

Rapid Proteomics Analysis Algorithm Development Using Proteome Discoverer: A Development Environment for Proteomics Scientists

Iman Mohtashemi; Kai Fritzsche; Hans Grensemann; Bernard Delanghe; David Horn; Torsten Ueckert

Thermo Fisher Scientific, San Jose, CA



Overview

Purpose: Here we present the Thermo Scientific Proteome Discoverer software, a node-based framework for proteomics data processing. We will explain the general structure of the software and provide multiple examples of creation of nodes for use in this system.

Methods: All data were acquired on a hybrid ion trap-Orbitrap[™] mass spectrometer system equipped with electron transfer dissociation (ETD). New nodes were developed in C# using Visual Studio 2010 and their utility were demonstrated by incorporating them into Proteome Discoverer[™] workflows.

Results: Development, deployment and analysis of various node implementations are discussed.

Introduction

The speed and ultra high resolution acquisition of Orbitrap-based technologies have made the characterization of increasingly larger number of previously undetectable features in a single run challenging. Data analysis bottlenecks are emerging both with respect to the speed of analysis and the depth of converting large data sets into meaningful information. While traditional database searching techniques still remain the mainstream of traditional proteomics analysis, other analysis tools are continuously improving in various areas. The breadth and depth of algorithms from academic labs to commercial organizations are daunting yet an infrastructure does not exist to rapidly develop, deploy and share unique algorithms amongst proteomics scientists. A node-based development environment is proposed to test and deploy analysis algorithms without unnecessary overhead.

Methods

Deploying New Nodes Into the Proteome Discoverer framework

Proteome Discoverer software is written in the Microsoft C#/NET environment and runs under Microsoft Windows. It consists of a server component for processing the data with user defined processing workflows and a client component for creating and scheduling processing workflows, and creating and displaying the final reports. The workflows consist of processing nodes that perform in a specific task in the data analysis pipeline. These processing nodes are implemented as plug-ins, making it easy to extend the data processing with new algorithms or functionality.

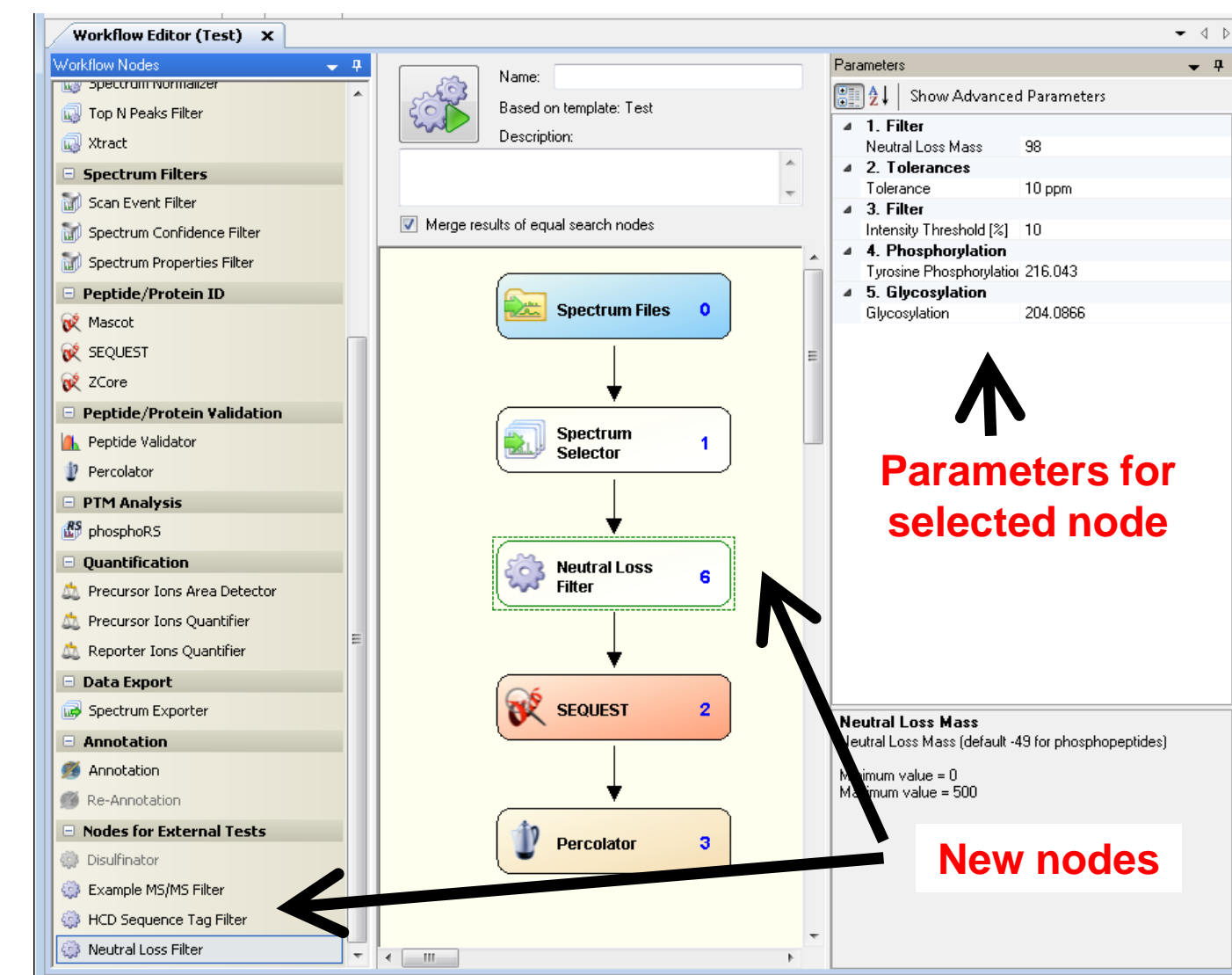
To create a new node for use in the Proteome Discoverer framework, a developer will start with a Microsoft Visual Studio 2010 solution that contains a template project that performs similar functionality to the node to be added. In the solution, the developer defines the node name, the input and output data types, whether or not the algorithm is tied to the standard license or a separate one, and the algorithm parameters and their acceptable limits. Figure 1 shows the code created for the neutral loss filter example, where the input and output for the node are MS/MS spectra and the appearance of the node is tied to the standard license in Proteome Discoverer software and thus will be available for all users. The code in Figure 1 show the parameters that are available for users in the Proteome Discoverer Workflow Editor user interface for that node and their acceptable limits. If the user sets any value outside these limits, the workflow will exit with an immediate error.

