

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

Iman Mohtashemi, Kai Fritzscheier, Hans Gensemann, Bernard Delanghe, David Horn, and Torsten Ueckert
Thermo Fisher Scientific, San Jose, CA, USA

Overview

Purpose: Here we present 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 Microsoft® C# using Visual Studio® 2010 and their utility was demonstrated by incorporating them into Proteome Discoverer™ software 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 numbers 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 in 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

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 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 software framework, a developer starts with a Microsoft Visual Studio 2010 solution containing a template project that performs a 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 is available for all users. The code on the right side of Figure 1 shows the parameters that are available for users in the Workflow Editor user interface of Proteome Discoverer software (Figure 2) for that node and their acceptable limits. If the user sets any value outside of 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("F1F0PCDE-2274-4129-40D3-1328C1806E41",
    Category = ProcessingNodeCategories.ExternalTest,
    DisplayName = "Neutral Loss Filter",
    Description = "Filters out non-phosphorylated peptides",
    MajorVersion = 0,
    MinorVersion = 4)]
[ProcessingNodeConstraints(UsageConstraint = UsageConstraint.Unrestricted)]
[ProcessingNodeConnectionPoint("SpectrumMSM",
    DataTypeSpecialization = StandardDataTypeSpecializations.AcceptAny,
    ConnectionDirection = ConnectionDirection.Incoming,
    ConnectionMultiplicity = ConnectionMultiplicity.Single,
    ConnectionMode = ConnectionMode.Manual,
    ConnectionRequirement = ConnectionRequirement.RequiredAtDesignTime,
    ConnectionDisplayName = ProcessingNodeCategories.SpectrumMSM.FeatureRetrieval)]
[ProcessingNodeConnectionPoint("SpectrumSource",
    DataTypeSpecialization = StandardDataTypeSpecializations.AcceptAny,
    ConnectionDirection = ConnectionDirection.Outgoing,
    ConnectionMultiplicity = ConnectionMultiplicity.Multiple,
    ConnectionMode = ConnectionMode.Manual,
    ConnectionRequirement = ConnectionRequirement.Optional)]
[LicenseFeature]
Feature = "Discoverer_Base",
Description = "Discoverer base license.",
ShowIfNotAvailable = false]

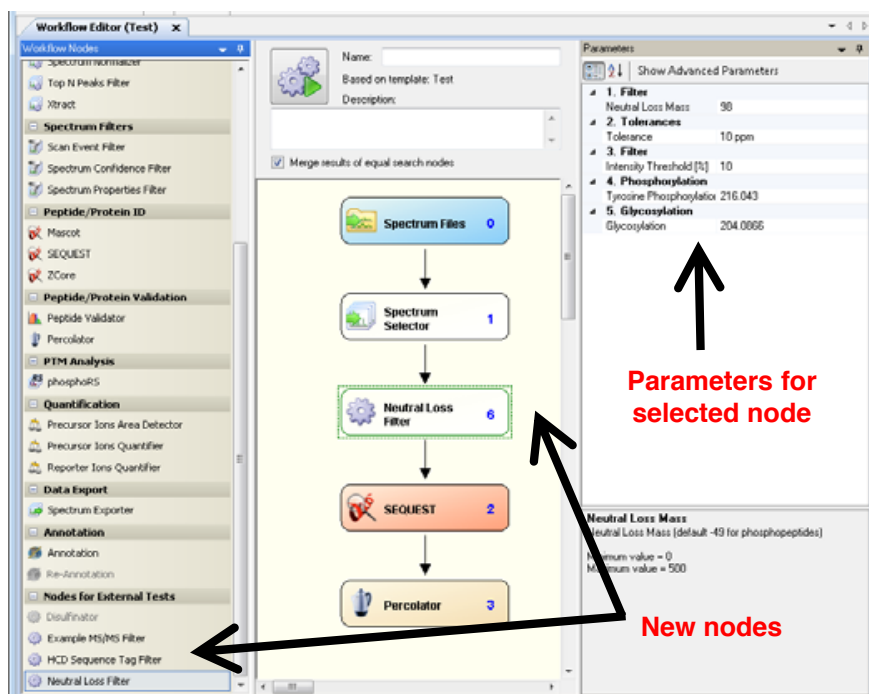
[ParameterInfo]
Category = "1. Filter",
DisplayName = "Neutral Loss Mass",
Description = "Neutral Loss Mass (default -40 for phosphopeptides)",
DefaultValue = "40",
MinimumValue = "0",
MaximumValue = "500",
Position = 1]
public double? Parameter_NeutralLoss;

