

ProCon – a Proteomics Conversion Tool for Standardized File Generation and Repository Submission

Gerhard Mayer, Christian Stephan, Helmut E. Meyer, Martin Eisenacher*
Medizinisches Proteom-Center, Ruhr-Universität Bochum, Germany, * e-Mail: martin.eisenacher@rub.de

Introduction

ProCon – A Proteomics Conversion Tool

In an EU FP6 project (ProDaC, www.fp6-prodac.eu) a major goal was the development of software tools to export, submit and import data sets using standard formats. In a current EU FP7 project (ProteomeXchange, www.proteomexchange.org) a standardized submission pipeline is developed for unique identifier generation and data exchange between large proteomics repositories. The ProCon tool, implemented at the Medizinisches Proteom-Center (MPC) in Bochum (<http://www.medizinisches-proteom-center.de/ProCon>), converts data into standard file formats (mzIdentML, PRIDE XML) enabling data submission through the ProteomeXchange pipeline.

ProCon User Interface Examples

The screenshot shows the ProCon application window. The main window has a menu bar (File, Help) and a toolbar with buttons for 'Sequest.out to mzIdentML', 'Proteome Discoverer to mzIdentML', 'ProteinScape 1.X Source', 'General Source', 'PRIDE XML', and 'ProteinScape 2.1 to mzIdentML'. Below the toolbar are four numbered steps: 1) Select folder with Sequest .out files, 2) Parse SEQUEST out folder, 3) Select mzIdentML output file, and 4) Export to mzIdentML file. A 'Ready.' status bar is visible. A secondary window titled 'Conversion of ProteinScape 2.1 -> mzIdentML 1.1' is open, showing connection data (Server: maldraunserver, User: sa, Password: ****), organization data (Researcher Name: Gerhard Mayer, Organization Name: Medizinisches Proteom Center (MPC)), and conversion parameters (Calculate theoretical m/z values, Calculate isoelectric points, pI calculation: Consensus). A table of search events is also displayed.

Mapping SloMo Phosphorylation Sites to PRIDE Files during ProteinScape 1.3 Conversion

Implementing a user request Phosphorylation information of the SloMo tool [2] can be integrated into a PRIDE file by loading specially formatted csv columns during ProteinScape 1.3 conversion.

The screenshot shows a table of phosphorylation sites with columns: Frag.Run, Cmpd, Charge, Peaks Per Window, Peptide Score, Peptide Seq, Score, Spectrum, Notes, Run, Cmpd, Charge, Fragmentation. Below the table is the ProCon interface with a dialog box titled 'Loaded instrument and reference information overwrites all corresponding sections in exported PRIDE XML files:'. The dialog has three sections: 'Load mzDataInstrument XML file' (file: /Instruments/MPC_OrbiTrap.xml), 'Load References' (file: /config/references.properties), and 'Load SloMo .csv' (file: /config/SloMo.csv, score threshold: 19.0).

ProCon Functionality

- conversion of Proteome Discoverer® (Thermo Fisher Scientific Inc., San Jose, CA) 1.2 and 1.3 .msf and .prot.xml file information into mzIdentML
- conversion of peptide information of Sequest™ .out files (Thermo Fisher Scientific Inc., Waltham, MA, USA) into mzIdentML
- conversion of ProteinScape® (Bruker Daltonik GmbH, Bremen, Germany) version 1.3 results (SearchEvents and gels) to PRIDE XML
- conversion of ProteinScape 2.1 results (SearchEvents) to mzIdentML

A Java API was developed for accessing the SQL database content of ProteinScape. By using this API the Analyte tree of a project is parsed and the user chooses the samples and results. Then the mzIdentML file is built up in the main memory and ultimately written by using the JAXB (Java Architecture for XML binding) marshaller functionality. During this process also the relevant CV terms are built in by accessing the .obo files via the Ontology Lookup Service [1] (OLS).

ProteomeDiscoverer® to mzIdentML Conversion

After (1) setting meta-project information (such as institute and contact e-mail), (2) specifying the .prot.xml and .msf file paths, and (3) other conversion parameters (such as method for isoelectric point calculation), the conversion into mzIdentML can be initiated and is processed automatically. Internally the Open Lookup Service (OLS) is used to get the needed CV terms and the NCBI taxonomy identifiers used are acquired by accessing the Entrez Programming Utilities. (<http://www.ncbi.nlm.nih.gov/books/NBK25501>).

The screenshot shows the 'Conversion parameters for PD -> mzIdentML conversion' dialog. It has fields for 'Proteome Discoverer *.msf file', 'Proteome Discoverer *.prot.xml file', and 'Name of mzIdentML 1.1 file to generate'. Below are 'Organization data' fields (Name, Contact address, Contact email) and 'Conversion parameters' (Calculate theoretical m/z values, Calculate isoelectric points, pI calculation: Consensus).

ProCon Perspectives

- The current ProteomeDiscoverer® conversion works both with ProteomeDiscoverer® versions 1.2 and 1.3 (both the .prot.xml and .msf files are needed for getting the information about the protein ambiguity groups / protein inference information). It's planned to accelerate the conversion for version 1.3, because here all the relevant information is contained and easily accessible in the .msf (SQLite) database. It will also integrate the ProteomeDiscoverer API [3].
- conversion of ProteinScape 2.1 gel results to mzIdentML
- conversion of ProteinScape 2.1 results to PRIDE XML
- mzQuantML conversion, e.g. for spectral counting workflows

Summary

ProCon complements other converters (e.g. <http://code.google.com/p/pride-converter-2>) to make data conversion easy even for laboratories without bioinformatics expertise. That supports journal submission and community data access. Thus meta-analyses of proteomics data and linkage to other 'omics data will lead to new insights regarding biomarkers or functional disease mechanisms.

ProCon-related Websites

<http://www.medizinisches-proteom-center.de/ProCon>

Submission: <http://www.proteomexchange.org>

Standards: <http://www.psdev.info>

References

1. Cote RG, Jones P, Martens L, Apweiler R, Hermjakob H. The Ontology Lookup Service: more data and better tools for controlled vocabulary queries. *Nucleic Acids Res.* 2008 Jul 1;36(Web Server issue):W372-6. Epub 2008 May 8 [PMID: 18467421]
2. Bailey CM, Sweet SM, Cunningham DL, Zeller M, Heath JK, Cooper HJ (2009) SloMo: automated site localization of modifications from ETD/ECD mass spectra. *J Proteome Res.* 2009 Apr;8(4):1965-71 [PubMed PMID: 19275241]
3. Colaert N, Bartsnes H, Vaudel M, Helsens K, Timmerman E, Sickmann A, Gevaert K, Martens L (2011) Thermo-msf-parser: an open source Java library to parse and visualize Thermo Proteome Discoverer msf files. *J Proteome Res.* 2011 Aug 5;10(8):3840-3 [PMID: 21714566]

Funding

This work is funded by the European Commission within the 7th European Union Research Framework Programme (Contract No.: 260558). ME is funded by P.U.R.E. (Protein Unit for Research in Europe), a project of North Rhine-Westphalia, a federal state of Germany.