

# Tutorial: Submission of MS/MS datasets to ProteomeXchange via PRIDE

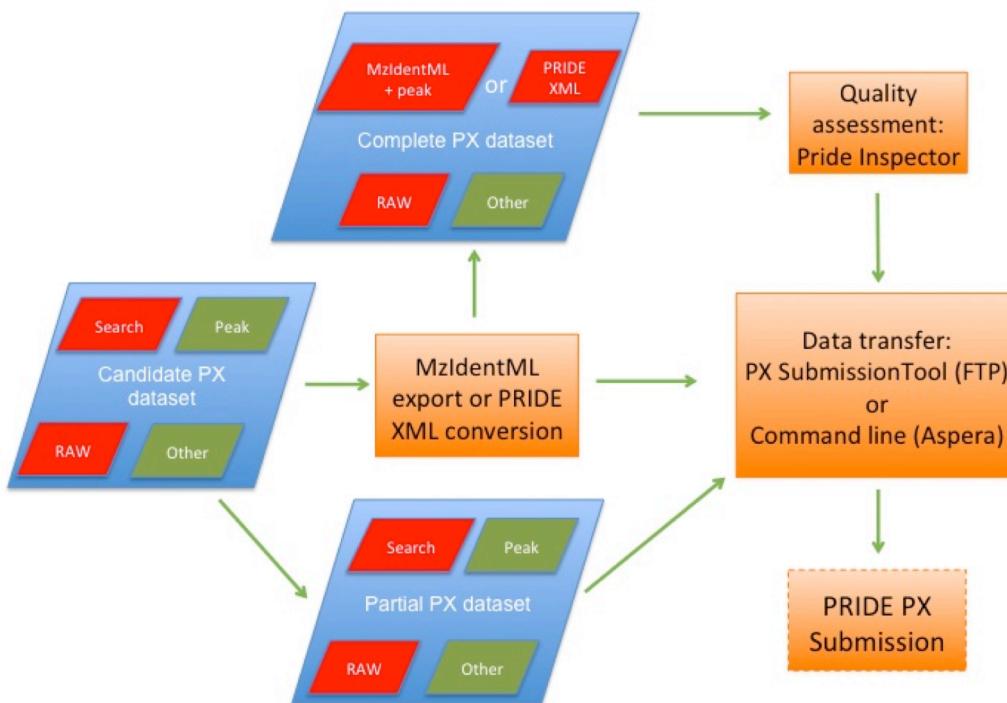
---

<b>1 Where do I start? Submission summary overview .....</b>	<b>3</b>
<b>2 Submission types: Complete and Partial Submissions .....</b>	<b>4</b>
2.1 Complete Submission .....	4
2.2 Partial Submissions .....	6
<b>3 Bulk Submissions.....</b>	<b>7</b>
<b>4 How to make complete submissions? .....</b>	<b>7</b>
<b>5 How to make Partial Submissions? .....</b>	<b>24</b>
<b>6 How to make bulk submissions?.....</b>	<b>36</b>
6.1 Creation of the PX Submission Summary File.....	36
6.2 Submission using the PX Submission tool .....	37
6.3 Fast upload option for big datasets using Aspera .....	38
<b>7 What happens after the submitter has uploaded all the data? .....</b>	<b>40</b>
<b>8 Accessing Private Data .....</b>	<b>40</b>
8.1 PRIDE Archive web page.....	40
8.2 PRIDE Inspector.....	41
<b>9 Post-submission steps .....</b>	<b>42</b>
9.1 How to do a resubmission of a dataset? .....	42
<i>9.1.1 Resubmission with the PX Submission Tool.....</i>	<i>42</i>
<i>9.1.2 Resubmission via Aspera.....</i>	<i>44</i>
9.2 Referencing the dataset in the paper .....	45
9.3 Public release of the dataset .....	45
<b>10 Appendix I: Definitions .....</b>	<b>47</b>
<b>11 Appendix II: Available tools to help you with the submission .....</b>	<b>50</b>
11.1 Creation of PRIDE XML files .....	50
<i>11.1.1 Tools developed by the PRIDE team.....</i>	<i>50</i>
<i>11.1.2 External tools developed by collaborators.....</i>	<i>50</i>
11.2 Creation of mzIdentML files.....	51
11.3 Checking the files before submission (initial quality assessment) .....	52
<i>11.3.1 Tool developed by the PRIDE team.....</i>	<i>52</i>
<i>11.3.2 External tool developed by collaborators.....</i>	<i>52</i>
11.4 File submission to PRIDE .....	52
<b>12 Appendix III: Summary of formats supported by PRIDE for PX MS/MS submissions .....</b>	<b>53</b>
<b>13 Appendix IV: Metadata requirements for MS/MS submissions .....</b>	<b>56</b>
<b>14 Appendix V: Recommended Partial Submission search engine identification results for particular software tools .....</b>	<b>58</b>
14.1 MaxQuant .....	58
14.2 ProteinPilot .....	58



## 1 Where do I start? Submission summary overview

The default PRIDE submission consists of the deposition of MS/MS proteomics datasets according to the guidelines of the ProteomeXchange (PX) consortium. Figure 1 shows the overall submission process (by December 2013).



**Figure 1:** Overview of the PRIDE/PX submission process with the two default submission types.

Each submitted PX dataset will contain:

- peptide/protein identification files (called ‘RESULT’),
- mass spectrometer output files (called ‘RAW’), which are either machine raw files or not heavily processed files in a XML standard format such as mzXML or mzML,
- optionally other files can be included like peak list files (called ‘PEAK’), search engine output files (called ‘SEARCH’, mandatory for “Partial submissions”, see below), quantification results (‘QUANT’), gel images (‘GEL’), and any other file types (‘OTHER’).

There are two different submission workflows depending on whether peptide/protein identification results can be submitted in a standard format that can be handled by PRIDE (both PRIDE XML and mzIdentML version 1.1 are now supported) or not. If PRIDE XML ‘RESULT’ files or mzIdentML plus the accompanying peak list/XML-based files containing the referenced spectra are provided, the ‘Complete’ Submission option is available. If ‘RESULT’ files are not

available in these formats, a ‘Partial’ Submission can be done, but only in case search engine output files/identification results can be provided (see next section).

It is important to highlight that the current version of pipeline does not support a full and standard representation of the quantification results, linked to the identification results (unless this information is provided in the PRIDE XML files). It is expected that data standards for quantitative proteomics data (mzQuantML, mzTab) will be supported in the future. However, any quantification result output files can be submitted as accompanying ‘QUANT’ files.

Before a submission is started it is necessary to have a PRIDE user account (please register at <http://www.ebi.ac.uk/pride/archive/register>). All submissions to ProteomeXchange are private by default, and the username and password are needed to access your data. Data will be made publicly available when the submitter notify us to do it or by default when the corresponding manuscript is made available (see Section 9.3).

## 2 Submission types: Complete and Partial Submissions

As summarized above, two main submission types/workflows are available: ‘Complete’ or ‘Partial’ Submissions. For all types of submissions to PX via PRIDE, the first option for the users is to use the Java stand-alone tool “PX Submission tool” (available at <http://www.proteomexchange.org/submission>).

### 2.1 Complete Submission

This is the recommended and preferred option. ‘RAW’ files need to be provided together with the ‘RESULT’ type supported file formats PRIDE XML or mzIdentML (version 1.1) files [22375074]. These are the two subtypes of ‘Complete’ submissions. Uploading peak list (‘PEAK’), search engine output (‘SEARCH’), quantification and other post processing files (‘OTHER’) can also be done in order to give a near complete coverage and representation of your data and it is recommended but not enforced. However, if the submitter chooses to submit the ‘RESULT’ files as mzIdentML, the corresponding peak list files (‘PEAK’) used in the search and referenced in the mzIdentML file/s need to be submitted as well. The reason behind is that otherwise, the mass spectra will not be submitted since mzIdentML, unlike PRIDE XML, only contains the peptide/protein identification results.

After the submission, you will be issued with not only a ProteomeXchange accession number but also with a permanent DOI (Digital Object Identifier) to uniquely identify your submission in the future.

Your submitted data will be fully accessible in PRIDE and allow full visualization of the data for private journal review support using PRIDE Inspector. Your data will be made available via FTP (<ftp://ftp.pride.ebi.ac.uk/>) to download once it has been made public.

The complete submission requires at least two sets of files in case of PRIDE XML based submissions, and three in case of mzIdentML based submissions:

- **Result files fully supported by PRIDE** (called ‘RESULT’): Two formats are currently supported:
  - PRIDE XML files, which must contain both the mass spectra and the identifications (see definitions, Appendix I). Many of the most popular search engine output files can be converted to PRIDE XML using the tool [PRIDE Converter 2](#). However, PRIDE XML files can also be produced by other tools (see Appendix II) and/or external pipelines.
  - mzIdentML version 1.1 files. mzIdentML is the Proteomics Standards Initiative (PSI) standard for peptide/protein identification data [22375074]. Many of the most popular search engine output files can be exported to mzIdentML 1.1 (see Appendix II or <http://www.psidev.info/tools-implementing-mzidentml>). Since the MS data is not included in mzIdentML, to have a complete submission it is also mandatory to submit the corresponding peak list files (‘PEAK’, see below). mzIdentML 1.0 files (the non-stable version of the standard) are not supported.

In both cases, in the PX Submission Tool both types of files should be tagged as ‘RESULT’ (for a comprehensive list of the formats supported by PRIDE, see Appendix III).

- **Mass spectrometer output files** (called ‘RAW’): Two options are possible: MS instrument binary output files, such as BRUKER .baf files, or not heavily processed files in XML format like mzXML or mzML files (see definitions, Appendix I). If your ‘RAW’ files are organized in directories instead of individual files, please compress them into one individual file (for instance to .zip) before upload. In the submission tool they should be tagged as ‘RAW’.
- **Peak list files** (called ‘PEAK’, only mandatory for mzIdentML ‘RESULT’ files, optional for PRIDE XML ‘RESULT’ files): The user can provide the exact version of the files that was used by the search engine to generate the experimental results, the ones that are referenced from the original mzIdentML files. In the submission tool they should be tagged as ‘PEAK’. Otherwise, it would be impossible to link the identifications to the corresponding spectra.

Although not required, other types of files can be submitted optionally:

- **Search engine output files** (called ‘SEARCH’): the original output files from your search engine or your analysis pipeline, such as Trans-Proteomic Pipeline (TPP) pep.xml and/or prot.xml files, or MaxQuant text output files, among many others. They should contain the peptide/protein identifications. In the submission tool they should be tagged as ‘SEARCH’.
- **Quantification related files**: In the PX Submission Tool they should be tagged as ‘QUANT’.

- **Gel images files:** In the PX Submission Tool they should be tagged as 'GEL'.
- **Any other files:** In the PX Submission Tool they should be tagged as 'OTHER'.

If the PX Submission Tool is not used to perform the submission, an extra file is needed. The file is generated automatically and submitted by the PX submission tool, so it does not need to be created independently if the PX Submission Tool is used.

- **PX submission summary file:** This file captures the descriptive information about a ProteomeXchange submission, such as: experimental metadata, included files, file mappings, etc. All the details about the data format can be found [here](#).

## 2.2 Partial Submissions

You should only choose this option if your search results cannot be converted/exported to PRIDE XML or mzIdentML v1.1 (plus the accompanying spectra). It is not the recommended option, since it will significantly reduce the reusability of your dataset. 'RAW' files need to be provided together with search engine output files ('SEARCH'). Uploading peak list ('PEAK'), quantification and other types of files ('QUANT', 'GEL' or 'OTHER') is possible but not enforced. As a result, you will be issued with a ProteomeXchange accession number but not with a DOI. Once it is made public, your dataset will be available to download via FTP.

The partial submission requires two sets of files:

- **Search engine result files:** (called 'SEARCH'): the original output files from your search engine or your analysis pipeline, Trans-Proteomic Pipeline (TPP) pep.xml and/or prot.xml files, or MaxQuant text output files, among many others. They should contain the peptide/protein identifications. In the submission tool they should be tagged as 'SEARCH'.
- **Mass spectrometer output files** (called 'RAW'): MS instrument binary output files, such as BRUKER .baf files or not heavily processed mzXML or mzML files (see definitions, Appendix I). If your 'RAW' files are organized in directories instead of individual files, please compress them into one individual file (for instance to .zip) before upload. In the submission tool they should be tagged as 'RAW'.

Again, although not required, other types of files can be submitted optionally:

- **Peak list files:** It is strongly recommended to provide the peak list files (eg. mgf files) that were used for the original search and these are different from the provided mandatory raw files. In the submission tool they should be tagged as 'PEAK'.

- **Quantification related files:** In the PX Submission Tool they should be tagged as ‘QUANT’.
- **Gel images files:** In the PX Submission Tool they should be tagged as ‘GEL’.
- **Any other files:** In the PX Submission Tool they should be tagged as ‘OTHER’.

As explained earlier, if the PX Submission Tool is not used to perform the submission, an extra file is needed. The file is generated automatically and submitted by the PX submission tool, so it does not need to be created independently if the PX Submission Tool is used.

- **PX submission summary file:** This file captures the descriptive information about a ProteomeXchange submission, such as: experimental metadata, included files, file mappings, etc. All the details about the data format can be found [here](#).

### 3 Bulk Submissions

Independently from being complete or partial, you can make a ‘**Bulk Submission**’ if you need to submit a large set of files. This path is envisioned for labs with some bioinformatics support since some scripting work is needed. **Both ‘Complete’ and ‘Partial’ Submissions can be performed through this mechanism.**

The “bulk submission” requires also two sets of information:

- **Experiment data files:** The files you want to submit to PRIDE via ProteomeXchange. See section 2 for the exact files needed for each submission type (either ‘Complete’ or ‘Partial’).
- **PX submission summary file:** Needed if the submission is not performed using the PX submission tool. This file captures the descriptive information about a ProteomeXchange submission, such as: experimental metadata, included files, file mappings, etc. All the details about the data format can be found [here](#). The file is automatically created by the ‘PX submission tool’.

### 4 How to make complete submissions?

As discussed earlier in Section 2.1 the two subtypes of ‘Complete’ submissions are either mzIdentML or PRIDE XML based. ‘Complete’ submissions mixing the two types of ‘RESULT’ files are not allowed.

Many of the submission steps are identical for the two subtypes so these steps are going to be discussed in a uniform manner. The differences will be

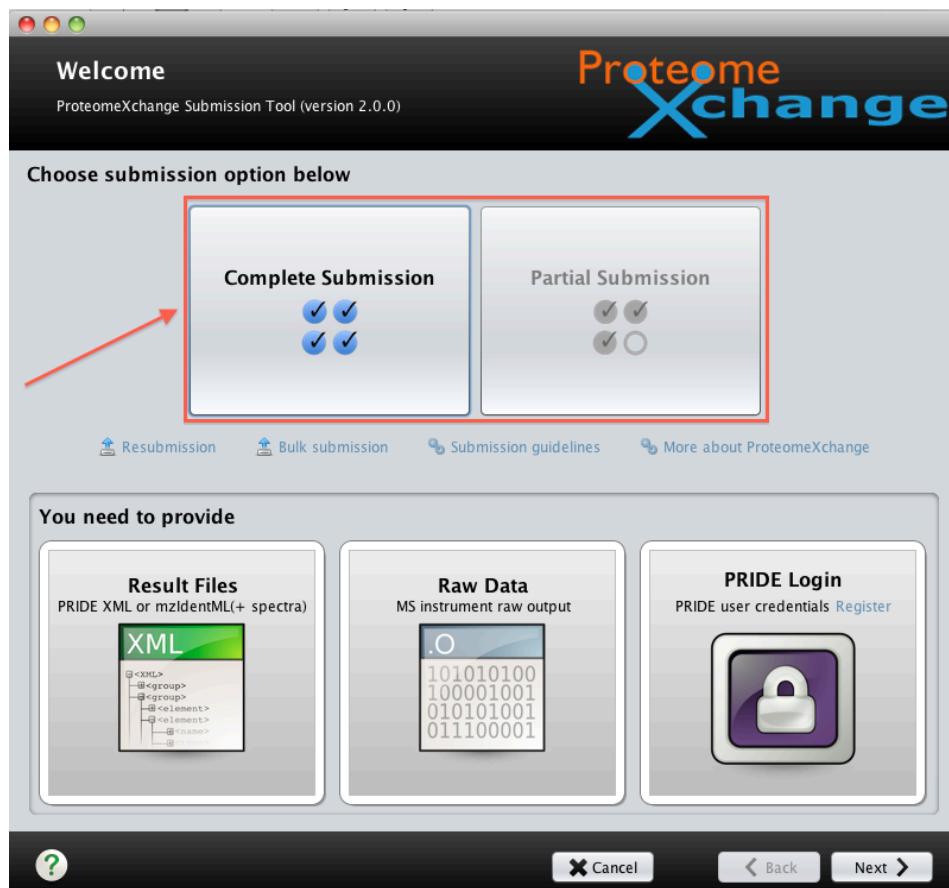
highlighted in case of those steps that are different. The different steps are the following: **Step 5**: ‘Add Files and assign file types’, and **Step 6**: ‘Assign relationships between the submitted files’.