[ParameterInfo]
Category = "2. Tolerances",
DisplayName = "Tolerance",
Description = "Tolerance",
Subset = "Da|mu|ppm",
DefaultValue = "10.0 ppm", MinimumValue = "0.0001 Da | 0.01 ppm", Max
//IntendedPurpose = ParameterPurpose.DynamicModification,
Position = 1]
public MassToleranceParameter PeptideTolerance;

[ParameterInfo]
Category = "3. Filter",
DisplayName = "Intensity Threshold [N]",
DefaultValue = "10",
MinimumValue = "0",
MaximumValue = "100",
Position = 1]

```

FIGURE 2. Proteome Discoverer software Workflow Editor with the new nodes.



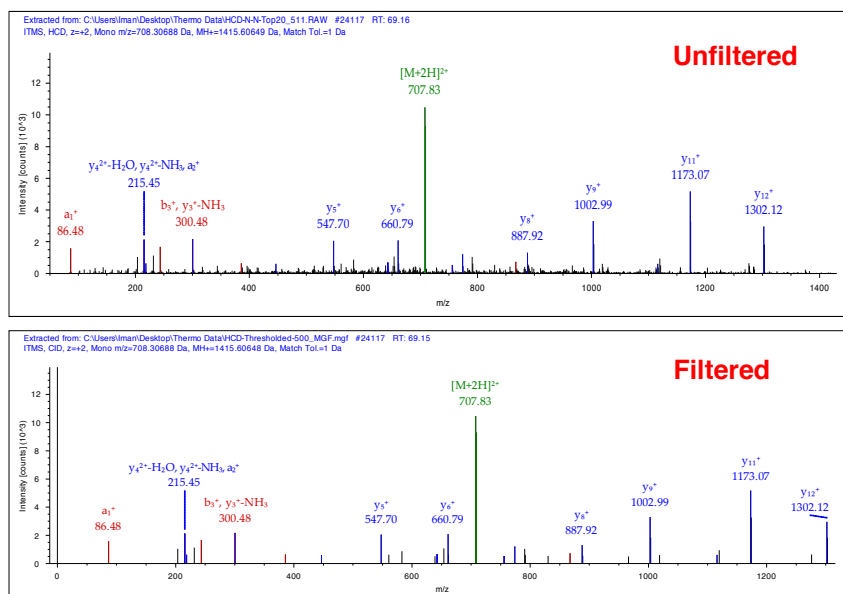
Results

The following four nodes were created for use in the Proteome Discoverer software application: a normalized MS² filter node, an HCD sequence tag filter, a CID MS² neutral loss filter and an ETD triggered MS³ node for disulfide mapping. The use and power of the custom nodes for the analysis of MS data is demonstrated briefly, showing the results for a typical set of data.

Normalized MS² Filter Node

Each MS/MS spectrum was normalized to the most intense peak within the spectrum. A threshold was then applied to achieve maximum identification rate with minimal search input (Figure 3).

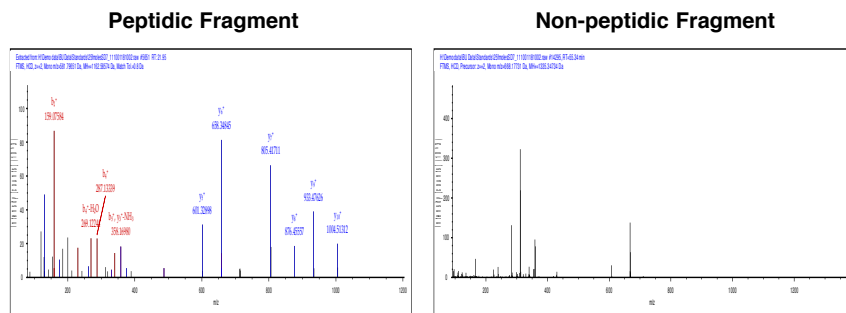
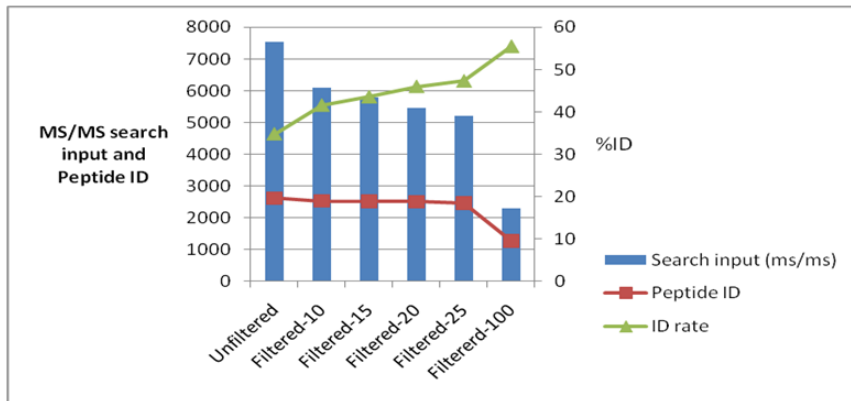
FIGURE 3. Filtered MS² spectrum.



HCD Sequence Tag Filter

A peptide sequence tag filter was applied to HCD MS² data to minimize search input for non-peptidic fragments (Figure 4).

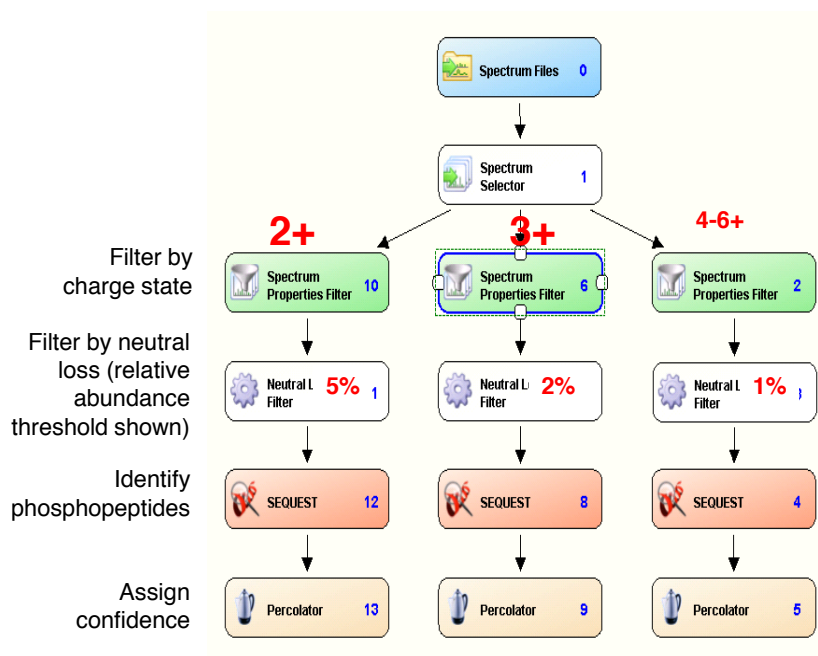
FIGURE 4. HCD sequence tag filter ID versus input