FIGURE 1. The code below defines the name of the node, the input and output data structures, and the parameters available for the users.

```

[ProcessingNode("Filter", 2274-4149-4003-15083C8B8447",
Category = "Filter",
Description = "Filter",
Definition = "Neutral loss filter",
Implementation = "Filter",
MaximumValue = "100",
MinimumValue = "0",
Position = 11)
[ParameterOf(
Category = "Filter",
Definition = "Neutral loss mass",
Description = "Neutral loss mass",
Implementation = "Filter",
MaximumValue = "99",
MinimumValue = "0",
Position = 12)
[ParameterOf(
Category = "Filter",
Definition = "Tolerance",
Description = "Tolerance",
Implementation = "Filter",
MaximumValue = "10",
MinimumValue = "0",
Position = 13)
[ParameterOf(
Category = "Filter",
Definition = "Intensity threshold",
Description = "Intensity threshold",
Implementation = "Filter",
MaximumValue = "10",
MinimumValue = "0",
Position = 14)
[ParameterOf(
Category = "Filter",
Definition = "Tolerance",
Description = "Tolerance",
Implementation = "Filter",
MaximumValue = "10",
MinimumValue = "0",
Position = 15)
    
```

Figure 2. Proteome Discoverer Workflow Editor with the new nodes

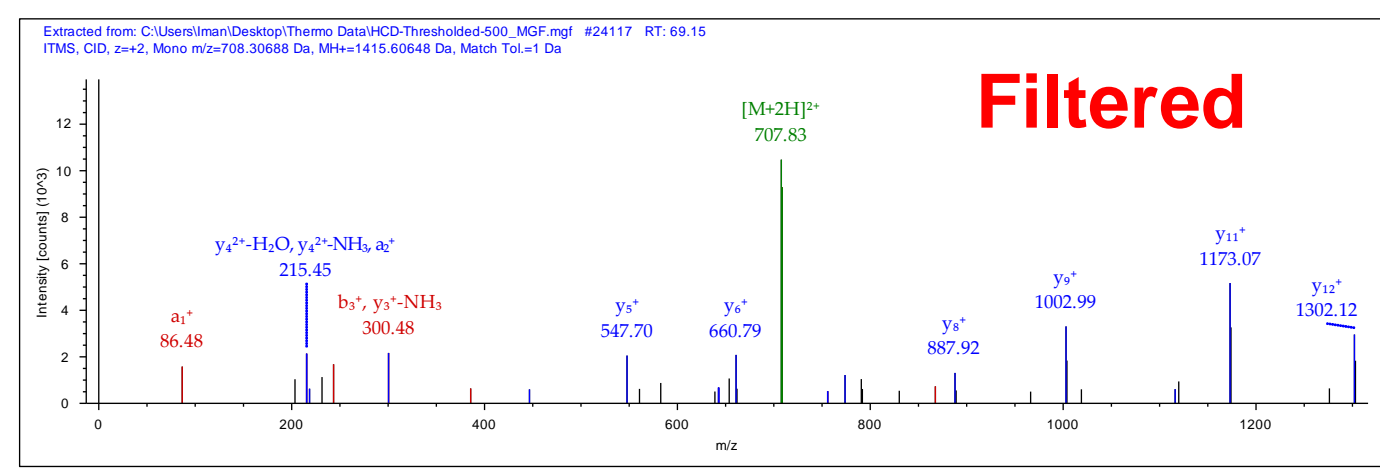
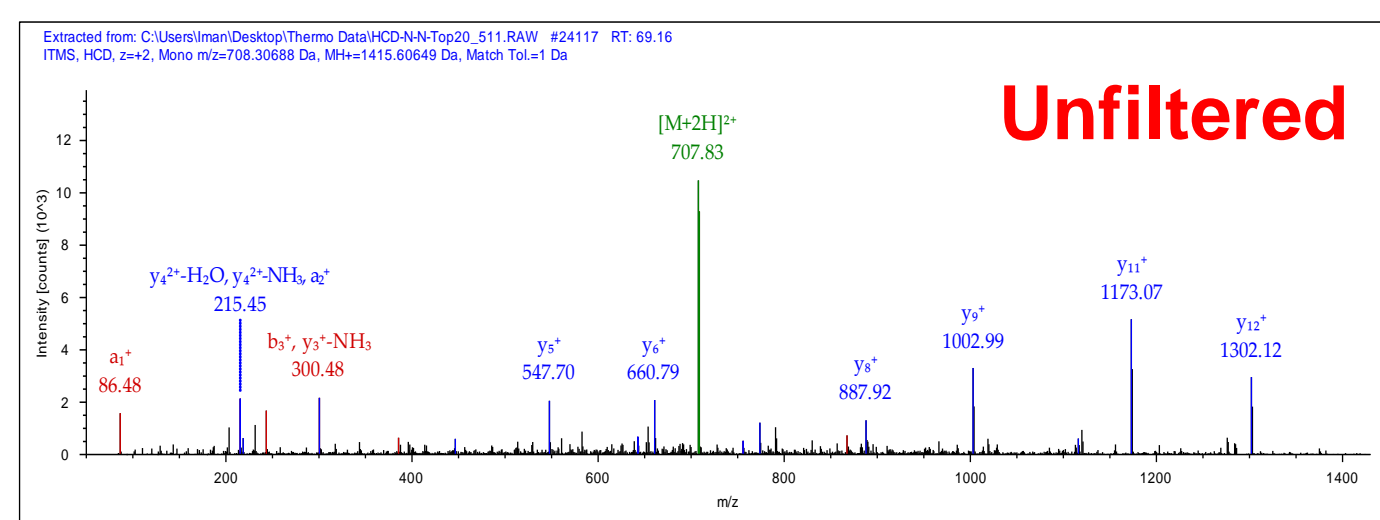


Results

Normalized ms² filter node

Each ms/ms spectra was normalized to the most intense peak within the spectrum. A threshold was then applied to achieve maximum identification rate with minimal search input (fig 3). Display of annotated spectra is seamlessly integrated without affecting any other components in Discoverer