### **Step 1: Launch PX Submission Tool**

First you need to install and launch the PX Submission Tool (available at <http://www.proteomexchange.org/submission>).

### **Step 2: Select Submission Type**

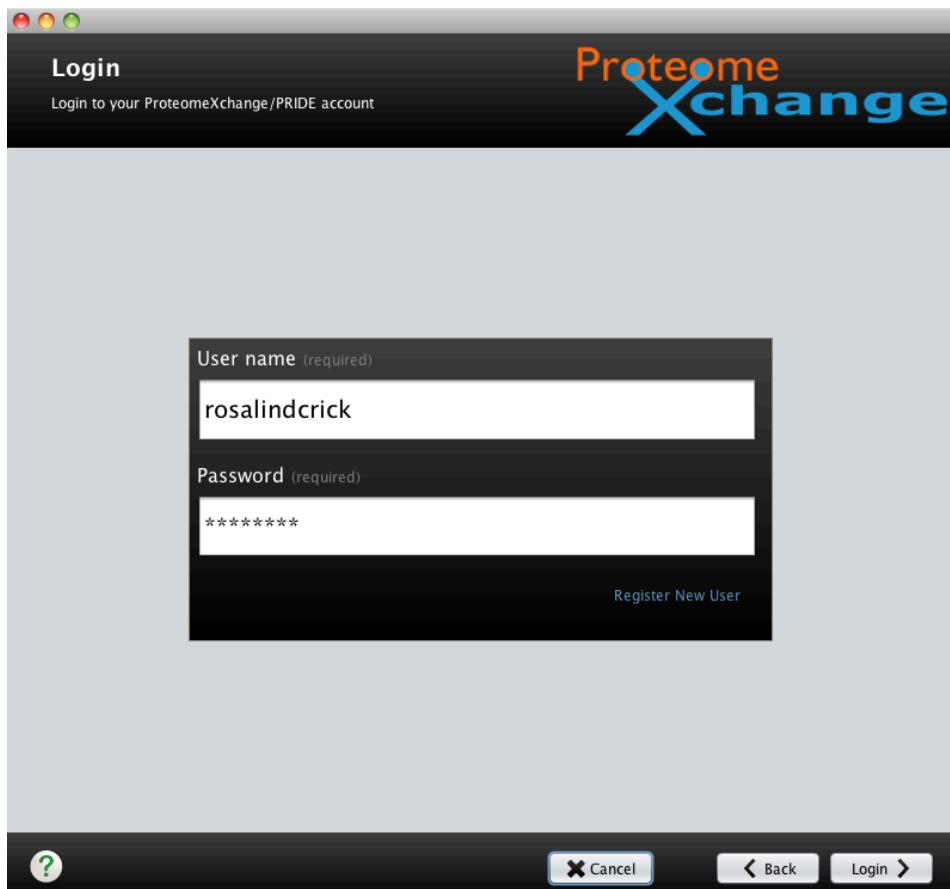
You then need to select ‘Complete Submission’ in the PX Submission Tool ‘Welcome’ screen (Figure 2).



**Figure 2:** ‘Welcome’ screen of the PX Submission Tool showing the two submission types.

### **Step 3: Login**

Please log in using your existing PRIDE account as shown in Figure 3:



**Figure 3** Login screen of the PX Submission tool.

#### **Step 4: Provide submission details**

The user is asked to provide some basic details about the uploaded dataset (Figure 4) such as the title, a list of keywords (in a comma separated format), and a brief description of the data (similar to the abstract of the corresponding publication) a sample processing and a data processing protocol. The user also picks a mass spectrometry experiment type from a drop-down menu.

The screenshot shows the 'Dataset Details' screen of the PX Submission tool. At the top right is the ProteomeXchange logo. A tip at the top right says 'Tip: Use Ctrl+C to copy, Ctrl+V to paste'. The form fields include:

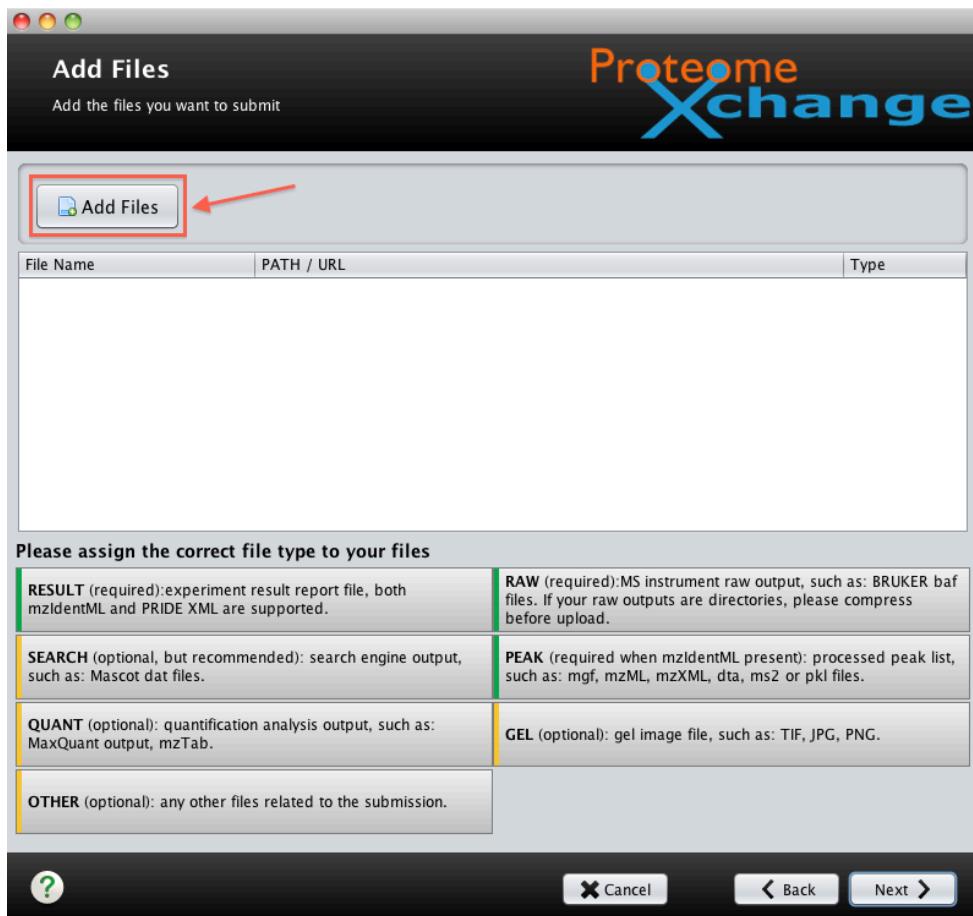
- Project title\***: i.e. Human liver LC-MSMS
- Keywords\***: i.e. Human, Liver, Plasma, LC-MSMS
- Project description\*** (50 to 5000 characters): Please provide an overall description of your study, think something similar in scope to the manuscript abstract.
- Sample processing protocol\*** (50 to 5000 characters): Please provide a short description on the sample preparation steps, separation, enrichment strategies and mass spectrometry protocols included.
- Data processing protocol\*** (50 to 5000 characters): Please provide a couple of sentences on the bioinformatics pipeline used, main search parameters, quantitative analysis, software tools and versions included. Think something similar in scope to the Data Analysis section of your manuscript.
- Experiment type\***: A dropdown menu showing options like 'Choose experiment type here', 'Shotgun proteomics', 'Cross-linking (CX-MS)', 'Affinity purification (AP-MS)', 'SRM/MMR', 'SWATH MS', and 'Other experiment type...'. The 'Choose experiment type here' option is currently selected.

At the bottom are buttons for 'Cancel', 'Back', and 'Next >'. A question mark icon is also present.

**Figure 4:** 'Dataset details' screen in the PX Submission tool.

### **Step 5: Add Files and assign file types**

In this stage, you should choose the files you would like submit. As shown in Figure 5, you can add files by clicking on the highlighted button.

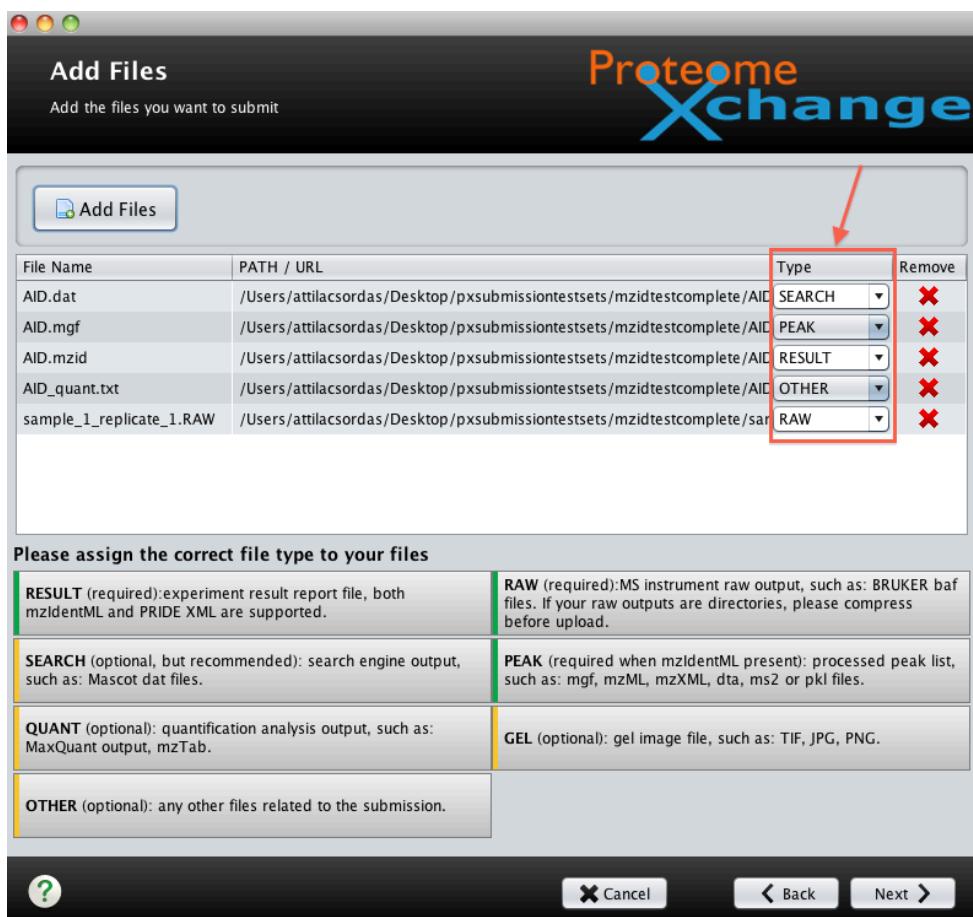


**Figure 5:** 'Add files' screen of the PX submission tool.

There are slight differences in this step between the two subtypes of submissions so we will discuss them separately.

### Step 5A: mzIdentML files

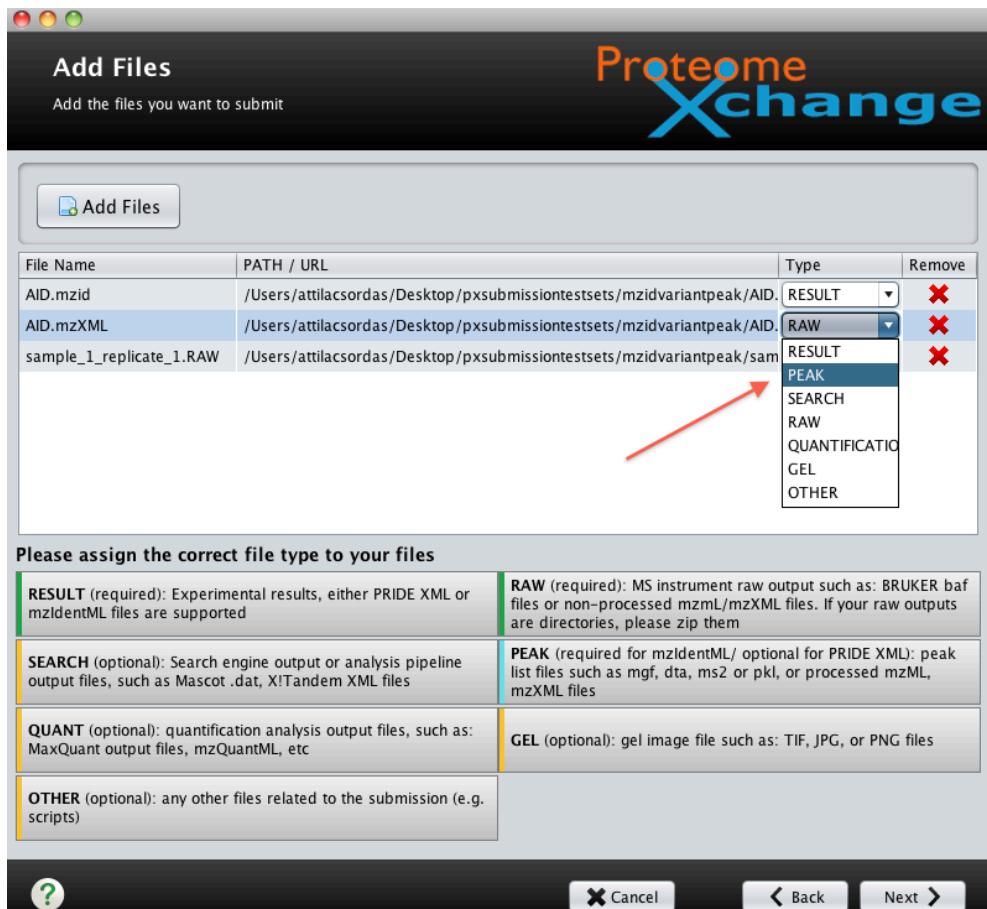
You have to make sure that at least 'RESULT' files, 'RAW' files and 'PEAK' files are selected. There could also be other file types included in the submission: 'SEARCH' (for search engine output files in case those were not mzIdentML files natively), 'QUANT', for quantification results, 'GEL', for gel images, or 'OTHER' (any other file eg. database fasta files, protein inference, post-search files). All the files need to be selected at this stage. Once they are added, double-check that they were assigned with the correct file type, as shown in Figure 6.



**Figure 6:** Adding files in case of an mzIdentML based ‘Complete’ submission: Assignment of the correct file types.

In the case of ‘PEAK’ files, the tool will check and validate that all the required file(s) that were referenced in the mzIdentML file’s <SpectraData> element are present. If your peak list files had an extension recognized by the tool (.mgf, .dta, .ms2, .pkl) then the tool will automatically annotate the type as ‘PEAK’ (see Figure 6) but in other cases you have to assign the file type yourself. For instance if the mzIdentML file references .mzXML files, the tool will recognize them as ‘RAW’ files, since they can be used as ‘RAW’ file replacements as well. In that case you have to change the file type manually and switch from ‘RAW’ to ‘PEAK’ (see Figure 7). The same applies if you are using a peak list files format that is not recognized by the tool as a ‘PEAK’ file but as an ‘OTHER’ file.

In case both the referenced ‘PEAK’ files and the ‘RAW’ files are the same files (in a XML-based format) then currently you need to provide them twice, as ‘RAW’ and as ‘PEAK’.