CID MS² Neutral Loss

CID MS² screened for phosphorylation via neutral loss monitoring (Figure 5).

FIGURE 5. Integration of flexible workflow



ETD Triggered MS³ for Disulfide Mapping

Custom data fields were used for data display for disulfide mapping (Figures 6 and 7). Each new field was assigned a new GUID. Figure 8 shows the disulfide map.

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)", ProcessingNodeNumber, searchNodeInfo.ProcessingNodeNumber, CustomDataTarget.Peptide, CustomDataType.String, CustomDataAccessMode.Read, DataVisibility.Visible);
CustomDataField massCrosslinkedPeptide = ProcessingServices.CustomDataService.GetOrCreateCustomDataField(workflow.ID, new Guid("4866912F-E945-4AD1-872C-88D934402AEF"), "Mass Crosslinked Peptide", ProcessingNodeNumber, searchNodeInfo.ProcessingNodeNumber, CustomDataTarget.Peptide, CustomDataType.Double, CustomDataAccessMode.Read, DataVisibility.Visible, plotType: PlotType.Numeric, format: "0.000");
CustomDataField calculatedMassCrosslinkedPeptide = ProcessingServices.CustomDataService.GetOrCreateCustomDataField(workflow.ID, new Guid("87C0A83E-7CBD-4554-A995-C77E3CDF8986"), "Calculated Mass Crosslinked Peptide", ProcessingNodeNumber, searchNodeInfo.ProcessingNodeNumber, 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-B4CD-CE3227E27897"), "Delta Mass Crosslinked Peptide [PPM]", ProcessingNodeNumber, searchNodeInfo.ProcessingNodeNumber, CustomDataTarget.Peptide, CustomDataType.Double, CustomDataAccessMode.Read, DataVisibility.Visible, plotType: PlotType.Numeric, format: "0.000");
    
```