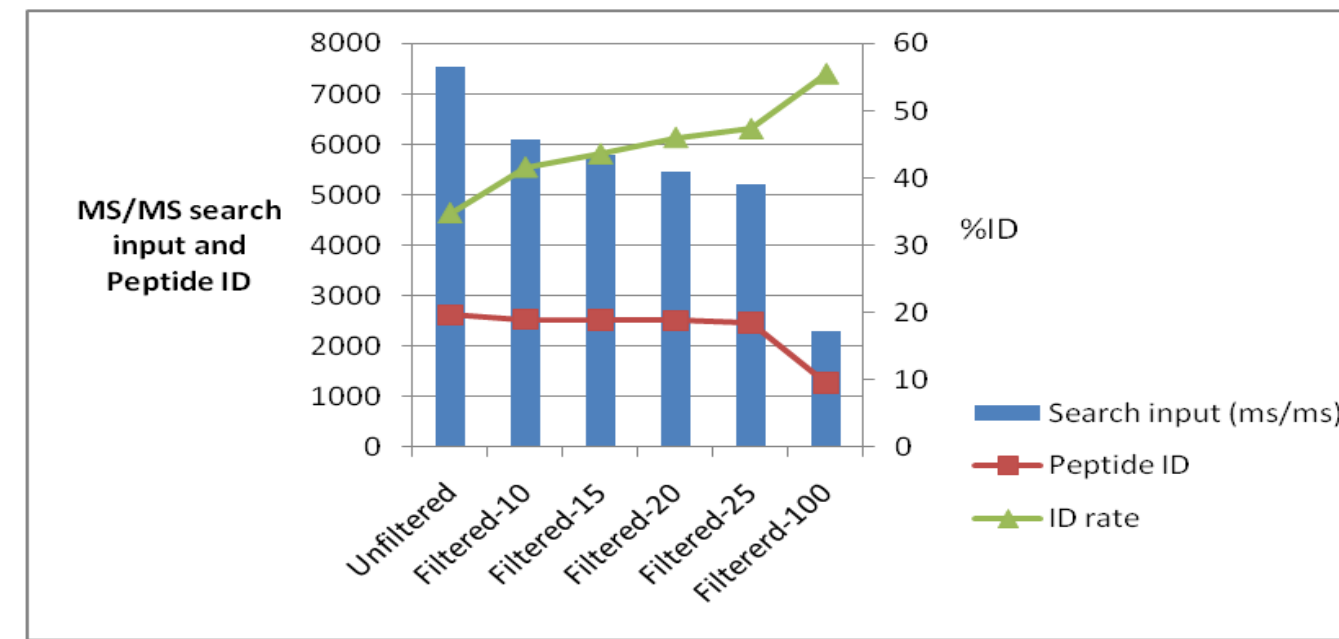
Figure 3. Filtered ms² spectrum



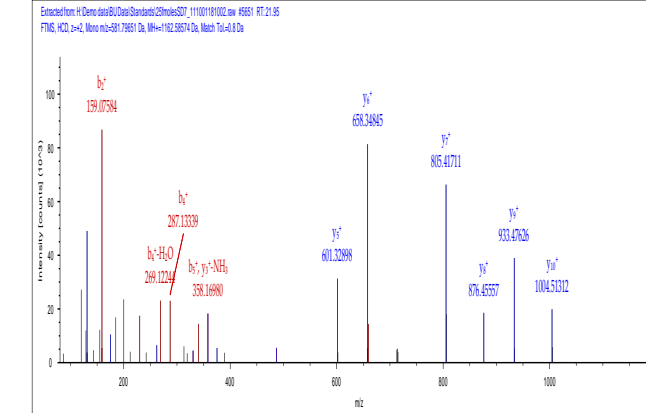
HCD Sequence tag filter

A peptide sequence tag filter was applied to HCD ms2 data to minimize search input for non-peptidic fragments (fig4)

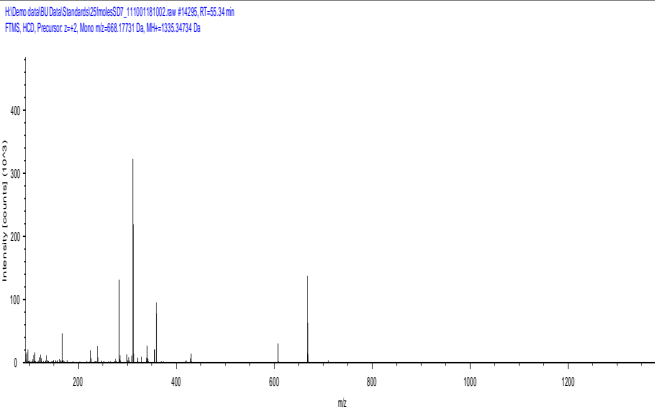
Figure 4. HCD sequence tag filter ID vs input



Peptidic Fragment



Non-peptidic Fragment



CID ms² neutral loss

CID ms2 screened for phosphorylation via neutral loss monitoring (fig5)