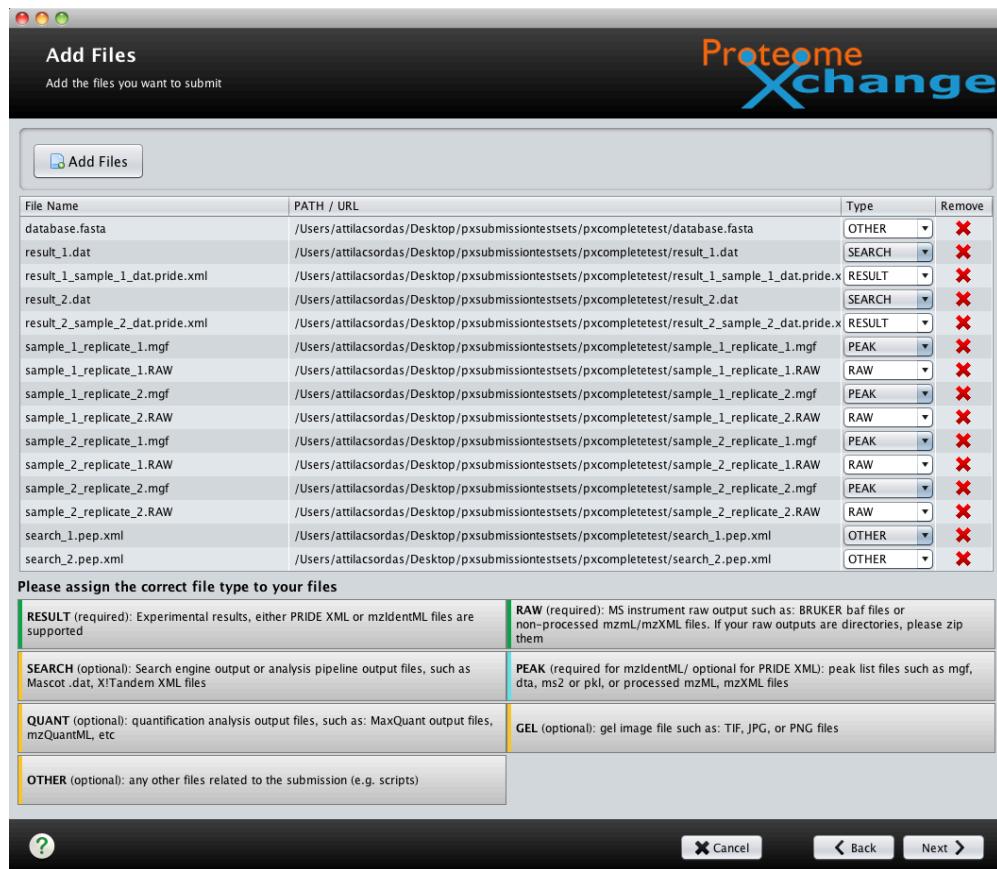


**Figure 7:** Switching the file type to the correct file type in case of an mzIdentML based 'Complete' submission.

## Step 5B: PRIDE XML files

When adding files please make sure that at least 'RESULT' files and the 'RAW' files are selected. The PRIDE XML result files do contain spectra data besides identifications so peak list files are not mandatory as opposed to mzIdentML based 'Complete' submissions. Once the files are added, double-check that they were assigned with the correct file type, as shown in Figure 8.

There could also be other files types included in the submission: 'SEARCH' (for search engine output files), 'PEAK' (for peak list files) or 'OTHER' (e.g. quantification files). All these files need to be selected at this stage.



**Figure 8:** Adding files in case of a PRIDE XML based ‘Complete’ submission: Assignment of the correct file types.

## Step 6: Assign relationships between the submitted files

This mapping step consist of assigning the relations between the ‘RESULT’ files and the other types of files included in the submission, for example, which ‘RAW’ (mandatory), ‘PEAK’ (mandatory for mzIdentML 1.1), ‘SEARCH’, ‘QUANT’, ‘GEL’ or ‘OTHER’ files can be linked to a given ‘RESULT’ file or are associated with it. This will enable others to understand how your data is connected and structured.

By default the tool makes an attempt to generate the mapping between the ‘RESULT’ and the other - most importantly ‘RAW’ - files. For instance if there has been only 1 ‘RESULT’ file found during the previous ‘Add Files’ step (Step 5) then all the other files will be mapped to this ‘RESULT’ file. If there are multiple ‘RESULT’ files the tool maps the other files – ‘RAW’, ‘PEAK’, ‘SEARCH’, ... - with the same file name prefix, but without the file extension, to the corresponding ‘RESULT’ files. This mapping is done even if the suffix part of the ‘RAW’ files contains different numbers (for instance indicating different replicates).

If the automatic mapping is partial only or does not apply, the submitter is asked to manually assign the relationships between the files.

Since there are differences in this step between the two subtypes we are going to discuss them separately.

## Step 6A: mzIdentML files

Each mzIdentML ‘RESULT’ file must have at least two files mapped to it: a ‘RAW’ and a ‘PEAK’ file. Make sure you assign the ‘PEAK’ type to the file(s) containing spectra information and referenced in the corresponding mzIdentML files, as discussed in the previous step (5A).

As shown in Figure 9 the file linking is done by clicking on the ‘Add Relation’ button. Many files can be assigned to the same ‘RESULT’ file.

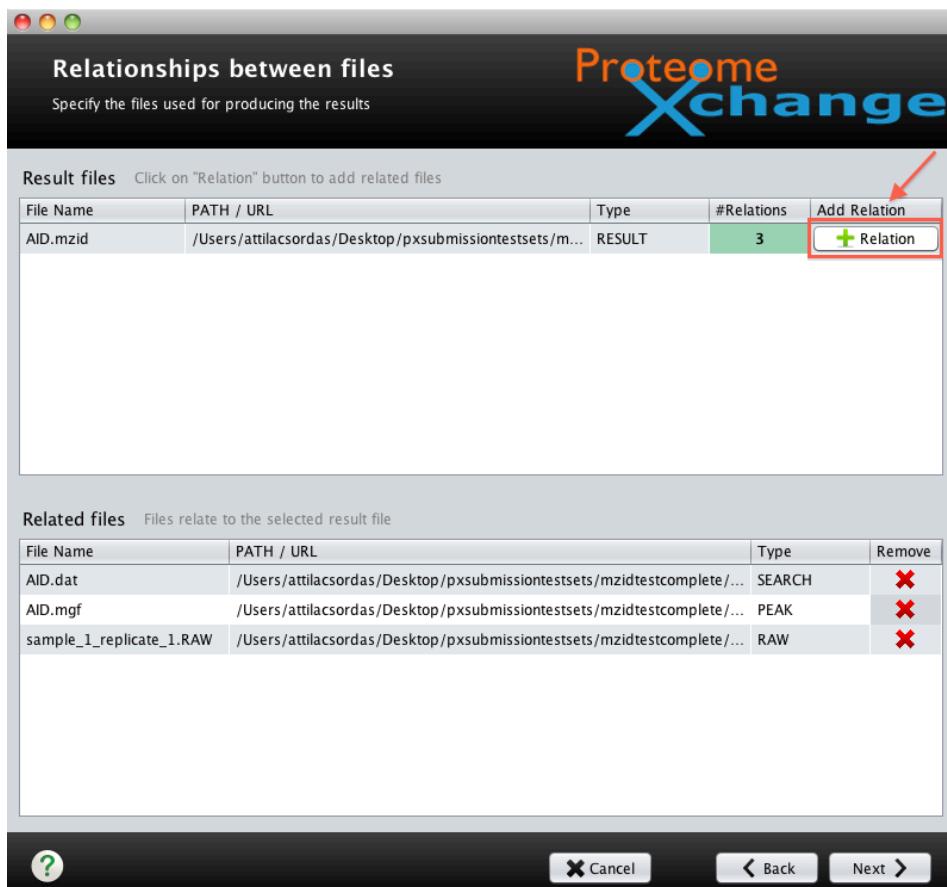
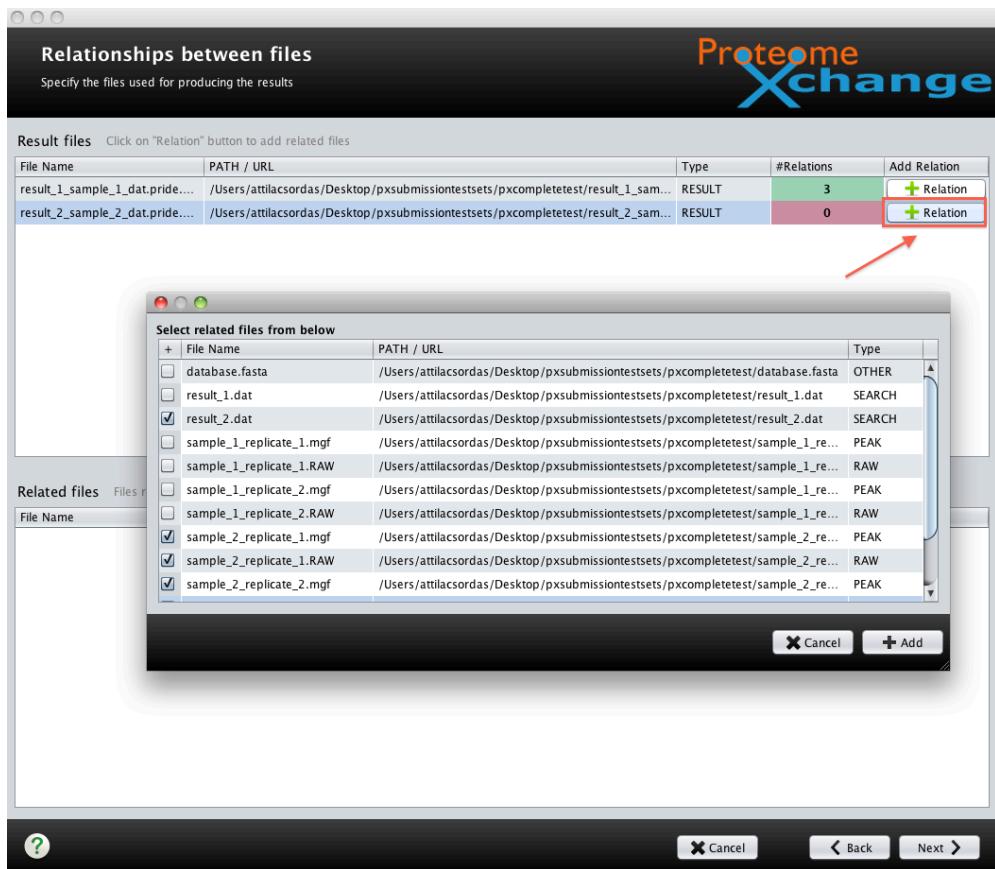


Figure 9: ‘Relationships between files’ screen of the PX Submission tool.

## Step 6B: PRIDE XML files

Each ‘RESULT’ file must have at least one ‘RAW’ file linked to it. Figure 10 shows the situation when ‘SEARCH’, ‘RAW’ and ‘PEAK’ files are added to a PRIDE XML file by clicking on the ‘Add Relation’ button. Different number of files can be assigned to the same ‘RESULT’ file.



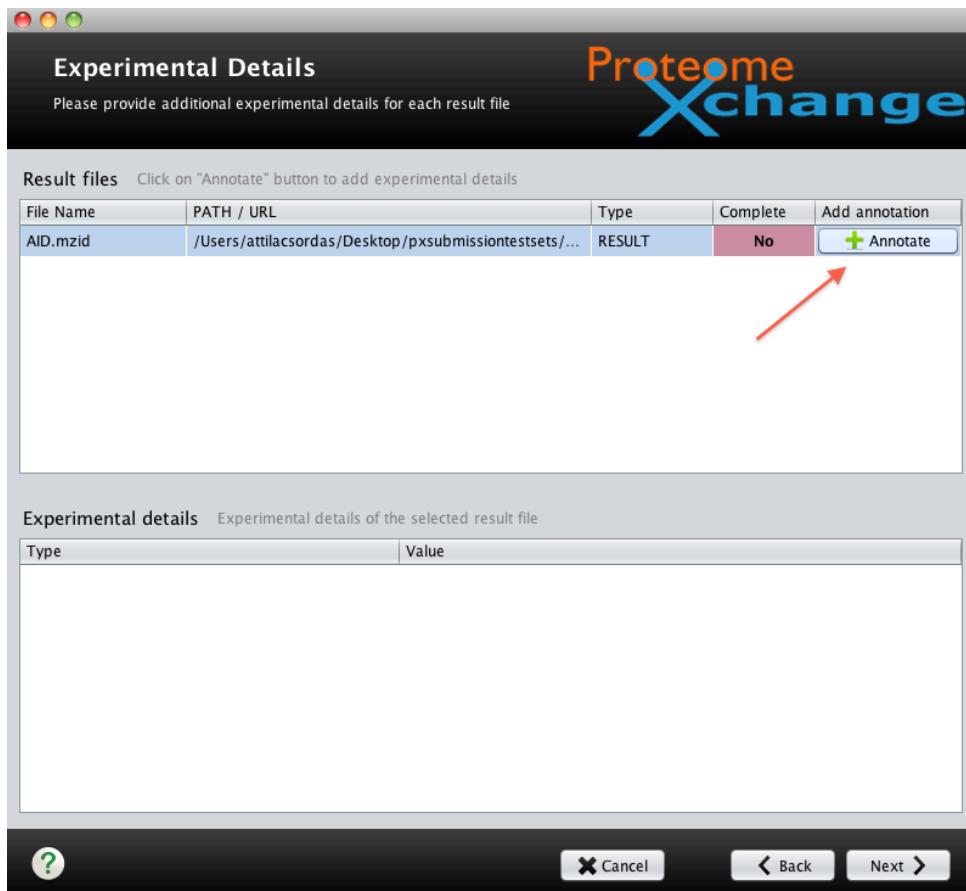
**Figure 10:** Assigning mappings between different and multiple file types on the 'Relationships between files' in the case of a PRIDE XML based 'Complete' submission.

### **Step 7: Provide additional experimental details for each result file**

Additional metadata need be provided for each 'RESULT' file in the case of a 'Complete' submission, and what is needed is the same for both subtypes of submissions (PRIDE XML and mzIdentML). Figure 11 shows the screen where the 'Annotate' button can be clicked for each 'RESULT' file. This information is usually imported automatically in the case of a PRIDE XML file (if the recommended CVs/ontologies are used). For mzIdentML, the information needs to be manually annotated.

The following additional metadata are required: species, tissue, and instrument information (provided as Controlled Vocabulary (CV) terms from a drop-down menu), and experimental factor information in a free text format (Figure 12). Optionally, providing information about the cell type, disease and quantification method (if applicable) is recommended.

If you have more than one 'RESULT' file, as it is usually the case, then you can pick the 'Apply to all' box for species and tissue information instead of doing this many times.



**Figure 11:** Please click the 'Annotate' button to add metadata to each result file.

This screenshot shows the annotation dialog box. It includes fields for 'Species\*', 'Instrument\*', 'Disease', 'Tissue\*', 'Cell type', and 'Quantification method', each with a dropdown menu and an 'Apply to all' checkbox. Below these are sections for 'Experimental factor' (with a note '(?)') and a text input field containing 'technical replicate 1'. At the bottom are 'Cancel' and 'Add' buttons.

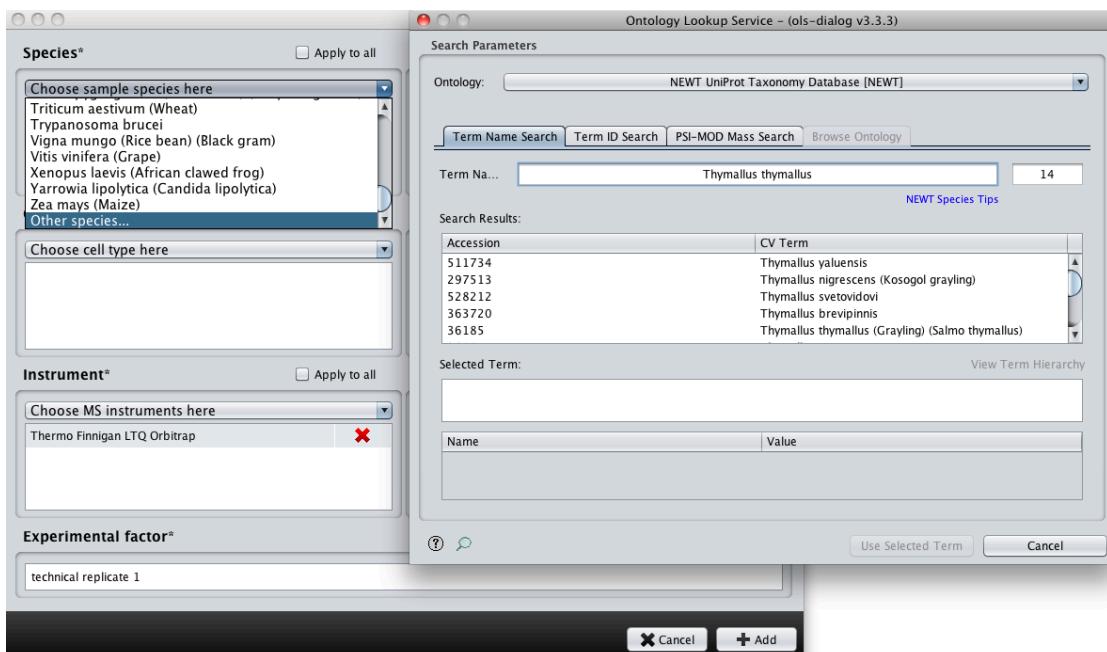
Species*		Instrument*		Disease		Tissue*		Cell type		Quantification method	
<input type="checkbox"/> Apply to all											
Choose sample species here	Homo sapiens (Human)	Choose MS instruments here	Thermo Scientific Q Exactive	Choose disease here	Acute leukemia	Choose tissue here	Blood	Choose cell type here	B cell	Choose quantification method here	Spectrum counting

