efforts have driven the development of universally adopted and stable data formats under the auspices of the HUPO Proteomics Standards Initiative (HUPO-PSI, <http://www.psidi.info/>) and have led to powerful data analysis strategies [9, 10]. Taken together, these advances have allowed the centralized aggregation of proteomics data and its reanalysis or meta-analysis, ultimately turning proteomics into a much more robust discipline in the life sciences. Several proteomics MS data repositories have been established so far, with GPMDB [11], Proteomics Identifications Database (PRIDE) [12], PeptideAtlas [13] and Proteopedia [14] among the most prominent ones at present [15, 16]. Additionally, the NCBI recently launched their PeptideMine (<http://www.ncbi.nlm.nih.gov/projects/peptideMine/>) as a centralized, public proteomics repository not dissimilar from PRIDE. The Tranche (<http://tranche.proteomcommons.org/>) system is used in the field as well, and essentially presents a data transfer layer relying on peer-to-peer Internet protocol technology. Apart from these large-scale efforts, there are also smaller, more

Correspondence: Dr. Lennart Martens, EMBL Outstation, European Bioinformatics Institute, Wellcome Trust Genome Campus, Hinxton, Cambridge, UK
E-mail: lennart.martens@ebi.ac.uk
Fax: +44-1223-494-484

Abbreviations: EBI, European Bioinformatics Institute; OLS, Ontology Lookup Service; PICR, Protein Identifier Cross-Referencing; PRIDE, Proteomics Identifications Database; PSI, Proteomics Standards Initiative

www.proteomics-journal.com
Weinheim

© 2009 WILEY-VCH Verlag GmbH & Co. KGaA, Weinheim

*To whom correspondence should be addressed. Tel: +44 1223 492 610; Fax: +44 1223 494 484; Email: lennart.martens@ebi.ac.uk

© The Authors 2009. Published by Oxford University Press.

This is an Open Access article distributed under the terms of the Creative Commons Attribution Non-Commercial License (<http://creativecommons.org/licenses/by-nc/2.5.uk/>) which permits unrestricted non-commercial use, distribution, and reproduction in any medium, provided the original work is properly cited.



© 2009 Nature America, Inc. All rights reserved.

598

VOLUME 27 NUMBER 7 JULY 2009 NATURE BIOTECHNOLOGY

Juan A. Vizcaino
juan@ebi.ac.uk



BSPR/EBI Educational Workshop
Hinxton, 16 July 2010



The PRIDE Team

Attila Csordas



Daniel Ríos



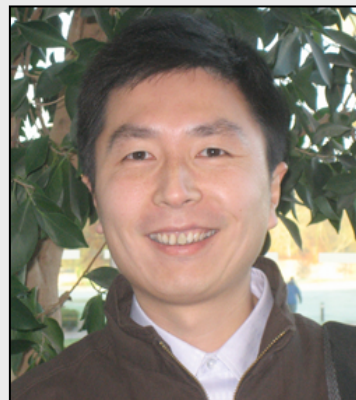
Richard Côté



Florian Reisinger



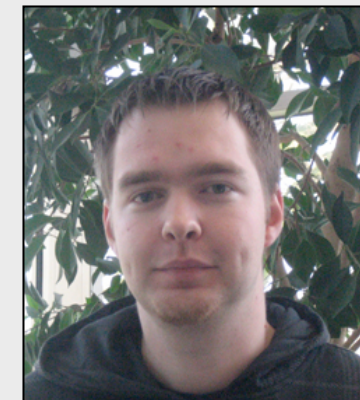
Henning Hermjakob



Rui Wang



Joe Foster
(Ph.D. student)



Andreas Schonegger
(Trainee)

Links, collaborations and funding

<http://www.psidev.info>

<http://www.ebi.ac.uk/ols>

<http://www.ebi.ac.uk/pride>

<http://www.ebi.ac.uk/tools/picr>

<http://www.ebi.ac.uk/pride/dod>

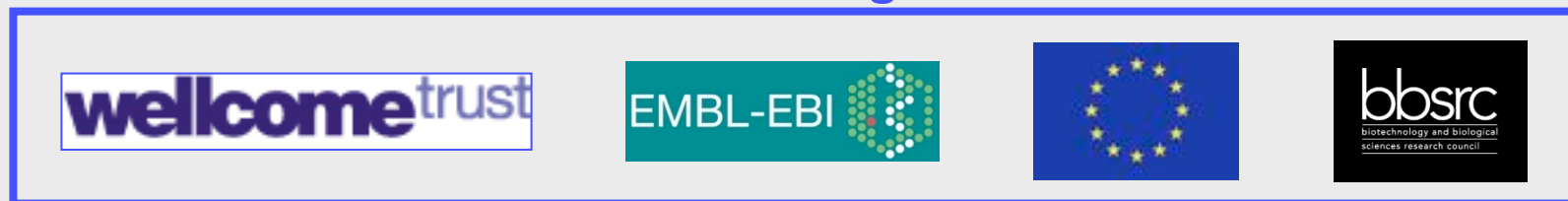
juan@ebi.ac.uk

pride-support@ebi.ac.uk

PRIDE collaborators



Funding



Juan A. Vizcaíno
juan@ebi.ac.uk



BSPR/EBI Educational Workshop
Hinxton, 16 July 2010



Thank you!

Questions?