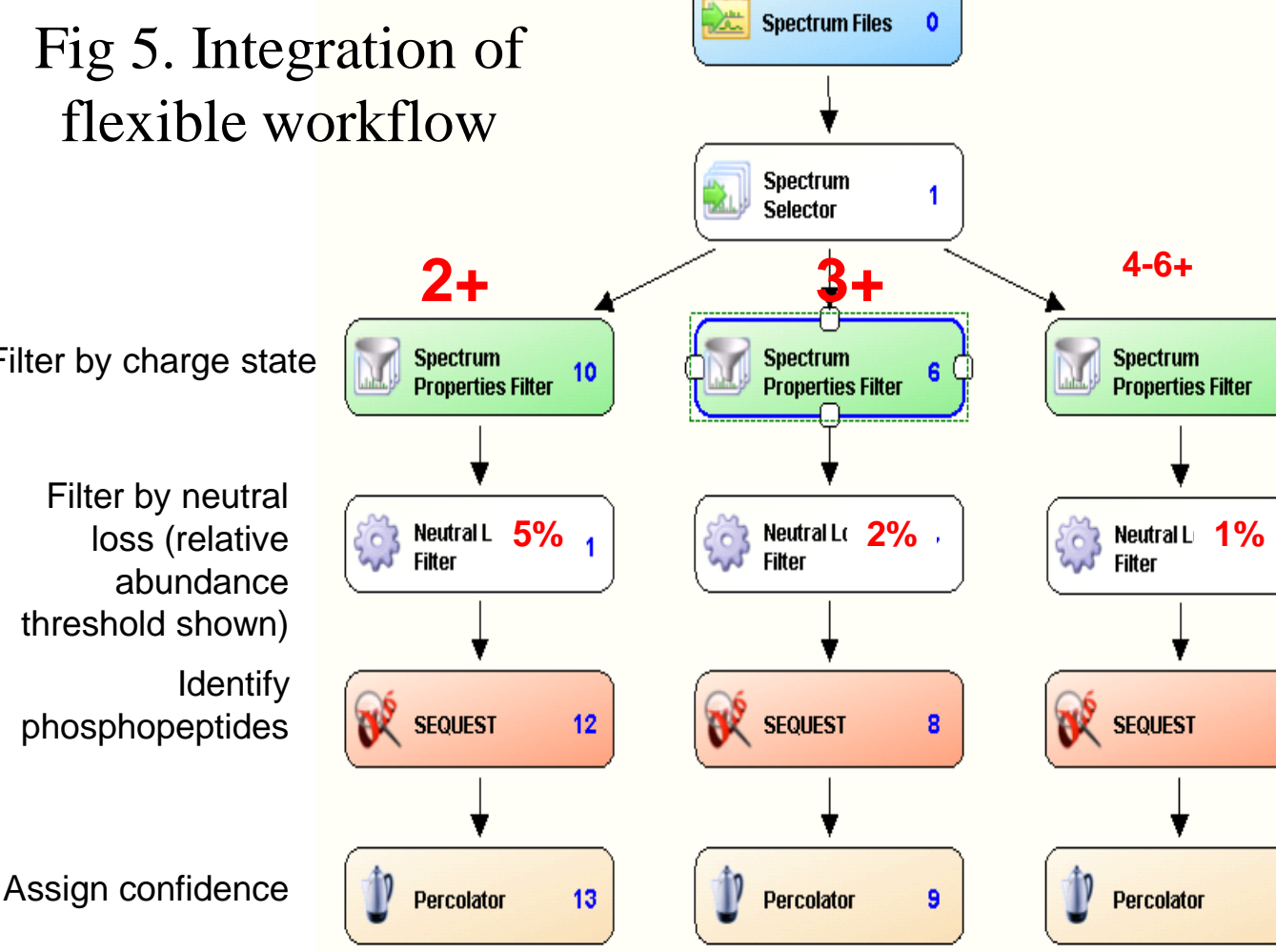


Fig 5. Integration of flexible workflow

Filter by charge state

Filter by neutral loss (relative abundance threshold shown)

Identify phosphopeptides

Assign confidence

ETD triggered ms³ for Disulfide Mapping

Custom data fields for data display for disulfide mapping (fig 6). Each new field is assigned a new GUID