Technical replicate 1

Cancel Add

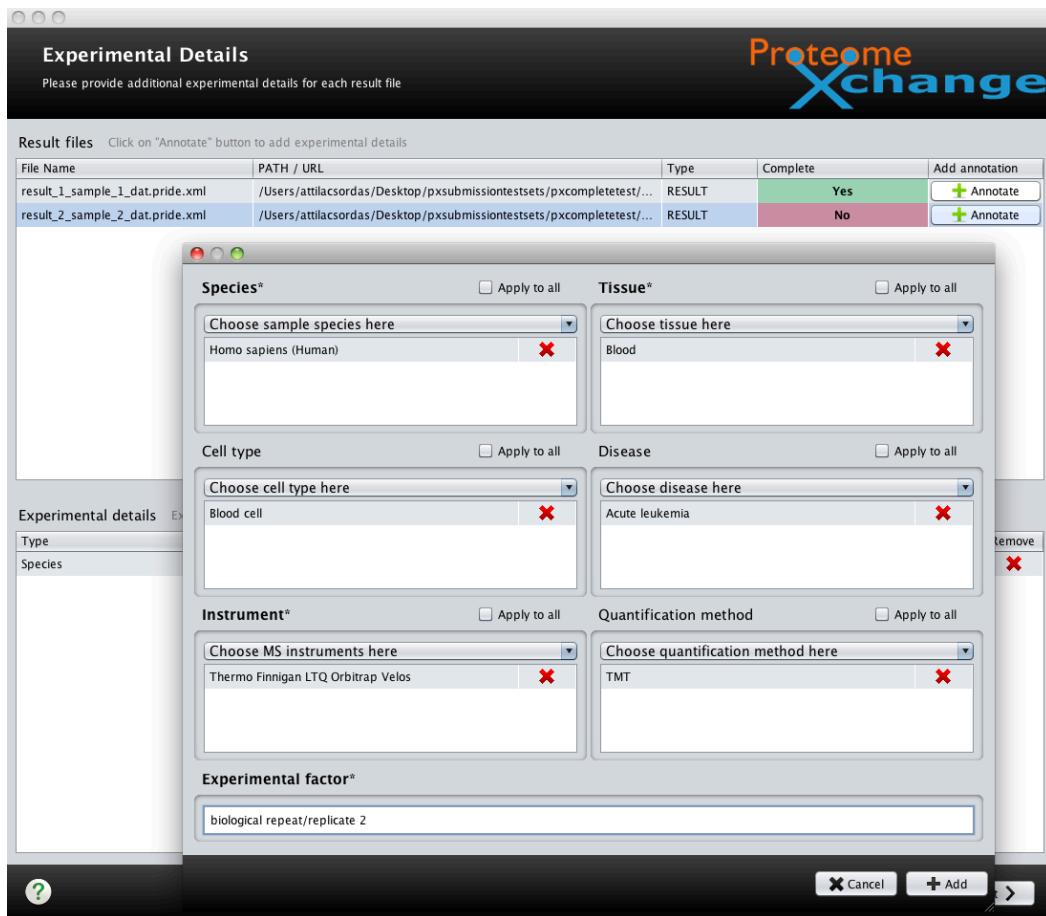
**Figure 12:** Annotating each result files with additional metadata.

In the majority of the cases you will find the metadata annotation you are looking for in the drop-down menu since the elements of the drop-down menus have been selected based on frequency. But sometimes the annotation you are looking for is not going to be available from the drop-down lists. If that's the case, you have to select to the OLS (Ontology Lookup Service) panel and search for the annotation you want to provide. For the more extensive search you need to click on the "other" options on the bottom of the drop-down menu. For instance, if you have samples from e.g. the fish Grayling (*Thymallus thymallus*) the species is not available from the drop-down list menu. You have to click on "Other species" and search for *Thymallus thymallus* in the OLS panel, see Figure 13.



**Figure 13:** Annotating a result file with additional metadata with the help of the OLS panel

In case you have multiple 'RESULT' files you have to provide data for all of them using the same panel, see Figure 14.



**Figure 14:** Annotating multiple result files.

### **Step 8: Add Lab Head**

Please provide contact details for the Lab Head/Principal Investigator of your study. Please do it in the recommended format, see Figure 15.

The screenshot shows a software window titled 'Lab Head' with a sub-instruction 'Please provide contact details of your lab head'. The main content area contains three input fields: 'Name (required)' with placeholder 'Lab head's first name and last name, i.e. John Smith'; 'Email (required)' with placeholder 'Lab head's email address'; and 'Affiliation (required)' with placeholder 'Lab head's affiliation, such as: department, lab, institute and country'. Below these fields is a note: 'NOTE: We are collecting this information for grouping submissions by lab and as a contact backup.' At the bottom of the window are standard Mac OS X interface elements: a question mark icon, a 'Cancel' button with a red 'X', a 'Back' button with a left arrow, and a 'Next >' button.

**Figure 15:** Providing contact details for the Lab Head.

### **Step 9: Optional metadata annotation**

In this panel it is recommended to provide additional metadata in four cases:

- your dataset is part of a bigger project/effort (for instance the Human Proteome Project or the EU project 'PRIME-XS'). It is a way to tag your dataset to enable grouping of datasets this way.
- there is already a PubMed ID associated with it (the data has been already published).
- your dataset represents a reanalysis of an earlier public PX dataset.
- there are other "omics" datasets (for instance transcriptomics, metabolomics data present in other repositories) that can be associated with it. In this case, please provide the accession number of the dataset in the corresponding repository.

The screenshot shows a software window titled 'Additional dataset details' for the 'ProteomeXchange' platform. The window has a dark header bar with the title and a logo. Below the header, there is a section for 'Parent project (optional)' containing a list of checkboxes for various projects. There are also sections for 'PubMed ID(s) (optional)', 'Reanalysis ProtomeXchange accession(s) (optional)', and 'Links to other 'Omics' datasets (optional)'. At the bottom, there are buttons for '?', 'Cancel', 'Back', and 'Next >'. A note at the top of the form asks users to provide additional details about their dataset.

**Parent project (optional)**  
If your project is part of a larger project, please select the parent project from the table below. If you would like to propose a new parent project, please contact us at: [pride-support@ebi.ac.uk](mailto:pride-support@ebi.ac.uk)

**Parent Project**

- Human Proteome Project
- Biology/Disease Based Human Proteome Project
- Chromosome Based Human Proteome Project
- PRIME-XS Project
- CPTAC Consortium

**PubMed ID(s) (optional)**  
Provide the PubMedID(s) if the dataset is associated with an existing publication (comma separated)

**Reanalysis ProtomeXchange accession(s) (optional)**  
Only applicable if your results are based on the reprocessing of one or several previously submitted PX dataset(s)

**Links to other 'Omics' datasets (optional)**  
Only applicable if proteomics results can be linked to other biological data submitted to other resources (e.g. ArrayExpress, GEO)

?

Cancel

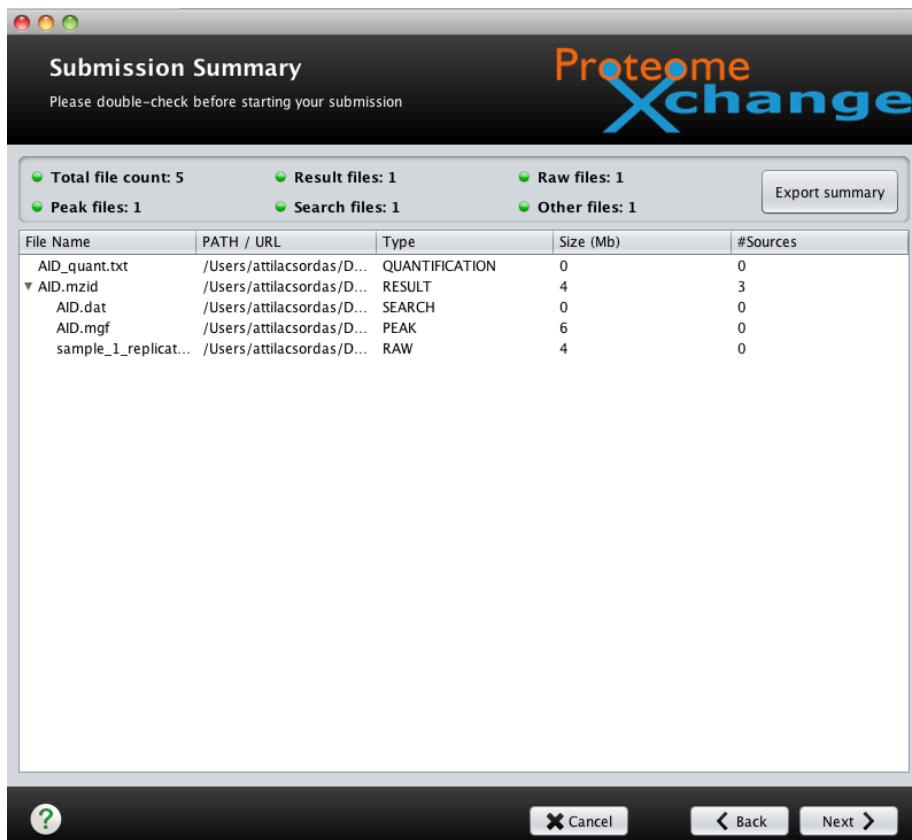
Back

Next >

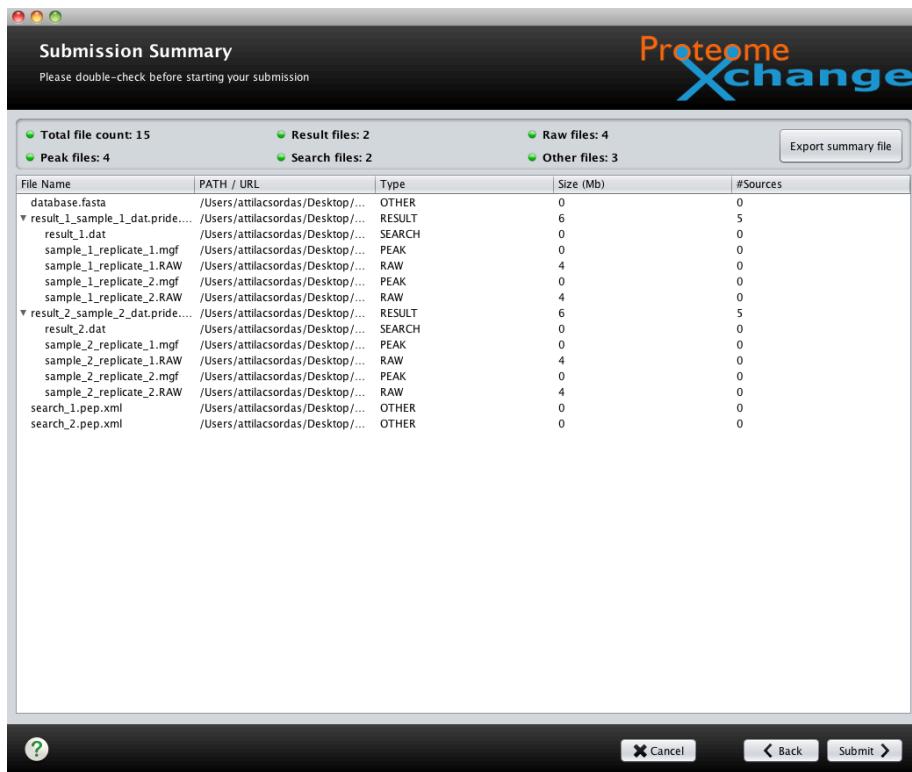
**Figure 16:** Providing additional, applicable metadata.

### **Step 10: Check before submission**

This is the last step before the file upload actually starts. You should double-check that all the necessary files are included in the submission summary before continuing to the upload step, see an example of an mzIdentML based 'complete' submission in Figure 17. Figure 18 shows the Submission Summary page with multiple result files in case of a PRIDE XML based 'Complete' submission.



**Figure 17:** 'Submission Summary' screen in the PX Submission Tool with a single 'RESULT' file.



**Figure 18:** 'Submission Summary' screen in the PX Submission Tool with multiple result files.

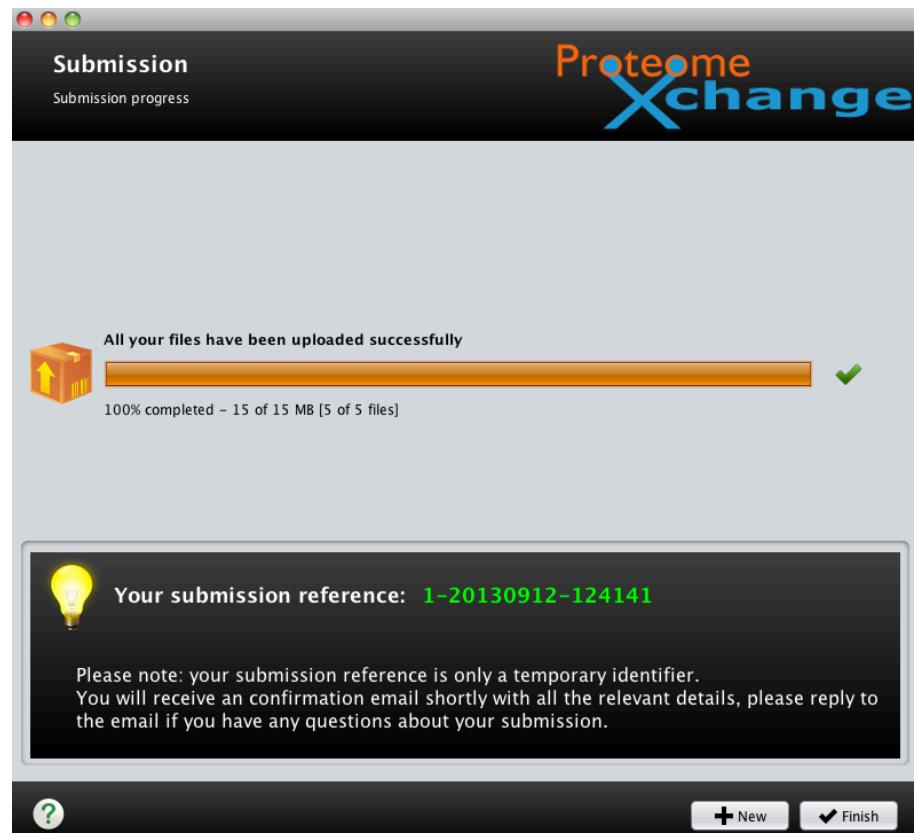
### **Step 11: File Submission**

This is the actual step when all your files are uploaded to PRIDE and ProteomeXchange (Figure 19). Once the upload is finished, an e-mail will be sent

to you to confirm that all your files have been uploaded successfully and that are waiting to be validated.

If for any reason the tools crashes at this point, the PX Submission Tool can be restarted and the file upload will restart in the same point before it crashed.

You will be also issued with a temporary submission reference, to help us to quickly identify and track your submission should you have any questions. This is not the PX accession number.



**Figure 19:** 'Submission' screen of the PX Submission Tool showing that a submission has been completed.

## 5 How to make Partial Submissions?

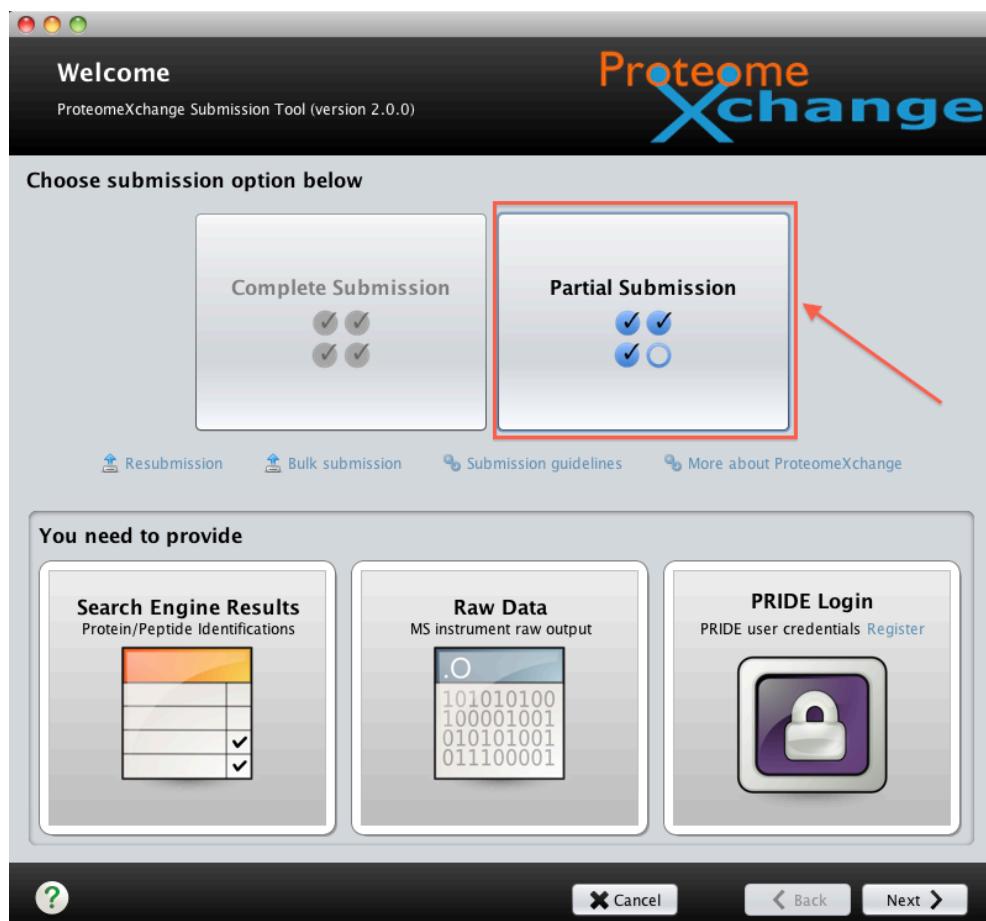
Remember that by default we recommended doing ‘Complete’ submissions. You should only use this option if your ‘RESULT’ files cannot be converted/exported to PRIDE XML or mzIdentML 1.1.

