



Protein Annotator: Freeware for the Generation of UniProt Formatted XML Files Containing Custom User-Specified Features

Joseph B. Greer, Ryan T. Fellers, Richard D. LeDuc, Bryan P. Early, Alexandra J. van Nispen, Paul M. Thomas, Neil L. Kelleher
Northwestern University, Evanston, IL

Background

- ✓ The National Resource for Translational and Developmental Proteomics (NRTDP) has developed a beta-version tool that allows users to create custom UniProt formatted XML files
- ✓ It allows users to add unannotated point and range features
- ✓ Exports UniProt formatted XML files that ProSightPC can use to create databases

Introduction

Protein Annotator is a free application allowing researchers to create UniProt formatted files with specific, custom features. When searching top-down data against shotgun-annotated databases, protein entries often do not contain the specific features such as custom PTMs, ADCs, glycans, disulfide bonds, or endogenous cleavages that researchers are interested in. This significantly limits their ability to search for novel proteoforms. Additionally, searching for a specific proteoform is difficult when an entry is heavily modified. Both situations require a candidate database created from custom entries. Unfortunately, manually creating UniProt formatted files containing custom entries is tedious and prone to syntax errors. Protein Annotator allows researchers to create custom, valid UniProt files – making it easier to perform targeted proteoform searches.

Application Workflow



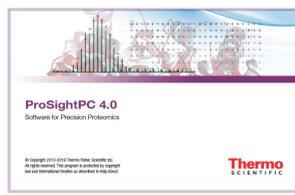
Input

- UniProt formatted FASTA, text or XML file

XML file

Output

- UniProt formatted XML containing user defined custom features that can be used as input in tools such as ProSightPC 4.0



User Interface

The screenshot shows the Protein Annotator application window. At the top, there are navigation buttons: Return to home, Add isoforms from file, Add isoform, Remove all isoforms, and Save isoforms to file. Below this is the 'Isoforms' table with columns for Accession, Description, and Modify. The main area displays a protein sequence with various amino acids highlighted in colored boxes representing features. At the bottom, there are two panels: 'Point Features' and 'Range Features'. The Point Features panel shows a list of common features like Acetylation, Monomethylation, Dimethylation, and Trimethylation. The Range Features panel shows a table of endogenous cleavages and disulfide bonds.

The user interface consists of four sections: the isoform table, the fragment map, the point feature panel, and the range feature panel. In this example, the input file contained entries associated with Thermo Fisher Scientific's Pierce Intact Protein Standard Mix.

Isoform Table

Accession	Description	Modify
PODN71	Protein ADP-ribosyltransferase	[edit] [delete] [copy]
PODN69	Invalid feature on entry PODN69. Feature cannot be mapped.	[edit] [delete] [copy]
PODN72	Lipoate--protein ligase	[edit] [delete] [copy]
PODN70	Protein-ADP-ribose hydrolase	[edit] [delete] [copy]
PODN73	Uncharacterized oxidoreductase SpyM50865	[edit] [delete] [copy]

Isoforms can be filtered and error messages are displayed.

Important Concepts

- .FASTA**
FASTA file
Text based file format. Application requires that FASTA files follow the UniProt format - entries contain a description line and the canonical protein sequence
- .txt**
Text file
Text based file format. Application requires that text files follow the UniProt format - entries contain the canonical protein sequence and annotated modifications
- .XML**
XML file
Extensible Markup Language based file format. Application requires that XML files follow the UniProt format - entries contain the canonical sequence and annotated modifications
- Point Feature**
Point Feature
Modification that only applies to a single amino acid such as PTMs and cSNPs
- PTM**
Post-Translational Modification
Processing event resulting from the addition or subtraction of a modifying group to one amino acid
- cSNP**
Coding Single Nucleotide Polymorphism
Single nucleotide variation at a specific position that changes the amino acid sequence
- Range Feature**
Range Feature
Modification that applies to a range of amino acids
- Endogenous Cleavage**
Endogenous Cleavage
Post translational cleavage event resulting in a new proteoform
- Pierce Standard**
Pierce Intact Protein Standard Mix™
Thermo Fisher Scientific developed mixture of six highly pure recombinant proteins used as an intact and top-down standard

Future Work

- Continue testing and work towards a stable public release
- Improve range feature visualization
- Expand point feature options

Acknowledgements

The authors would like to thank the members of the Top-Down Proteomics Development Team at Northwestern University. This research was supported by the National Institute of General Medical Sciences, and the National Institutes of Health under grant P41 GM108569. **Conflict Statement:** Some of the authors are involved with software commercialization.

<http://proteinannotator.northwestern.edu>

Point Features

The screenshot shows the Point Features panel with tabs for PTMs 4/5, Custom (0), cSNPs (0), and Glycans (0). It includes a list of common features (Acetylation, Monomethylation, Dimethylation, Trimethylation) and uncommon features (Phosphorylation). Below this is a grid of amino acid letters (A-Z) with a selected 'G'. Another section shows a table for custom modifications with columns for Name, Mono Mass, Ave Mass, and Formula.

Range Features

The screenshot shows the Range Features panel with tabs for Endogenous Cleavages (4) and Disulfide Bonds (6). It includes a table with columns for Name, Start, and End. Below this is another table for custom modifications with columns for Name, Start, and End.