Figure 6. Custom data fields for disulfinator node

```

CustomDataField crosslinkedPeptides = ProcessingServices.CustomDataService.GetOrCreateCustomDataField(workflow.ID, new Guid("6C382FDA-C6AA-4CDB-A390-AB4D9B0011C8"), "Crosslinked Peptide(s)", ProcessingModeNumber, searchModeInfo.ProcessingModeNumber, CustomDataTarget.Peptide, CustomDataType.String, CustomDataAccessMode.Read, DataVisibility.Visible);
CustomDataField massCrosslinkedPeptide = ProcessingServices.CustomDataService.GetOrCreateCustomDataField(workflow.ID, new Guid("4866912F-E945-4AD1-872C-880934402AEF"), "Mass Crosslinked Peptide", ProcessingModeNumber, searchModeInfo.ProcessingModeNumber, CustomDataTarget.Peptide, CustomDataType.Double, CustomDataAccessMode.Read, DataVisibility.Visible, plotType: PlotType.Numeric, format: "0.000");
CustomDataField calculatedMassCrosslinkedPeptide = ProcessingServices.CustomDataService.GetOrCreateCustomDataField(workflow.ID, new Guid("87C0A83E-7C8D-4554-A995-C77E3CD78966"), "Calculated Mass Crosslinked Peptide", ProcessingModeNumber, searchModeInfo.ProcessingModeNumber, CustomDataTarget.Peptide, CustomDataType.Double, CustomDataAccessMode.Read, DataVisibility.Visible, plotType: PlotType.Numeric, format: "0.000");
CustomDataField deltaMassCrosslinkedPeptide = ProcessingServices.CustomDataService.GetOrCreateCustomDataField(workflow.ID, new Guid("D80C2033-7C3D-4B66-84CD-CE3227E27897"), "Delta Mass Crosslinked Peptide [PPM]", ProcessingModeNumber, searchModeInfo.ProcessingModeNumber, CustomDataTarget.Peptide, CustomDataType.Double, CustomDataAccessMode.Read, DataVisibility.Visible, plotType: PlotType.Numeric, format: "0.000");
    
```

Figure 7. Custom data fields for disulfinator

Sequence	m/z [D]	MH+ [D]	ΔM [ppm]	Order	XCX	Crosslinked Peptide(s)	Mass Obs	Calculated	Delta Mass	Peptide 1
MPCTEDYLSILNR	834.41071	1667.81413	8.55	337	2.59	MPCTEDYLSILNR-EKCFVGRK	2735.292	2735.283	-0.902	EKCFVGRK
LPFHADICLPDTEK	925.95300	1859.89972	-0.36	335	6.65	MPCTEDYLSILNR-LPFHADICLPDTEK	3515.702	3515.690	-1.595	MPCTEDYLSILNR
MPCTEDYLSILNR	834.40985	1667.81243	-0.47	333	2.49	MPCTEDYLSILNR-LPFHADICLPDTEK	3515.702	3515.690	-1.595	LPFHADICLPDTEK