### **Step 1: Launch PX Submission Tool**

Please install and launch the PX Submission Tool (available at <http://www.proteomexchange.org/submission>).

### **Step 2: Select Submission Type**

Select ‘Partial Submission’ in the PX Submission Tool ‘Welcome’ screen (Figure 20).



**Figure 20:** Selecting Partial Submission in the ‘Welcome’ screen of the PX Submission tool.

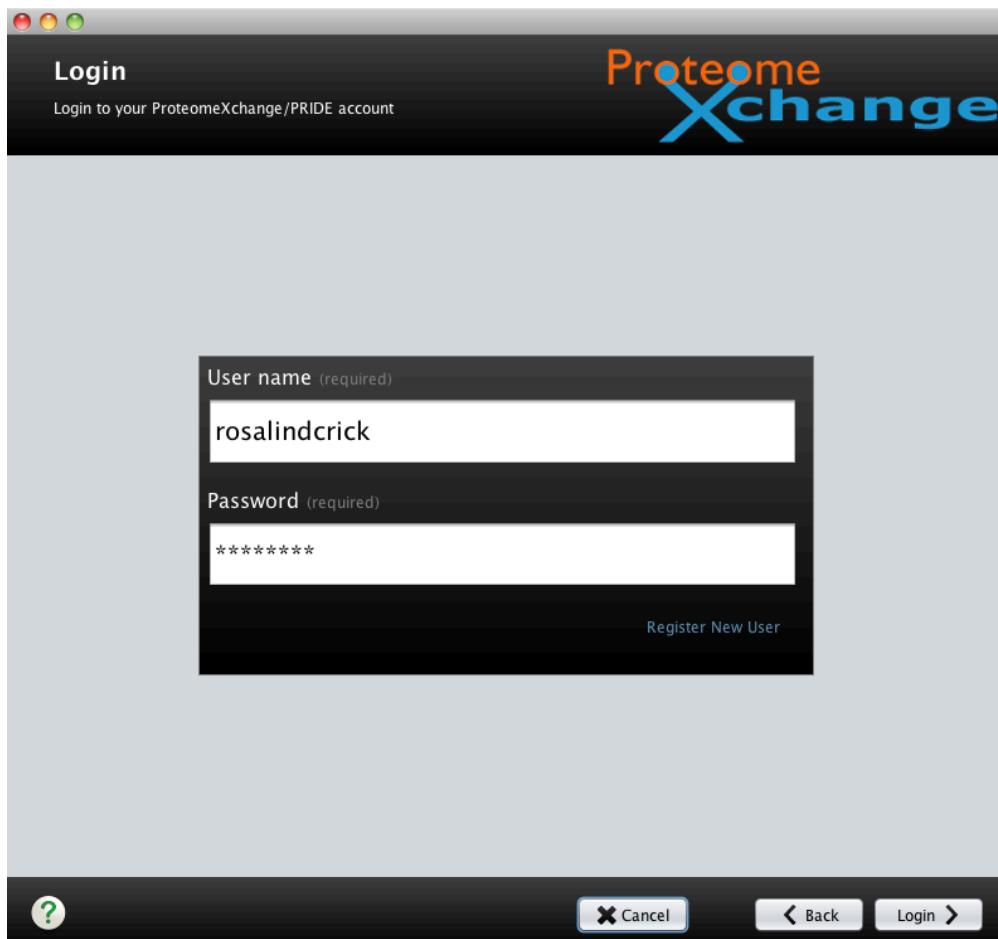
Upon selecting this option a warning will pop up, see Figure 21. Continue with clicking ‘Yes’.



**Figure 21:** Warning concerning Partial Submissions in the PX Submission tool.

### **Step 3: Login**

Please log in using your existing PRIDE account as shown in Figure 22.



**Figure 22.** Login screen of the PX Submission tool.

#### **Step 4: Provide submission details**

The user is asked to provide some basic details about the uploaded dataset (Figure 23) such as the title, a list of keywords (in a comma separated format), and a brief description of the data (similar to the abstract of the corresponding publication) a sample processing and a data processing protocol. The user also picks a mass spectrometry experiment type from a drop-down menu.

The screenshot shows the 'Dataset Details' screen of the PX Submission tool. At the top right is the 'Proteome Xchange' logo. A tip at the top right says 'Tip: Use Ctrl+C to copy, Ctrl+V to paste'. The main form fields include:

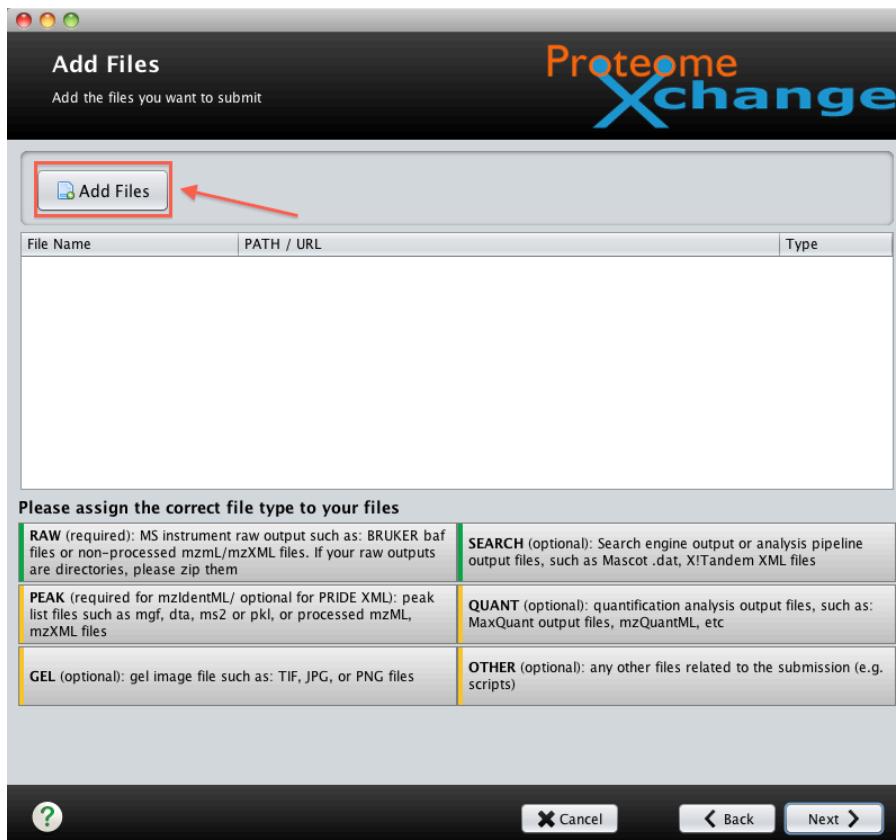
- Project title\***: i.e. Human liver LC-MSMS
- Keywords**: i.e. Human, Liver, Plasma, LC-MSMS
- Project description\*** (50 to 5000 characters): Please provide an overall description of your study, think something similar in scope to the manuscript abstract
- Sample processing protocol\*** (50 to 5000 characters): Please provide a short description on the sample preparation steps, separation, enrichment strategies and mass spectrometry protocols included
- Data processing protocol\*** (50 to 5000 characters): Please provide a couple of sentences on the bioinformatics pipeline used, main search parameters, quantitative analysis, software tools and versions included. Think something similar in scope to the Data Analysis section of your manuscript
- Experiment type\***: A dropdown menu titled 'Choose experiment type here' with options:
  - Choose experiment type here
  - Shotgun proteomics
  - Cross-linking (CX-MS)
  - Affinity purification (AP-MS)
  - SRM/MRM
  - SWATH MS
  - Other experiment type...

At the bottom are buttons for **Cancel**, **Back**, and **Next >**.

**Figure 23:** 'Dataset details' screen in the PX Submission tool.

### Step 5: Add Files and assign file types

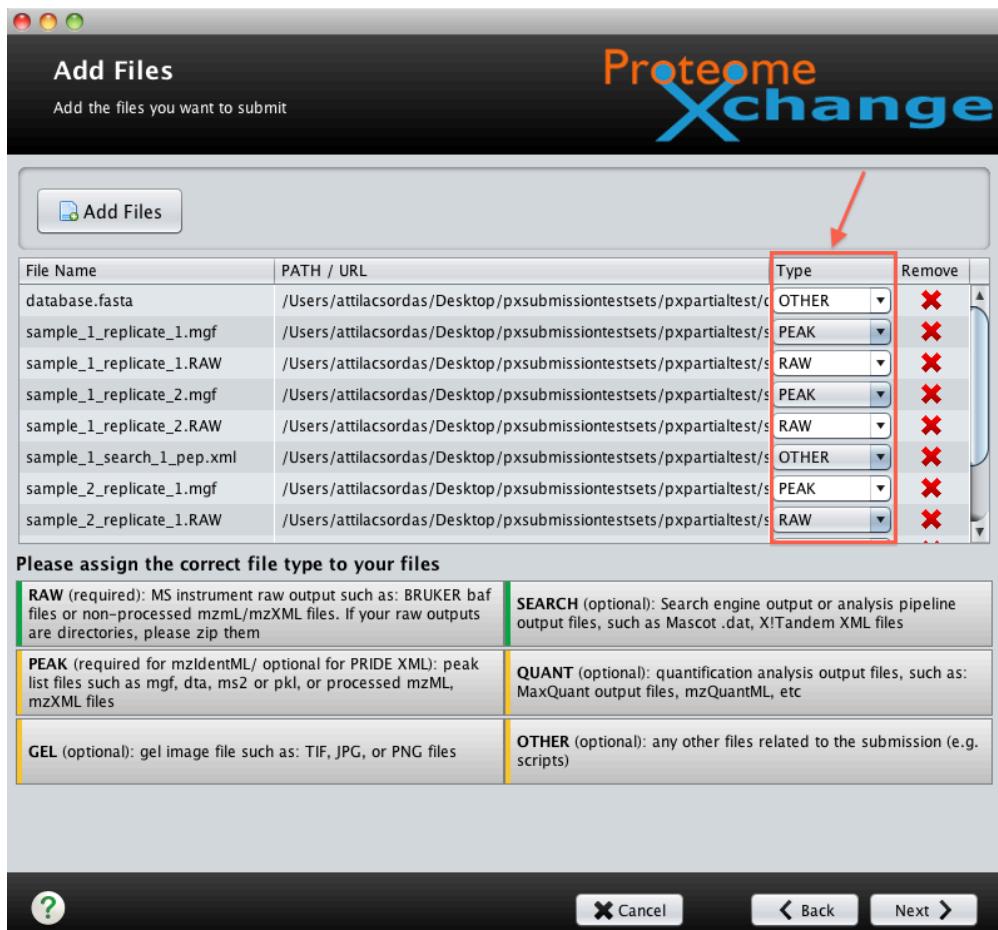
You should choose the files you would like submit in this step. As shown in Figure 24 you can add files by clicking on the highlighted button.



**Figure 24:** 'Add files' screen of the PX submission tool.

You should make sure that both the 'SEARCH' search engine output files and the 'RAW' files are selected. There could also be other files types included in the submission: 'PEAK' (for peak list files), 'QUANT', for quantification results, 'GEL', for gel images, or 'OTHER' (any other file). All the files need to be selected at this stage.

Once the files are added, double-check them to make sure they were assigned with the correct file types. For instance in Figure 25, the pep.xml 'SEARCH' file has been recognized as 'OTHER' file and this need to be changed by selecting 'SEARCH' from the drop-down menu.



**Figure 25:** PX Submission Tool ‘Add Files’ screen: Assignment of the correct file types.

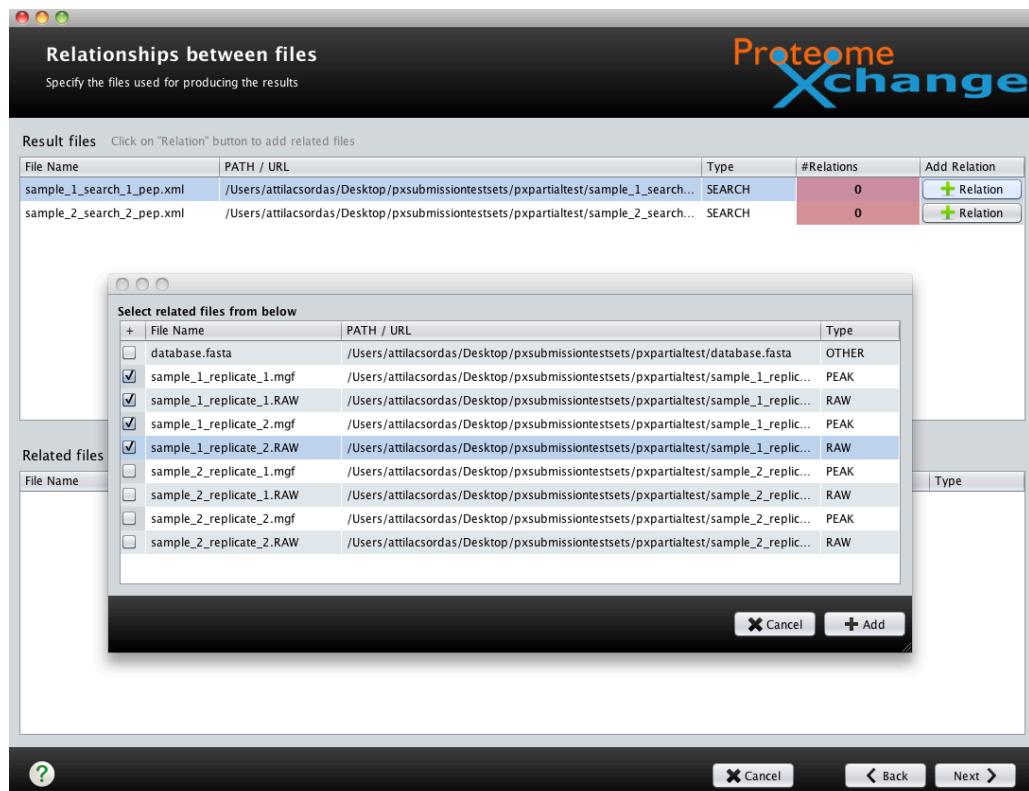
## Step 6: Assign relationships between the submitted files

This mapping step consists of assigning the relations between the ‘SEARCH’ files and the other file types included in the submission, for example, which ‘RAW’ (mandatory) or ‘PEAK’ files have been used to produce the search engine output files (‘SEARCH’). ‘QUANT’ or ‘OTHER’ files can also be added. This will enable others to understand how your files are connected.

By default the tool makes an attempt to generate the mapping between the ‘SEARCH’ and the other - most importantly ‘RAW’ - files. For instance if there has been only 1 ‘SEARCH’ file found during the previous ‘Add Files’ step (Step 5) then all the other files will be mapped to this ‘SEARCH’ file. If there are multiple ‘SEARCH’ files the tool maps the other files – ‘RAW’, ‘PEAK’, ... - with the same file name prefix, but without the file extension, to the corresponding ‘SEARCH’ files. This mapping is done even if the suffix part of the ‘RAW’ files contains different numbers (for instance indicating different replicates) or the prefix contains spaces or underscores.

If the automatic mapping is partial only or does not apply, the submitter is asked to manually assign the relationships between the files.

Each 'SEARCH' file must have at least one file linked to it. As shown in Figure 26, this is done by clicking on the 'Add Relation' button. Many files can be assigned to the same 'SEARCH' file.



**Figure 26:** Assigning mappings between different file types on the 'Relationships between files' screen in the PX Submission tool.

### Step 7: Provide additional experimental details

In order to increase the reusability of the dataset, some additional experimental details are needed such as species, tissue, cell type, disease, MS instrument and a list of the post-translational modifications (PTMs) present in the dataset.