FIGURE 7. Custom data fields for disulfinator node.

Sequence	m/z [Da]	MH+ [Da]	ΔM [ppm]	Onsld	XCorr	Crosslinked Peptide(s)	Mass Onsl	Calculated	Delta Mas	Peptide 1
MPCTEDYLSLILNR	834.41071	1667.81413	0.55	337	2.59	MPCTEDYLSLILNR--EACFAVEGRK	2715.292	2715.283	-3.602	EACFAVEGRK
LFTFHADICTLPDTEK	925.95300	1850.89873	-0.36	335	0.65	MPCTEDYLSLILNR--LFTFHADICTLPDTEK	3515.702	3515.690	-3.565	MPCTEDYLSLILNR
LFTFHADICTLPDTEK	925.95258	1850.89787	-0.82	333	1.88	MPCTEDYLSLILNR--LFTFHADICTLPDTEK	3515.702	3515.690	-3.565	MPCTEDYLSLILNR
MPCTEDYLSLILNR	834.40985	1667.81243	-0.47	333	2.49	MPCTEDYLSLILNR--LFTFHADICTLPDTEK	3515.702	3515.690	-3.565	LFTFHADICTLPD...
MPCTEDYLSLILNR	834.41162	1667.81597	1.65	329	1.06	MPCTEDYLSLILNR--RPFCSALTPDTEYWK	3488.705	3488.690	-4.346	RPFCSALTPDTEY...
MPCTEDYLSLILNR	834.41138	1667.81548	1.36	327	0.81	MPCTEDYLSLILNR--RPFCSALTPDTEYWK	3488.703	3488.690	-3.647	RPFCSALTPDTEY...
RPFCSALTPDTEYVK	912.45618	1823.90508	2.94	327	0.89	MPCTEDYLSLILNR--RPFCSALTPDTEYWK	3488.703	3488.690	-3.647	MPCTEDYLSLILNR
MPCTEDYLSLILNR	834.41022	1667.81316	-0.03	323	3.27	MPCTEDYLSLILNR--LQVLHFK	2506.256	2506.250	-2.412	LQVLHFK
MPCTEDYLSLILNR	834.41205	1667.81682	2.16	315	2.19	QNCDQFEK--MPCTEDYLSLILNR	2676.219	2676.210	-3.431	QNCDQFEK
MPCTEDYLSLILNR	834.41046	1667.81365	0.26	311	2.07	QNCDQFEK--MPCTEDYLSLILNR	2676.217	2676.210	-2.427	QNCDQFEK
DAIPENLPRLTADFAEDKVK	1201.09058	2401.17388	6.07	305	0.81	DAIPENLPRLTADFAEDKVK--LFTFHADICTLPDTEK	4249.044	4249.036	-1.932	LFTFHADICTLPD...
DAIPENLPRLTADFAEDKVK	1201.09058	2401.17388	6.07	301	0.82	DAIPENLPRLTADFAEDKVK--LQVLHFK	3239.602	3239.596	-1.940	LQVLHFK
DAIPENLPRLTADFAEDKVK	1201.08899	2401.17070	4.75	299	0.69	DAIPENLPRLTADFAEDKVK--QNCDQFEK	3409.563	3409.556	-2.029	QNCDQFEK
EYEATLECCAKDDPHACSTVFK	1494.08496	2867.16265	-7.65	283	0.42	DAIPENLPRLTADFAEDKVK--EYEATLECCAKDDPHACSTVFK	5263.313	5263.305	-1.485	DAIPENLPRLTAD...
DAIPENLPRLTADFAEDKVK	1201.07312	2401.13896	-8.47	281	0.80	DAIPENLPRLTADFAEDKVK--EYEATLECCAKDDPHACSTVFK	5263.313	5263.305	-1.485	EYEATLECCAK...
SHCIAEVEK	508.24765	1015.48802	0.27	267	1.44	SHCIAEVEK--EACFAVEGRK	2062.959	2062.957	-1.053	EACFAVEGRK
EACFAVEGRK	525.75061	1050.49394	1.35	267	1.44	SHCIAEVEK--EACFAVEGRK	2062.959	2062.957	-1.053	SHCIAEVEK
RPFCSALTPDTEYVK	912.45319	1823.89910	-0.34	263	1.21	SHCIAEVEK--RPFCSALTPDTEYWK	2836.370	2836.364	-1.894	SHCIAEVEK
SHCIAEVEK	508.24637	1015.48546	-2.26	261	1.03	SHCIAEVEK--LQVLHFK	1853.928	1853.925	-1.582	LQVLHFK
LQVLHFK	841.46179	841.46179	2.04	259	1.08	SHCIAEVEK--LQVLHFK	1853.928	1853.925	-1.582	SHCIAEVEK
SHCIAEVEK	508.24704	1015.48680	-0.93	259	0.88	SHCIAEVEK--LQVLHFK	1853.928	1853.925	-1.582	LQVLHFK
MPCTEDYLSLILNR	834.41034	1667.81340	0.11	255	2.53	SHCIAEVEK--MPCTEDYLSLILNR	2680.284	2680.278	-2.433	SHCIAEVEK
MPCTEDYLSLILNR	834.41144	1667.81560	1.43	251	3.22	SHCIAEVEK--MPCTEDYLSLILNR	2680.287	2680.278	-3.527	SHCIAEVEK
QNCDQFEK	506.21542	1011.42357	3.47	247	0.71	SHCIAEVEK--QNCDQFEK	2023.886	2023.885	-0.766	SHCIAEVEK