MPCTEDYLSILNR	834.41162	1667.81597	1.65	329	1.16	MPCTEDYLSILNR-RPFSALTPDEYIK	3488.705	3488.690	-0.346	RPFSALTPDEYIK
MPCTEDYLSILNR	834.41138	1667.81548	1.98	327	0.91	MPCTEDYLSILNR-RPFSALTPDEYIK	3488.705	3488.690	-0.647	RPFSALTPDEYIK
RPFSALTPDEYIK	812.45518	1623.90508	2.94	327	0.91	MPCTEDYLSILNR-RPFSALTPDEYIK	3488.703	3488.690	-0.647	MPCTEDYLSILNR
MPCTEDYLSILNR	834.41022	1667.81316	-0.03	327	1.44	MPCTEDYLSILNR-LQLHEK	2516.256	2516.250	-0.412	LQLHEK
MPCTEDYLSILNR	834.41205	1667.81682	1.16	315	2.15	QNCQDEK-MPCTEDYLSILNR	2676.219	2676.210	-0.411	QNCQDEK
MPCTEDYLSILNR	834.41046	1667.81355	0.26	311	2.17	QNCQDEK-MPCTEDYLSILNR	2676.217	2676.210	-0.427	QNCQDEK
DAPENLRPLTADFAEDK	1201.09050	2401.17388	6.07	305	0.81	DAPENLRPLTADFAEDK-LPFHADICLPDTEK	4249.044	4249.036	-1.932	LPFHADICLPDTEK
DAPENLRPLTADFAEDK	1201.09050	2401.17388	6.07	301	0.82	DAPENLRPLTADFAEDK-LQVHEK	3229.602	3229.596	-1.940	LQVHEK
DAPENLRPLTADFAEDK	1201.08950	2401.17070	4.75	299	0.89	DAPENLRPLTADFAEDK-QNCQDEK	3499.653	3499.556	-2.029	QNCQDEK
EYATLEECACDDPHACYSYVQK	1434.09496	2867.16265	-7.65	283	0.42	DAPENLRPLTADFAEDK-EYATLEECACDDPHACYSYVQK	5263.213	5263.205	-1.485	DAPENLRPLTADFAEDK
DAPENLRPLTADFAEDK	1201.07912	2401.15896	-8.47	281	0.90	DAPENLRPLTADFAEDK-EYATLEECACDDPHACYSYVQK	5263.213	5263.205	-1.485	EYATLEECACDDPHACYSYVQK
SHCFAVEYK	508.24785	1015.49802	0.27	267	1.44	SHCFAVEYK-EKCFVGRK	2062.959	2062.957	-1.053	EKCFVGRK
EKCFVGRK	525.75951	1050.46394	1.35	267	1.44	SHCFAVEYK-EKCFVGRK	2062.959	2062.957	-1.053	SHCFAVEYK
SHCFAVEYK	508.24637	1015.48546	-2.26	261	1.02	SHCFAVEYK-LQVHEK	1853.928	1853.925	-1.592	LQVHEK
LQVHEK	841.49179	841.49179	2.04	259	0.29	SHCFAVEYK-LQVHEK	1853.928	1853.925	-1.592	SHCFAVEYK
SHCFAVEYK	508.24794	1015.48800	-0.93	259	0.89	SHCFAVEYK-LQVHEK	1853.928	1853.925	-1.592	LQVHEK
MPCTEDYLSILNR	834.41034	1667.81340	0.11	255	2.33	SHCFAVEYK-MPCTEDYLSILNR	2680.284	2680.278	-2.433	SHCFAVEYK
MPCTEDYLSILNR	834.41144	1667.81550	1.43	251	3.28	SHCFAVEYK-MPCTEDYLSILNR	2680.287	2680.278	-3.527	SHCFAVEYK
QNCQDEK	505.21542	1011.42557	3.47	247	0.71	SHCFAVEYK-QNCQDEK	2022.868	2022.885	0.766	SHCFAVEYK