**Additional Details**  
Please give additional details about your submission

**Species\***: Choose sample species here  
Homo sapiens (Human) X

**Tissue\***: Choose tissue here

**Modification\***: Choose modifications here  
Phosphorylation X  
Oxidation X

**Instrument\***: Choose MS instruments here  
Thermo Scientific Q Exactive X

**Cell type**: Choose cell type here  
Blood cell X

**Disease**: Choose disease here  
Acute leukemia X

**Quantification method**: Choose quantification method here  
Spectrum counting X

?

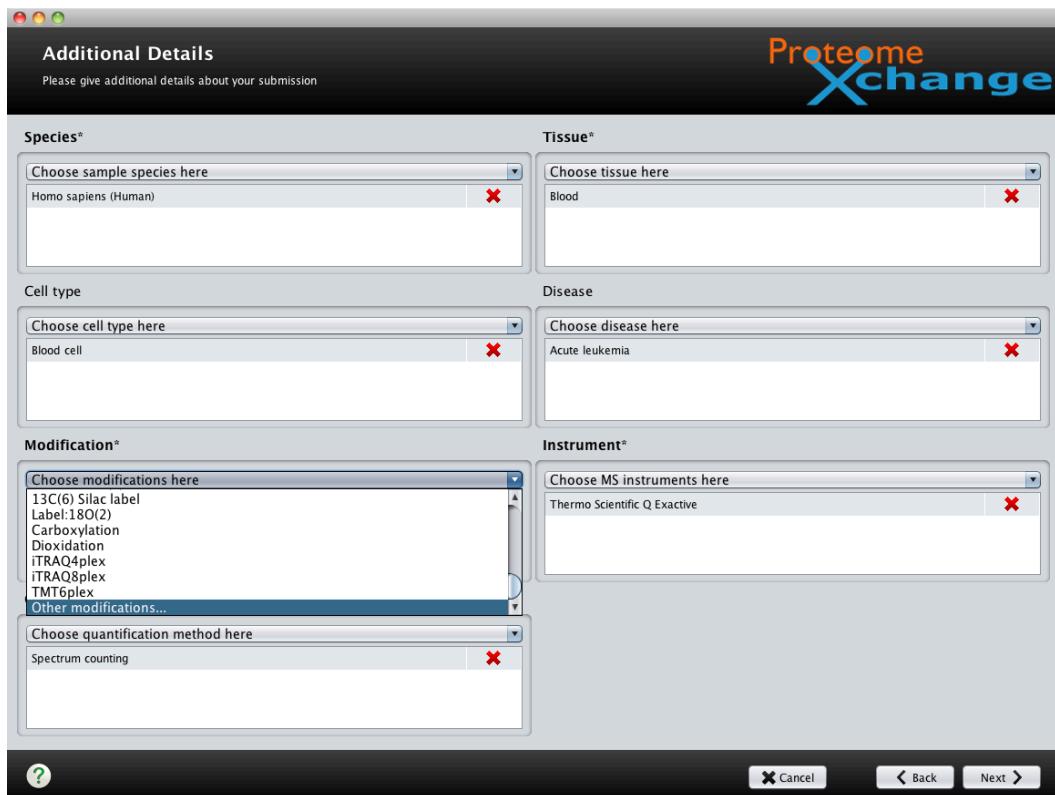
Cancel

Back

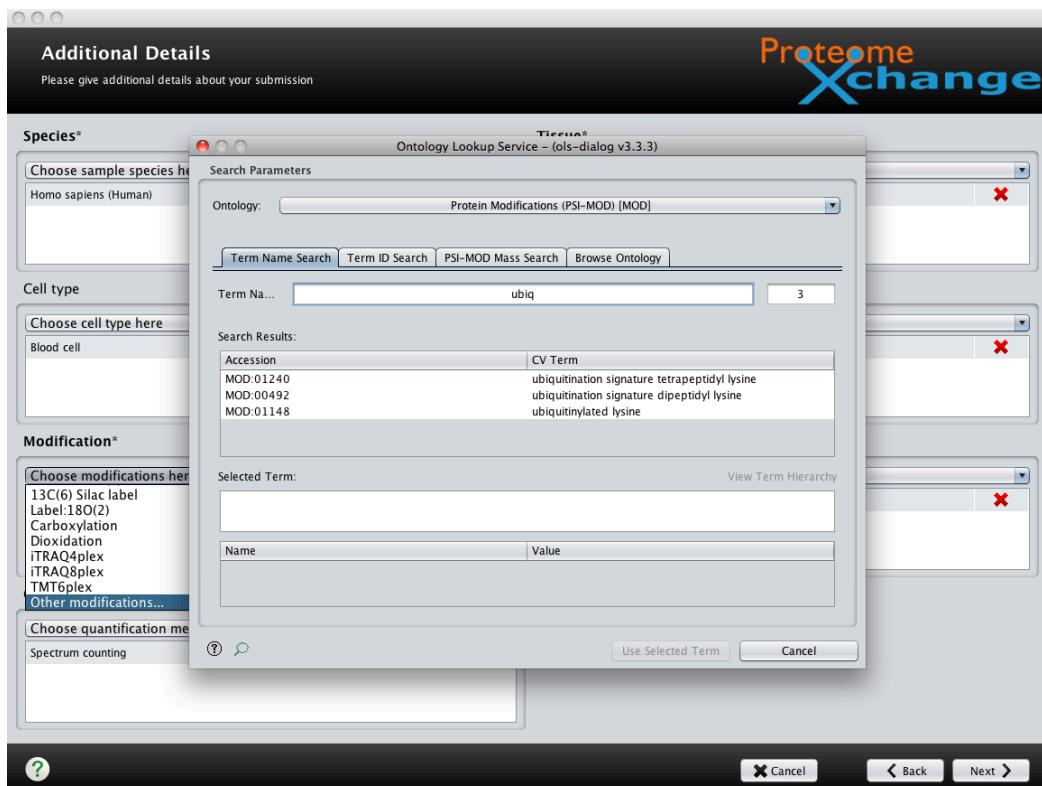
Next >

**Figure 27:** 'Additional details' screen in the PX Submission Tool for Partial Submissions.

For each type of required experimental details, the submission tool provides a short list of commonly used values (Figure 27). If this list doesn't contain your experimental specific details, you should choose the 'Other' option, as shown in Figure 28 for modifications. If that option is selected, a pop-up window will appear providing access to the 'Ontology Lookup Service' (OLS, <http://www.ebi.ac.uk/ontology-lookup/>). There, you can search for a specific term from a controlled vocabulary or ontology, please see Figure 29.



**Figure 28:** Screenshot of the PX Submission Tool showing how to choose 'other' modifications.



**Figure 29:** Screenshot with the 'Ontology Lookup Service' (OLS) pop-up window in the PX submission tool.

## Step 8: Add Lab Head

Please provide contact details for the Lab Head/Principal Investigator of your study (Figure 30).

The screenshot shows a software window titled 'Lab Head' with the 'Proteome Xchange' logo at the top right. A sub-instruction 'Please provide contact details of your lab head' is displayed above the form fields. The form consists of three text input fields: 'Name (required)' containing 'Lab head's first name and last name, i.e. John Smith'; 'Email (required)' containing 'Lab head's email address'; and 'Affiliation (required)' containing 'Lab head's affiliation, such as: department, lab, institute and country'. Below these fields is a note: 'NOTE: We are collecting this information for grouping submissions by lab and as a contact backup.' At the bottom of the window are standard navigation buttons: a question mark icon, 'Cancel' (with an 'X'), 'Back' (with a left arrow), and 'Next >' (with a right arrow).

**Figure 30:** Providing contact details for the Lab Head of your project.

### **Step 9: Optional metadata annotation**

In this panel it is recommended to provide additional metadata in four cases:

- your dataset is part of a bigger project/effort (for instance the Human Proteome Project or the EU project 'PRIME-XS'). It is a way to tag your dataset to enable grouping this way.
- there is already a PubMed ID associated with it (the data has been already published).
- your dataset represents a reanalysis of an earlier public PX dataset.
- there are other "omics" datasets (for instance transcriptomics, metabolomics data present in other repositories) that can be associated with it. In this case, you need to provide the accession number of the dataset in the corresponding repository.

The screenshot shows a Mac OS X window titled 'Additional dataset details' for the 'ProteomeXchange' platform. The window contains several optional input fields:

- Parent project (optional):** A list of checkboxes for various projects:
  - Human Proteome Project
  - Biology/Disease Based Human Proteome Project
  - Chromosome Based Human Proteome Project
  - PRIME-XS Project
  - CPTAC Consortium
- PubMed ID(s) (optional):** A text input field for comma-separated PubMed IDs.
- Reanalysis ProteomeXchange accession(s) (optional):** A text input field for previously submitted PX dataset accessions.
- Links to other 'Omics' datasets (optional):** A text input field for linking to other biological data resources like ArrayExpress or GEO.

At the bottom of the window are standard Mac OS X interface elements: a question mark icon, a 'Cancel' button, a 'Back' button, and a 'Next >' button.

**Figure 31:** Providing additional, applicable metadata.

### **Step 10: Check before submission**

This is the last step before the file upload actually starts. You should double-check that all the necessary files are included in the submission summary before continuing to the upload step, please see Figure 32.

The screenshot shows the 'Submission Summary' screen of the PX Submission tool. At the top, it displays file counts: Total file count: 11, Result files: 0, Peak files: 4, Search files: 2, Raw files: 4, and Other files: 1. There is a button to 'Export summary file'. Below this is a table listing the uploaded files:

File Name	PATH / URL	Type	Size (Mb)	#Sources
database.fasta	/Users/attilacsordas/Desktop...	OTHER	0	0
▼ sample_1_search_1_pep.xml	/Users/attilacsordas/Desktop...	SEARCH	0	4
sample_1_replicate_1.mgf	/Users/attilacsordas/Desktop...	PEAK	0	0
sample_1_replicate_1.RAW	/Users/attilacsordas/Desktop...	RAW	4	0
sample_1_replicate_2.mgf	/Users/attilacsordas/Desktop...	PEAK	0	0
sample_1_replicate_2.RAW	/Users/attilacsordas/Desktop...	RAW	4	0
▼ sample_2_search_2_pep.xml	/Users/attilacsordas/Desktop...	SEARCH	0	4
sample_2_replicate_1.mgf	/Users/attilacsordas/Desktop...	PEAK	0	0
sample_2_replicate_1.RAW	/Users/attilacsordas/Desktop...	RAW	4	0
sample_2_replicate_2.mgf	/Users/attilacsordas/Desktop...	PEAK	0	0
sample_2_replicate_2.RAW	/Users/attilacsordas/Desktop...	RAW	4	0

At the bottom, there are buttons for '?', 'Cancel', 'Back', and 'Submit'.

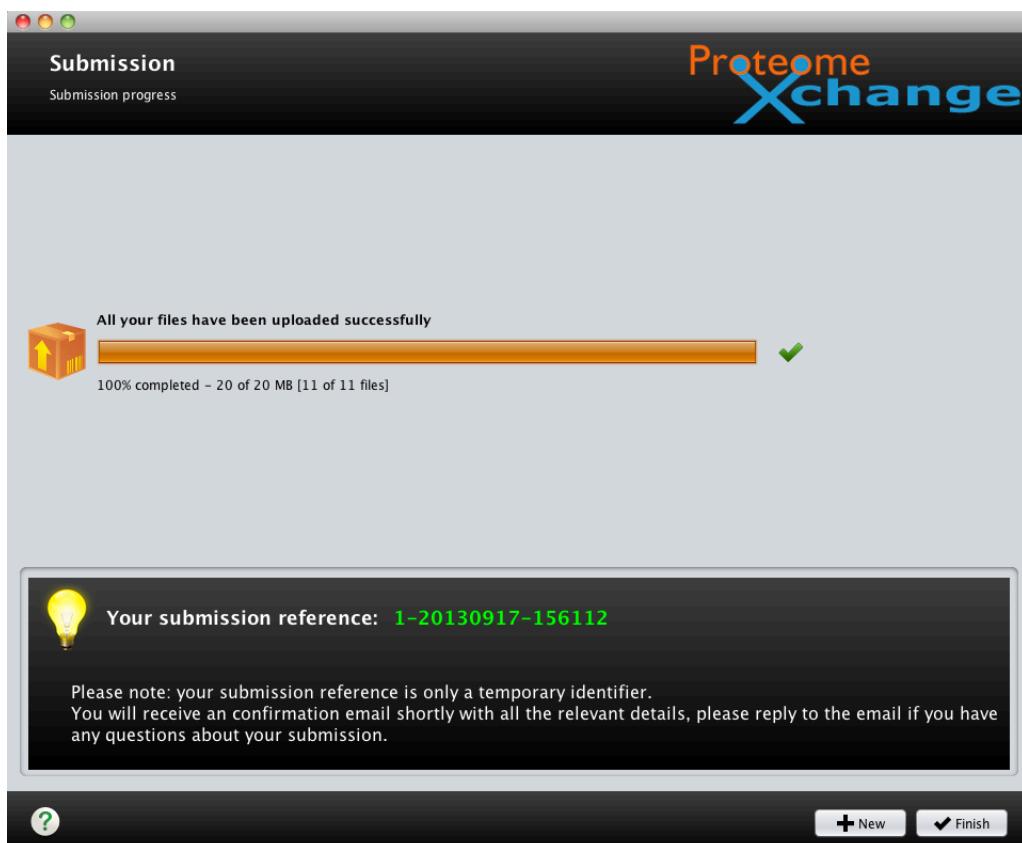
**Figure 32:** 'Submission Summary' screen for a 'Partial' submission in the PX Submission tool.

### **Step 11: File Submission**

This is the actual step when all your files are uploaded to PRIDE and ProteomeXchange. Once the upload is finished, an email will be sent to you to confirm that all your files have been uploaded successfully and that are waiting to be validated.

If for any reason the tool crashes at this point, the PX Submission Tool can be restarted and the file upload will restart in the same point before it crashed.

You will be also issued with a temporary submission reference, to help us to quickly identify and track your submission should you have any questions. This is neither the final PX accession number, nor a temporary one. As such it should not be used in the manuscript.



**Figure 33:** ‘Submission’ screen of the PX Submission Tool showing that a submission has been completed.

For particular examples of partial submissions (e.g. software like MaxQuant or ProteinPilot), see Appendix V.

## 6 How to make bulk submissions?

Two steps are required: ‘Creation of the PX submission summary file’, and ‘Submission using the PX submission tool’.

### 6.1 Creation of the PX Submission Summary File

A submission summary file (submission.px) contains two types of information needed for any PX submission:

- **Metadata:** general experimental metadata like experiment description, sample taxonomy information, instruments and modifications used, experimental tags, contact information, etc.
- **Mapping between the uploaded files:** for instance between the ‘RAW’ files and the corresponding ‘RESULT’ or search engine output files (‘SEARCH’).

There are two ways to create the file:

A) Generating the file independently from the PX submission tool. Some scripting work is needed. Details about the tab delimited PX submission format can be found [here](#).

B) Using the PX Submission Tool: This is the recommended option if there are not many files, so the metadata and the file mappings can be provided with the tool but the actual FTP upload is not performed at that point. Instead the submitters can upload their files in an alternative way (see Section 6.3). For these cases the PX Submission Tool provides an 'Export Summary' functionality. You can use this functionality when reaching the 'Submission Summary' screen, at the end of the submission process, please see Figure 34. The summary file can then be stored locally (usually with the extension .px).

File Name	PATH / URL	Type	Size (Mb)	#Sources
database.fasta	/Users/attilacsordas/Desktop...	OTHER	0	0
sample_1_search_1_pep.xml	/Users/attilacsordas/Desktop...	SEARCH	0	4
sample_1_replicate_1.mgf	/Users/attilacsordas/Desktop...	PEAK	0	0
sample_1_replicate_1.RAW	/Users/attilacsordas/Desktop...	RAW	4	0
sample_1_replicate_2.mgf	/Users/attilacsordas/Desktop...	PEAK	0	0
sample_1_replicate_2.RAW	/Users/attilacsordas/Desktop...	RAW	4	0
sample_2_search_2_pep.xml	/Users/attilacsordas/Desktop...	SEARCH	0	4
sample_2_replicate_1.mgf	/Users/attilacsordas/Desktop...	PEAK	0	0
sample_2_replicate_1.RAW	/Users/attilacsordas/Desktop...	RAW	4	0
sample_2_replicate_2.mgf	/Users/attilacsordas/Desktop...	PEAK	0	0
sample_2_replicate_2.RAW	/Users/attilacsordas/Desktop...	RAW	4	0

**Figure 34:** 'Submission Summary' screen in the PX Submission tool, highlighting how to export and store locally the PX summary file.

## 6.2 Submission using the PX Submission tool

You have already created a PX submission summary file for your dataset by scripting. But the actual size of the dataset is only around a few GBs despite the many files included. In this case you can use the PX Submission Tool to perform the submission. In the 'Welcome' screen of the PX submission tool, please select the option 'Bulk submission' highlighted in Figure 35, and proceed as indicated by the tool. You will need to load the created PX summary file.