Custom data fields →

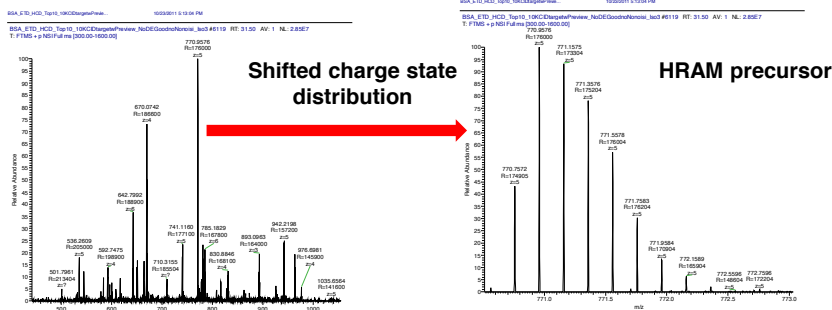
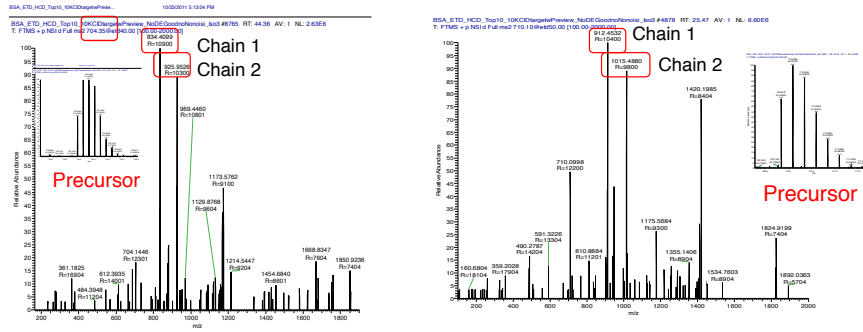


FIGURE 8. Disulfide map.