Custom date Fields

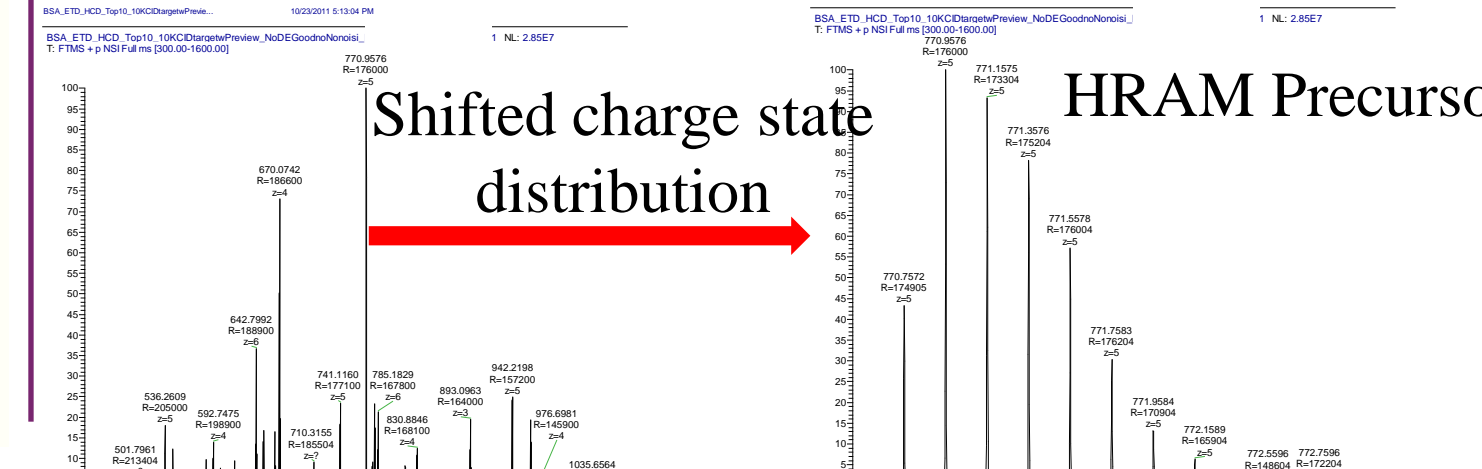
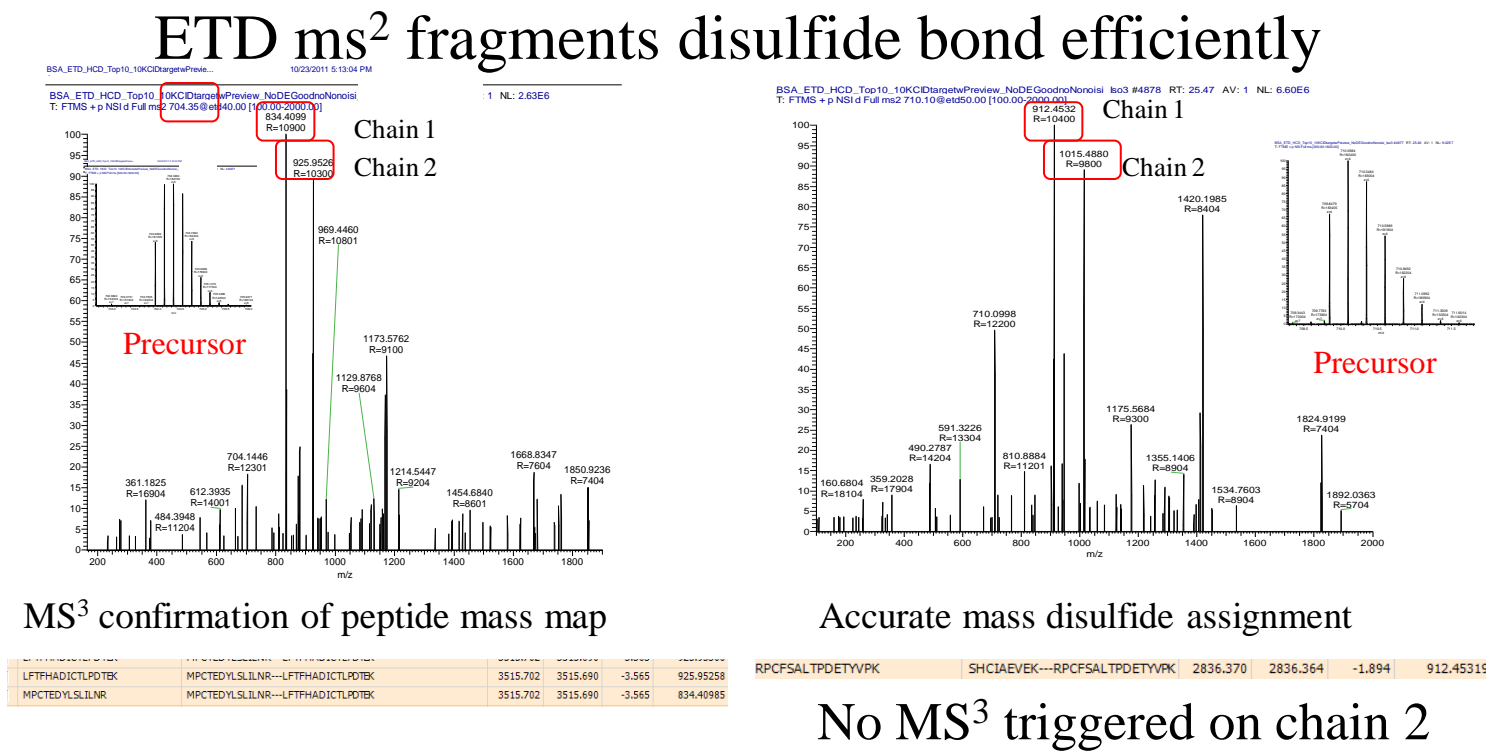


FIGURE 8. Disulfide Map.



Conclusion

- A node development environment is implemented in Proteome Discoverer that allows rapid deployment of custom algorithms without the need to develop graphical user interfaces and other input/output infrastructure. This allows for algorithms to be deployed and shared among the mass spectrometry community
- Multiple nodes were implemented in PD 1.3 and easily integrated into the workflows. For example, the neutral loss filter (fig 5) was integrated with multiple node connections I/O to search under variable parameter settings without affecting any other part of Discoverer
- The external nodes developed for this poster can be used as templates for other development.
- Non-traditional proteomics workflows such as disulfide mapping presented here can also be easily integrated and the data can be displayed via the custom data fields properties of Proteome Discoverer
- Nodes, source code and templates will be available at www.PD-Nodes.org

References

- Wu et al. 2010. Identification of the Unpaired Cysteine Status and Complete Mapping the 17 Disulfides of Recombinant Tissue Plasminogen Activator Using LC-MS with ETD/CID