**Figure 35:** 'Welcome screen' of the PX Submission Tool highlighting the 'Bulk submission' mode.

### 6.3 Fast upload option for big datasets using Aspera.

If the size of your dataset is over a few GB and/or the FTP upload process is considerably slow we can provide the Aspera protocol (<http://www.asperasoft.com/>), at present *via* a command line upload option. Some command line skills are needed in order to use the Aspera upload option. In the future we plan to include this functionality in the PX submission tool. If you wish to upload your datasets via Aspera, please follow the steps below.

**Requirements:** Please download the Aspera Connect Web Browser Plug-in. Although you download a Browser Plug-in you will be using the 'ascp' command line transfer program distributed with it.

Operating System: Windows XP / 2003 / Vista / 2008 / 7 / 8, Mac OS Intel 10.5 / 10.6 / 10.7 / 10.8

You don't have to register in order to download the Browser Plug-in and the download is free of charge.

- Check the command line transfer usage for more configuration details. This is the location of the 'ascp' program in the file system:

- Mac: on the desktop go

```
cd /Applications/Aspera\ Connect.app/Contents/Resources/  
there you'll see the command line utilities where you're going to use 'ascp'.
```

- Windows: the downloaded files are a bit hidden. For instance in Windows 7 the ascp.exe is located in the users home directory at: AppData\Local\Programs\Aspera\Aspera Connect\bin\ascp.exe

### **How to upload a directory of files**

**Step 1.** Ask PRIDE support (at pride-support@ebi.ac.uk) for a target directory and a password.

The PRIDE curators will specify a target directory for you, see <name-of-target-dir-specified-by-PRIDE> in the following commands, and you will be provided with this information.

**Step 2.** The upload command and process.

When preparing your dataset please be sure to unambiguously assign a unique file name to all of your files. Please also upload the submission summary file into the same folder.

- Mac: ./ascp -QT -l500m --file-manifest=text -k 2 <path-to-folder-to-be-uploaded> pride-drop-006@ah01.ebi.ac.uk:<name-of-target-dir-specified-by-PRIDE>

- Windows: ascp.exe -QT -l500m --file-manifest=text -k 2 <path-to-folder-to-be-uploaded> pride-drop-006@ah01.ebi.ac.uk:<name-of-target-dir-specified-by-PRIDE>

The <path-to-folder-to-be-uploaded> should not have any blank spaces in it.

Please set the '--file-manifest=text -k 2' flags as well.

This will generate an Aspera progress file on your side that will contain a report on the files that have been uploaded, also you can interrupt the process and then it will only upload the ones that were not there so no more overwriting files. It will also skip the ones that are already in the target directory.

If -l500m ~ 500 Mb/s is unstable and leads to timeouts then we suggest to go back to -l250m as the maximum transfer rate, even that is fast enough to transfer theoretically 2 TBs within a day.

Once upload has been finished you will be prompted to enter the password provided earlier.

### **Step 3. Notify the PRIDE Team**

E-mail [pride-support@ebi.ac.uk](mailto:pride-support@ebi.ac.uk) in case your upload has been successfully finished.

## **7 What happens after the submitter has uploaded all the data?**

Once your dataset has been uploaded into the EBI, the PRIDE/ProteomeXchange internal submission pipeline will validate your files. The results of the validation will be checked by a curator and, if no problems are found, the dataset will be submitted to PRIDE and the relevant information will be stored. The process varies for ‘complete’ and ‘partial’ submissions. As a result, you will be issued with a ProteomeXchange accession number.

In addition, a DOI will also be assigned if a ‘complete’ submission was performed. PRIDE assay accession numbers will also be provided for PRIDE XML and mzIdentML result files in case of ‘complete’ submissions. A confirmation e-mail will be sent to you with all the relevant details once your submission is complete, including a username and password for potential journal reviewers and editors to be able to access your data privately. Please note all submissions are private by default.

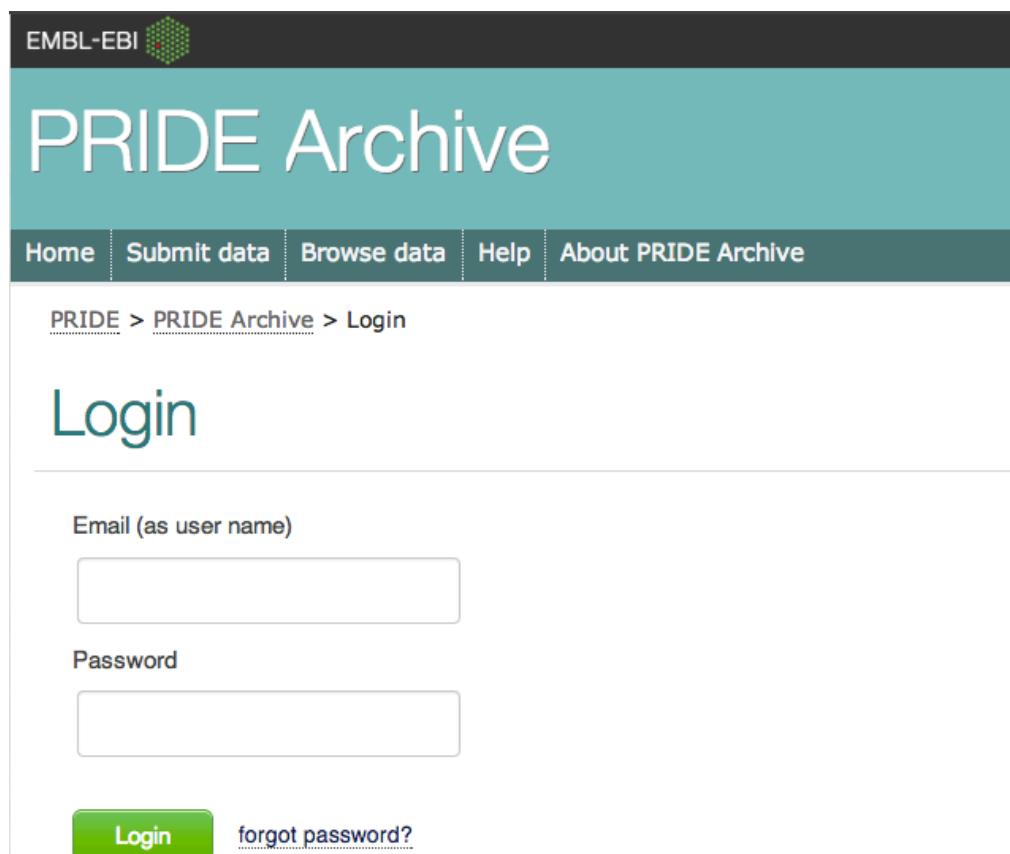
## **8 Accessing Private Data**

Submitted datasets are private by default, which means you need to be logged-in to view your data. We will however also create a PX reviewer account and a password for your dataset, which you should include in your manuscript. Again, the PX reviewer account will give you access to all of the files belonging to your submission. For that you can use the new PRIDE Archive web site or the PRIDE Inspector tool.

### **8.1 PRIDE Archive web page**

The new PRIDE Archive web site is available at <http://www.ebi.ac.uk/pride/archive>. Registered submitters can use their personal accounts or the reviewer accounts to access and download the individual PX datasets. For every submission there is a separate reviewer account generated.

Please navigate first to the login page available at <http://www.ebi.ac.uk/pride/archive/login> (see Figure 36):



The screenshot shows the PRIDE Archive login page. At the top left is the EMBL-EBI logo. The main title 'PRIDE Archive' is centered above a navigation bar with links for Home, Submit data, Browse data, Help, and About PRIDE Archive. Below the navigation bar is a breadcrumb trail: PRIDE > PRIDE Archive > Login. The main section is titled 'Login'. It contains two input fields: 'Email (as user name)' and 'Password', each with a corresponding text input box. Below the password field is a 'Forgot password?' link. At the bottom left is a green 'Login' button.

**Figure 36:** PRIDE Archive ‘Login’ page.

Once logged in with your registered User (the e-mail account you used to register in PRIDE) or an issued Reviewer Account you are going to see the private dataset/s listed.

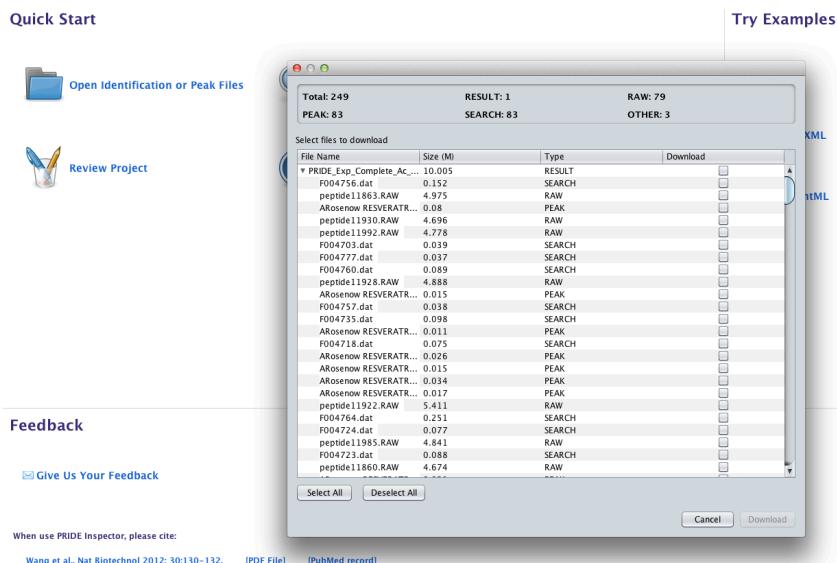
## 8.2 PRIDE Inspector

PRIDE Inspector is a stand-alone tool developed by the PRIDE team. It can be downloaded here:

<http://code.google.com/p/pride-toolsuite/wiki/PRIDEInspector>,

for further information please see Appendix 2.

In order to access private datasets, first open PRIDE Inspector by clicking on the pride-inspector-<version-number>.jar file in the tool's working directory and go to Review Project-> Reviewer account details. You can open the mzIdentML (plus spectra files) or PRIDE XML result files with PRIDE Inspector or just download all the files that you wish to investigate.



**Figure 37:** Downloading data with the reviewer account using PRIDE Inspector private download option.

Alternatively you can launch PRIDE Inspector with a Java Web Start URL provided in the automatic "Submission Complete" e-mail. This option is for downloading the mzIdentML and PRIDE XML files only into a target folder and so it is available only for 'complete' submissions. In order to use the PRIDE Inspector Java Web Start option to display your data there is a waiting period of up to one day upon getting the automatic "Submission Complete" e-mail.

## 9 Post-submission steps

### 9.1 How to do a resubmission of a dataset?

While the data is still private (during the manuscript review process) it is possible to resubmit the whole dataset by keeping the previously issued PX identifier. Data resubmissions consisting in a subset of the previous submission are not currently supported.

#### 9.1.1 Resubmission with the PX Submission Tool

Install and launch the PX Submission Tool as explained before (available at <http://www.proteomexchange.org/submission>).

##### **Step 1: Click resubmission on the 'Welcome' page**

The option is highlighted in Figure 38.



**Figure 38:** 'Welcome screen' of the PX Submission Tool highlighting the resubmission mode.

### **Step 2: Enable resubmission and provide resubmission details**

In the pop-up dialog box please provide your PRIDE login details and select the PX identifier of the dataset you want to resubmit, please see Figure 39.



**Figure 39:** Screenshot showing how to select the dataset that needs to be resubmitted.

After these two steps the resubmission follows the same steps described for a regular submission.

### 9.1.2 Resubmission via Aspera

If you have done a bulk submission using the command line Aspera fast transfer option resubmission of the whole dataset is possible via Aspera again. You will upload the whole modified dataset with the submission summary file into the same target directory again. You can use the PX Submission Tool to export the summary file as explained before but in that case you need to use the “Resubmission” option of the tool and specify the PX Identifier that will be used for resubmission, please see the 9.1.1 section above. This way the summary file will contain the required resubmission information.

In case you are generating the summary file using scripting (see section 6.1) the following line need be added to the Metadata section of the submission.px file to indicate that the dataset is a resubmission of an earlier submitted one:

```
MTD resubmission_px PXD000444
```

## 9.2 Referencing the dataset in the paper

By default we recommend to add the following formula to your manuscript (typically in "Material and Methods" or just before/in the "Acknowledgements"):

The mass spectrometry proteomics data have been deposited to the ProteomeXchange Consortium (<http://proteomecentral.proteomexchange.org>) via the PRIDE partner repository [1] with the dataset identifier <PXD000xxx>."

[1] and also for general PRIDE reference, please use: Vizcaino JA, Cote RG, Csordas A, Dianes JA, Fabregat A, Foster JM, Griss J, Alpi E, Birim M, Contell J, O'Kelly G, Schoenegger A, Ovelleiro D, Perez-Riverol Y, Reisinger F, Rios D, Wang R, Hermjakob H. The Proteomics Identifications (PRIDE) database and associated tools: status in 2013. Nucleic Acids Res. 2013 Jan 1;41(D1):D1063-9. doi: 10.1093/nar/gks1262. Epub 2012 Nov 29. PubMed PMID:23203882.

Additionally and if it is feasible we'd like to ask our submitters to reference the dataset in a much abridged form in the abstract itself, like this: "The data have been deposited to the ProteomeXchange with identifier <PXD000xxx>."

See for example this Chromosome-Centric Human Proteome Project dataset and paper: <http://www.ncbi.nlm.nih.gov/pubmed/?term=23312004>, and other examples on PubMed. In our experience, a PX Identifier in the abstract makes the dataset much more visible and accessible.

## 9.3 Public release of the dataset

By default, your data will be made publicly available after your manuscript has been accepted, or when we have your instructions to do so. While we may also receive acceptance notifications from some journals, we would like to ask all submitters to kindly notify us separately. Otherwise, it can happen that we don't know that your manuscript is already published. You can notify us two ways:

A) Via the new PRIDE Archive web site (<http://www.ebi.ac.uk/pride/archive>). Once you have logged in with your user account at <http://www.ebi.ac.uk/pride/archive/login> you can click the green "Publish" buttons located next to your unpublished datasets. Here you can provide details for your dataset and submit a web form, please see Figure 40.

## Publish Project: PXD000223

Please provide us with some publication details if available

PubMed(s) (Optional, comma separated)

DOI(s) (Optional, comma separated)

Reference details, such as: title, authors (Optional)

**Confirm**

**Cancel**

**Figure 40:** Providing publication details using the PRIDE Archive web.

B) Contacting pride-support@ebi.ac.uk.

Upon making the project public, a project page will be released over at ProteomeCentral (<http://proteomecentral.proteomexchange.org>) and from a particular dataset page an FTP location will be available.

## 10 Appendix I: Definitions

Proteomics data come in a variety of forms, which are defined here:

- **Mass spectrometer output files:** the data and metadata generated by mass spectrometers, usually one file per run (although some instruments put multiple runs per file). The data may be the original profile mode scans or may already have had some basic processing like centroiding applied. They may be:
  - o i) raw data (as described below).
  - o ii) peak list spectra in a standardized format such as mzML, mzXML or mzData (see below), but they cannot be ‘processed peak lists’ (see below).

However, it is important that all of the scans that were generated are included with applicable metadata.

- **Raw data:** the binary, vendor-specific output files directly created by the instrument software. These files are typically large (several gigabytes) and require specialized software in order to be read.
- **Standardized MS data formats:** There are currently three widely known mass spectrometry data formats in Proteomics: mzXML (developed at the Institute of Systems Biology (ISB), Seattle, USA), mzData (now made obsolete, originally developed by the HUPO Proteomics Standards Initiative (PSI)), and the successor to both of the above: mzML (currently v1.1, jointly developed by the ISB and PSI, <http://www.psidev.info/mzml>). These data formats can be used to represent processed peak lists, as well as raw data. In addition to the mass spectra, they contain detailed metadata that provide context to the measurements.
- **Processed peak lists:** Heavily processed form of mass spectrometry data, usually derived from the raw data files through various (semi-)automatic steps, e.g.: centroiding, deisotoping, and charge deconvolution. These files are formatted in plain text, with typical formats like dta, pkl, ms2 or mgf. They usually contain only a subset of only the MS2 scans (MS1 scans are excluded), and are missing significant amounts of metadata that were present in the source format.
- **Protein/peptide identifications:** Proteomics mass spectra can be matched to peptides or proteins, resulting in identifications for those spectra. Typically a spectrum is considered identified if the score attributed to a peptide or protein match qualifies against an *a priori* or *a posteriori* defined threshold. In the case of fragmentation spectra, the initial identification will consist of a peptide sequence; subsequent steps will derive a list of proteins from the identified peptides. The protein assembly step can be a discernible process with its own input and output files, or it can be implicit in the overall identification software. This information can be represented by a variety of data formats called search engine output files (see below).