ETD MS² fragments disulfide bond efficiently

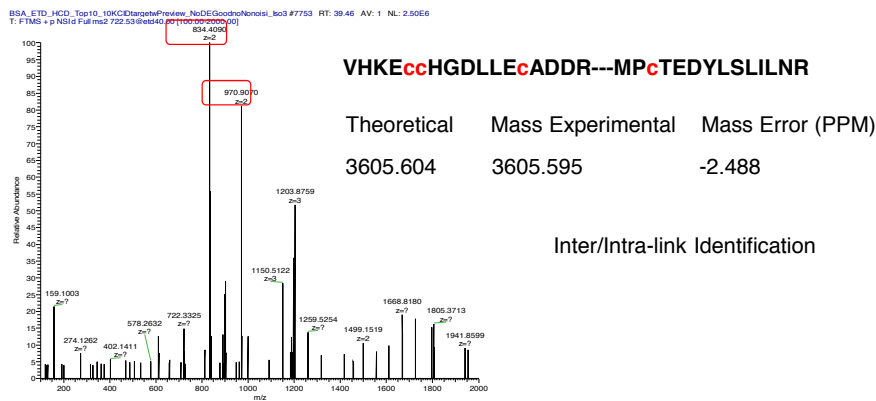


MS³ confirmation of peptide mass map

Accurate mass disulfide assignment

LPFSAALTCDEYK	SHCAEVEK--RPFSAALTCDEYK	2836.275	2836.284	-0.884	812.4531
MPCTEDYLSLILNR	SHCAEVEK--RPFSAALTCDEYK	3605.604	3605.595	-0.884	812.4531

No MS³ triggered on chain 2



Conclusion

- A node development environment is implemented in Proteome Discoverer software 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 Proteome Discoverer software version 1.3 and easily integrated into the workflows. For example, the neutral loss filter (Figure 5) was integrated with multiple node connections input/output to independently search under variable parameter settings.
- The external nodes developed for this poster can be used as templates for other node development efforts.
- Non-traditional proteomics workflows, such as disulfide mapping presented here, can also be easily integrated. The data can be displayed via the custom data fields properties of the Proteome Discoverer software graphical user interface.
- Nodes, source code and templates will be available at www.PD-Nodes.org

Microsoft, Windows, and Visual Studio are registered trademarks of Microsoft Corporation. All other trademarks are the property of Thermo Fisher Scientific and its subsidiaries.

This information is not intended to encourage use of these products in any manners that might infringe the intellectual property rights of others.

www.thermoscientific.com

©2012 Thermo Fisher Scientific Inc. All rights reserved. ISO is a trademark of the International Standards Organization. All other trademarks are the property of Thermo Fisher Scientific Inc. and its subsidiaries. This information is presented as an example of the capabilities of Thermo Fisher Scientific Inc. products. It is not intended to encourage use of these products in any manners that might infringe the intellectual property rights of others. Specifications, terms and pricing are subject to change. Not all products are available in all countries. Please consult your local sales representative for details.



Africa-Other +27 11 570 1840	Europe-Other +43 1 333 50 34 0	Japan +81 45 453 9100	Spain +34 914 845 965
Australia +61 3 9757 4300	Finland/Norway/Sweden +46 8 556 468 00	Latin America +1 561 688 8700	Switzerland +41 61 716 77 00
Austria +43 1 333 50 34 0	France +33 1 60 92 48 00	Middle East +43 1 333 50 34 0	UK +44 1442 233555
Belgium +32 53 73 42 41	Germany +49 6103 408 1014	Netherlands +31 76 579 55 55	USA +1 800 532 4752
Canada +1 800 530 8447	India +91 22 6742 9434	New Zealand +64 9 980 6700	
China +86 10 8419 3588	Italy +39 02 950 591	Russia/CIS +43 1 333 50 34 0	
Denmark +45 70 23 62 60		South Africa +27 11 570 1840	

PN63563_E 06/12S

Thermo
SCIENTIFIC
Part of Thermo Fisher Scientific