- **Protein/peptide quantification:** Protein/peptide expression values can also be obtained from a MS-based proteomics experiment. There is a high diversity of approaches that result in the existence of very heterogeneous software and data analysis pipelines. Some search engines are able to perform both identification and quantification, and produce 'search engine output files' containing both types of data. However, there is software that only performs the quantification part of the analysis and the generated data is represented in quantification software output files (see below).
- **Search engine output files:** They contain the data and metadata generated by the software (usually called search engines) used for performing the identification and quantification of peptides and proteins. Each search engine has its own specific output file. The formats are typically formatted in either plain text or XML, with typical formats like mascot.dat, OMSSA xml, etc. In addition to each specific format, a data standard format called mzIdentML (currently v1.1, <http://www.psidev.info/mzidentml>) has been developed by the PSI to represent this kind of information. Some search engine output files can represent as well quantification results, but this is not the case of mzIdentML. A second standard data format called mzTab (<http://code.google.com/p/mztab/>), currently under development, can represent both identification and basic quantification results.
- **Supported identification results:** This definition includes all protein/peptide identification processed data that can be fully represented by the receiving repository. For the PRIDE database, as the PX submission point for tandem MS/MS datasets, the data formats supported are PRIDE XML and mzIdentML version 1.1. It can represent both mass spectra data and protein/peptide identifications, and for some use cases in PRIDE XML, basic quantification information. Search engine output files need to be converted/exported to PRIDE XML or mzIdentML 1.1 to allow a full representation of the processed results in the PRIDE database and in the PX consortium.
- **Quantification software output files:** the data and metadata generated by the software used for performing exclusively the quantification analysis of peptides and proteins. In addition to each specific format from each software tool, a data standard format called mzQuantML (currently v1.0, <http://www.psidev.info/mzquantml>) is released by the PSI to represent this kind of information. As mentioned before, a second data format called mzTab (<http://code.google.com/p/mztab/>) can represent basic quantification results, although is currently not yet fully ratified.
- **Gel image files:** In case two-dimensional gel electrophoresis has been used as a separation method the gel image files generated.

**Metadata:** Whereas mass spectra present the core output of any mass spectrometer, a simple collection of spectra does not provide sufficient information for confident interpretation. Something similar happens for the

peptide and protein identifications and their expression values. This lack of context can be solved by providing relevant metadata along with the spectra and/or the identifications and quantification data. Mass spectrometer, search engine, and quantification software output files (see above) typically accommodate this information.

## 11 Appendix II: Available tools to help you with the submission

### 11.1 Creation of PRIDE XML files

#### 11.1.1 Tools developed by the PRIDE team

PRIDE Converter 2 (<http://code.google.com/p/pride-converter-2/>) is the most recent conversion tool developed by the team. It can work in batch mode and it can be integrated into automatic pipelines due to its modular software architecture. It is composed of 4 independent applications:

- The *PRIDE Converter 2* application will convert MS search result files containing identification and spectra into PRIDE XML.
  - The *PRIDE mzTab Generator* will produce skeleton mzTab files from MS search results files. At present, these skeleton files require either manual or scripted editing to add quantitation and/or gel information, but will be updated for automated insertion of quantitation results from different community file formats when the mzTab format is finalized.
  - The *PRIDE XML Filter* will remove identifications or spectra from PRIDE XML files based on a series of configurable filters.
  - The *PRIDE XML Merger* will combine several PRIDE XML files into a single one.
- List of the formats supported by PRIDE Converter 2 by November 2013 (Table 1).

Format Name	File Type	Data Content
Mascot	.dat	Spectra and Identifications
X!Tandem	.xml	Spectra and Identifications
OMSSA	.csv	Spectra and Identifications
SpectraST	.txt	Spectra and Identifications
CRUX	.txt	Spectra and Identifications
MSGF	.txt	Spectra and Identifications
Proteome Discoverer	.msf	Spectra and Identifications
DTA	.dta	Spectra Only
MGF	.mgf	Spectra Only
mzData	.xml	Spectra Only
mzXML	.xml	Spectra Only
PKL	.pkl	Spectra Only

**Table 1:** List of formats supported by PRIDE Converter 2.

Tutorials for general users and developers are available at the PRIDE Converter 2 Google Code page (<http://code.google.com/p/pride-converter-2/>).

#### 11.1.2 External tools developed by collaborators

1)- PeptideShaker ([peptide-shaker.googlecode.com/](http://peptide-shaker.googlecode.com/)). It can use as input Mascot .dat, X!Tandem XML and OMSSA .omx files.

2)- ProteinLynx Global Server (PLGS, Waters Corporation). It has an exporter to PRIDE XML from version 2.4 but with several limitations (metadata is not

properly annotated for some fields like submitter, species, etc). Improved support from version 3.0.

- 3)- OmicsHub Proteomics (Integromics).
- 4)- hEIDI (<http://biodev.extra.cea.fr/docs/heidi>). Local LIMS.
- 5)- Proteios (<http://www.proteios.org/>). A LIMS system developed by F. Levander's group (PubMedID: 19354269).
- 6)- EasyProt (<http://easyprot.unige.ch/>).
- 7) Protein Scape (Bruker).
- 8)- The ProteoRed MIAPE Extractor tool (<http://www.proteored.org/MIAPEExtractor>). It is able to generate fully MIAPE compliant (MS-MSI) PRIDE XML files containing much more detailed metadata than the minimal required by a ProteomeXchange submission.

## 11.2 Creation of mzIdentML files

As mentioned in the previous section, mzIdentML is the HUPO-PSI standard for protein/peptide identifications coming from MS-based proteomics approaches. The stable version is 1.1, which is supported by PRIDE and PX. It does not contain the mass spectra, which must be provided in external files referenced from the mzIdentML files (XML based files like mzML, mzXML or mzData, or peak lists like mgf, dta, ms2, or pkl).

At the time of writing, this is the list of software that can export mzIdentML v1.1 (see an updated list at <http://www.psidev.info/tools-implementing-mzidentml>):

- 1- Mascot (Matrix Science, <http://www.matrixscience.com/>). From version 2.4.
- 2- MS-GF+ (<http://proteomics.ucsd.edu/Software/MSGFPlus.html#pubs>).
- 3- Phenyx (GeneBio, <http://www.genebio.com/products/phenyx/>).
- 4- ProCon: Converter for Sequest .out, ProteomeDiscoverer (Thermo) v1.2/1.3/1.4 .msf files and ProteinScape 2.1 (Bruker) database content (<http://www.medizinisches-proteom-center.de/procon>).
- 5- TPP (pep.xml and prot.xml files): The idConvert tool from can be downloaded from ProteoWizard, or is bundled with the TPP directly starting with version 4.6.3.
- 6- X!Tandem and OMSSA: Using the mzidLibrary (<https://code.google.com/p/mzidentml-lib/>).
- 7- Scaffold (Proteome Software, <http://www.proteomesoftware.com/products/scaffold/>). From version 4.0.
- 8- OpenMS
- 9- MIAPE MSI Extractor (<http://proteored.org/miape/>, ProteoRed, Madrid)
- 10- PAnalyzer: Tool to perform protein inference analysis (<https://code.google.com/p/ehu-bio/wiki/PAnalyzer>).

11- Tools from D. Tabb lab: Myrimatch, Pepitome (spectral library search)<sup>15</sup> , TagRecon and IDPicker.

### 11.3 Checking the files before submission (initial quality assessment)

#### 11.3.1 Tool developed by the PRIDE team

PRIDE Inspector (<http://code.google.com/p/pride-toolsuite/wiki/PRIDEInspector>). This is an open source rich client application for inspecting MS-based proteomics data. Experiments can be examined based on different views emphasising either metadata, identified proteins or peptides, mass spectra, or quantification results.

Apart from its powerful visualization features, the major strength of PRIDE Inspector is the possibility to perform a first assessment of data quality using e.g. the 'Summary charts', which are generated based on different aspects of the data. Currently, PRIDE Inspector supports fast data retrieval on standard file formats: mzML, mzIdentML (plus the corresponding peak list files) and PRIDE XML. In addition, it also gives the user direct access to a PRIDE public database instance. As a key point, it provides journal reviewers/editors access to (privately available) experiments during the review process.

#### 11.3.2 External tool developed by collaborators

- 1) PRIDE Viewer (<http://proteo.cnb.csic.es/prideviewer/>). It can visualize PRIDE XML files.
- 2) mzML validator (link to Java Web Start to be done if necessary): a Java-based tool to validate semantics and MIAPE compliance of mzML files.
- 3) mzIdentML validator (<http://psi.pi.googlecode.com/svn/trunk/validator/trunk/mzid-validator.html>): a Java-based tool to validate semantics and MIAPE compliance of mzIdentML files.
- 4) ProteoRed MIAPE Extractor tool workflow (<http://www.proteored.org/MIAPEExtractor>): After the MIAPE information, data can be integrated, inspected and validated before the PRIDE XML creation.

### 11.4 File submission to PRIDE

As described before in this tutorial, the PX Submission Tool can be used (<http://www.proteomexchange.org/submission>). It creates the relations between the different types that can be part of a dataset and uploads the data into PRIDE via FTP.

## 12 Appendix III: Summary of formats supported by PRIDE for PX MS/MS submissions

### a) as raw data

Formats supported:

- mzML, mzXML, mzData. These files must not be heavily processed to be considered ‘raw’.
- Thermo .RAW, ABSCIEX .wiff, .wiff.scan, Agilent .d/, Waters .raw/
- imzML, Shimadzu .run/, Bruker .yep, Bruker .baf

All peak lists formats (mgf, dta, ms2, pkl) can be supported but they will not be considered raw data. They will be considered as ‘peak list processed files’ or simply ‘peak’.

### b) as processed identification results'

Two formats are now supported: PRIDE XML and mzIdentML.

b.1) PRIDE XML: Different search engine output files need to be converted to PRIDE XML using existing tools like PRIDE Converter 2 (<http://code.google.com/p/pride-converter-2/>) and others (see Appendix 2). Formats supported:

- Tandem XML
- OMSSA .csv.
- Mascot .dat
- Sequest Crux .txt
- SpectraST .xls
- ProteomeDiscoverer .msf files.
- All accompanying peak lists formats.

b.2) mzIdentML (version 1.1): There are a number of tools that can export mzIdentML 1.1 (see Appendix 1). Formats supported this way:

- Tandem XML (using mzidLibrary, <https://code.google.com/p/mzidentml-lib/>)
- OMSSA .csv (using mzidLibrary, <https://code.google.com/p/mzidentml-lib/>).
- Mascot .dat ( direct export functionality available from Mascot 2.4).
- Sequest .out files (using the ProCon tool, <http://www.medizinisches-proteom-center.de/procon>).
- ProteomeDiscoverer .msf files (using the ProCon tool, <http://www.medizinisches-proteom-center.de/procon>).
- ProteinScape 2.1 (Bruker) database content (using the ProCon tool, <http://www.medizinisches-proteom-center.de/procon>).

- MS-GF+ (direct export functionality available).
- Phenyx (direct export functionality available).
- Trans-Proteomic Pipeline (pep.xml files). The idConvert tool from can be downloaded from ProteoWizard, or is bundled with the TPP directly starting with version 4.6.3.
- Scaffold (direct export functionality available). From version 4.0.
- OpenMS output.
- MIAPE MSI Extractor output (<http://proteored.org/miape/>, ProteoRed, Madrid)
- PAnalyzer output: Tool to perform protein inference analysis (<https://code.google.com/p/ehu-bio/wiki/PAnalyzer>).
- Output files from Myrimatch, Pepitome (spectral library search) , TagRecon and IDPicker.
- All accompanying peak lists formats.

### **c) as search engine output files**

Only those data formats that cannot be converted/exported to PRIDE XML/mzIdentML are considered to be ‘unsupported formats’ and can use this alternative approach (datasets type B, Datasets containing raw data and search engine output files). At present, there are no reliable converters to PRIDE XML/mzIdentML for the following formats amongst others:

- MaxQuant output files,
- ProteinPilot .group files

### **d) as quantification results**

The current version of pipeline does not support a full and standard representation of the quantification results, linked to the identification results (unless this information is provided in PRIDE XML files. This can be done using PRIDE Converter 2). It is expected that data standards for quantitative proteomics data (mzQuantML, mzTab) will be supported in the future. However, any quantification result output files can be submitted as accompanying ‘QUANT’ files.

### **e) as gel images**

Gel images (in any format) tagged as ‘GEL’ can be included in the submission.

### **f) as others**

Any other type of files are optional and can be supported as part of a PX submission together with the other files.



## 13 Appendix IV: Metadata requirements for MS/MS submissions

Proteomics data are substantially enriched when sufficient metadata are provided. Metadata will be as extensive as possible and will aim to comply with the MIAPE (Minimum Information About a Proteomics Experiment) guidelines. However, the presence of the metadata required in this Appendix will be enforced for any PX submission (they are mandatory in the PX Summary File format). They can be provided using the PX Submission tool.

The user will need to provide:

- Contact name and e-mail for the submission. The contact details of the data submitters need to be provided, allowing interested users to contact the original authors if desired.
- Lab Head or Principal Investigator.
- Name of the PX dataset.
- Project description: it could be considered as the abstract information of the dataset (provided as free text).
- Summary of the Sample Protocol (provided as free text).
- Summary of the Data analysis Protocol (provided as free text).
- Experiment type. Chosen from a drop-down menu.
- Keywords: A list of keywords that describe the content and type of the experiment being submitted. Multiple entries should be comma separated.
- Sample annotation: species. At least one NEWT Controlled Vocabulary (CV) term is mandatory per dataset.
- Sample annotation: tissue. Using the BRENDA Tissue ontology (BTO), accessible at  
<http://obo.cvs.sourceforge.net/obo/obo/ontology/anatomy/BrendaTissue.obo>
- Instrument details. Using the PSI-MS CV. It is accessible at  
<http://psidev.cvs.sourceforge.net/viewvc/psidev/psi/psi-ms/mzML/controlledVocabulary/psi-ms.obo>.
- Quantification method (if applicable).
- Protein post-transcriptional modifications (PTMs). They are reported using the PSI-MOD ontology (accessible at  
<http://psidev.cvs.sourceforge.net/psidev/psi/mod/data/PSI-MOD.obo>).

### Optional information:

- Sample annotation: cell type. Use the “Cell Type” ontology.
- Sample annotation: Disease. Use the “Human Disease” ontology (DOID).
- Dataset optional details:
  - o your dataset is part of a bigger project/effort (for instance the Human Proteome Project or ‘PRIME-XS’). It is a way to tag your dataset to enable grouping this way.

- there is already a PubMed ID associated with it (the data has been already published).
- your dataset represents a reanalysis of an earlier public PX dataset
- there are other “omics” datasets (for instance transcriptomics, metabolomics data present in other repositories) that can be associated with it. In this case, please provide the accession number of the dataset in the corresponding repository.

## 14 Appendix V: Recommended Partial Submission search engine identification results for particular software tools

There are software tools and workflows with search results for which there are not available exporters to PRIDE XML. In these case search/peptide/protein identification results can be provided in the form of partial submissions.

Here we describe the workflow for two popular tools: MaxQuant (PubMed ID: 19029910) and ProteinPilot™ (AB SCIEX).

### 14.1 MaxQuant

If you are using the latest version of MaxQuant (1.3.0.5) there is a txt folder generated and by default you can just zip this text folder and upload as a 'SEARCH' file.

If this is complicated, we would recommend uploading the following particular text output files:

parameters.txt  
 peptides.txt  
 modifiedPeptides.txt  
 proteinGroups.txt

and your 'Experimental Design Template file' saved as a tab delimited file.

### 14.2 ProteinPilot

For ProteinPilot as peptide/protein identification files we strongly recommend providing human readable files instead of the binary '.group' file. Please export the group files into XML files using:

<http://www.absciex.com/products/software/proteinpilot-software>

"Command Line Control and Open Results. To support users and third-party software vendors that want to integrate ProteinPilot Software, it is possible to script searches *via* command line and decrypt the '.group' file results into clear XML for full access to all the data it contains."

Here is a 'how to' on the conversion process from one of our submitters:

1. Create a txt file in Notepad entitled say "group2XML\_Example.bat.txt" and save it in the ProteinPilot folder (where the group2xml.exe is located).
2. Rename "group2XML\_Example.bat.txt" to "group2XML\_Example.bat", giving it a Windows batch file extension.

3. Open this batch file in 'Notepad' and type in the following command line instructions:

group2XML.exe XML <full path to the .group file to be converted> <full path to the .xml file the .group file will be converted into>

for instance

```
group2XML.exe XML "C:\AB SCIEX\ProteinPilot Data\Results\Example.group"  
"C:\AB SCIEX\ProteinPilot Data\Results\Example.xml"
```

The command has the following argument structure: group2XML.exe <Type> <Result.group> <Output.file>

where:

- <Type> specifies the type of output.
- <Result.group> is a .group file created by ProteinPilot Software.
- <Output.file> is the name of the file to be created.

4. Save and close the file.

5. Double-click on the file to run the conversion.