

Comprehensive Analysis of a Published Dataset: Identifying more of your data using the latest Protein Prospector

Peter R. Baker, Robert J. Chalkley, Katalin F. Medzihradzky and Alma L. Burlingame

Mass Spectrometry Facility, Dept. of Pharmaceutical Chemistry, University of California, San Francisco, USA

Summary

A public dataset (Klimek et al 2007):

- File: QS20060131_S_18mix_02.wiff
- Mixture of roughly 30 proteins.
- Tryptic digest (microwave-assisted)
- 1867 spectra

The aim was to use the tools in the latest version (v5.0) of Protein Prospector (**freely available on the web at <http://prospector2.ucsf.edu>**) to try to identify as many as possible of the spectra.

Experimental

Database searching is performed as a 2 step process.

- 1) First the proteins present are identified using very strict criteria.
- 2) Data are searched against the list of identified proteins using:
 - A selection of common variable modifications.
 - Mass modification searching (Baker et al 2006)
 - No enzyme specificity.

Semi-manual analysis was then performed using various options and programs within Protein Prospector:

- The precursor data was inspected for each spectrum to determine the correct precursor m/z and charge.
- Each spectrum was also individually searched using MS-Tag from File to allow checking for multiple hits and adjustment of the precursor m/z and charge.
- A separate search was done to look for up to 2 amino acid substitutions using the homology tick box. Any substitutions identified were checked by searching NCBItr and UniprotKB using MS-Pattern.

Batch-Tag Options in PPv5

The screenshot displays the 'Batch-Tag Web' interface. Key elements include:

- Database:** SwissProt 2007 12 04
- Digest:** Trypsin
- Max. Missed Cleavages:** 1
- Species:** All
- Variable Mods:** Acetyl (O), Aspy (Protein N-term), Aspy+Oxidation (Protein N-term)
- Mass Modifications:** Range (Da) 100 to 300, Defect 0.00048
- Matrix Modifications:** Unknown Amino Acid, Single Base Change, Homology
- Upload Data From File:** Includes a 'Browse...' button for file selection.

Non-specific options (No enzyme also available on Digest menu)

Search can be limited by accession number

Mass modification options can be combined with variable modifications

Neutral loss option

Homology option

Wiff file can be uploaded and will be automatically centroided. The raw data precursor and MSMS data can thus be viewed.

Search Compare Options in PPv5

Output format can be HTML or tab delimited text

Score filters

Composition filters

Report Columns

Raw data/quantitation options

Option that allows averaging of the precursor data

If this option is selected all subsequent MSMS spectra will be displayed using all the peaks in the centroid file

Links from Search Compare Peptide Report

720.3563⁺³

Precursor scan can be accessed here

IFDGVNSAFHLWC(210.1008)NGR⁺³

Search Compare report sorted by start amino acid shows several cysteine modifications

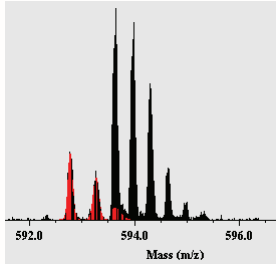
Spectra can be searched individually with different parameters

MS-Tag from File

MS-Product

The peptide match can be viewed in MS-Product. Users have control of various aspects of how the match is displayed

Contributions from Neighboring Peaks

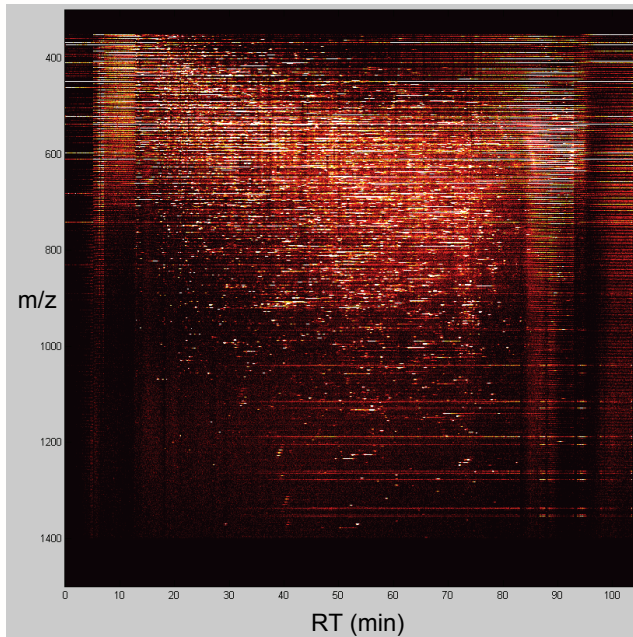
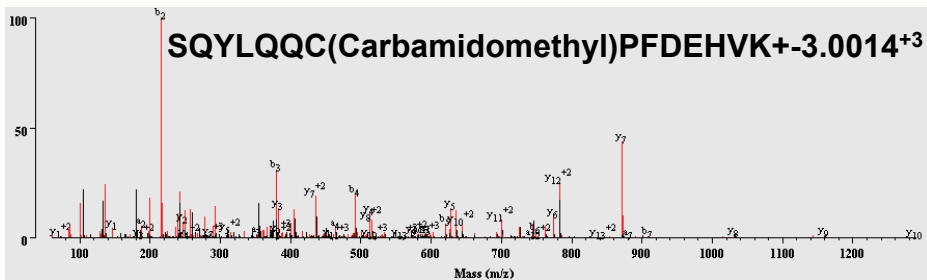


• A lot of the MSMS spectra contain more than one component. These are either from a closely co-eluting peptide or chemical background peaks.

• 60 spectra had two identifiable peptides.

• The centroiding script believes the precursor is doubly charged peak (in red on the left), but most of the fragmentation is from the +3 peak.

• Prospector will find the +3 peak via the neutral loss option in mass modification searching. Around 140 similar extra hits were found via the neutral loss loss option representing cases where the major fragmentation is not at the m/z recorded in the centroid file.



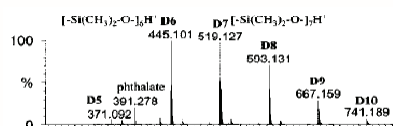
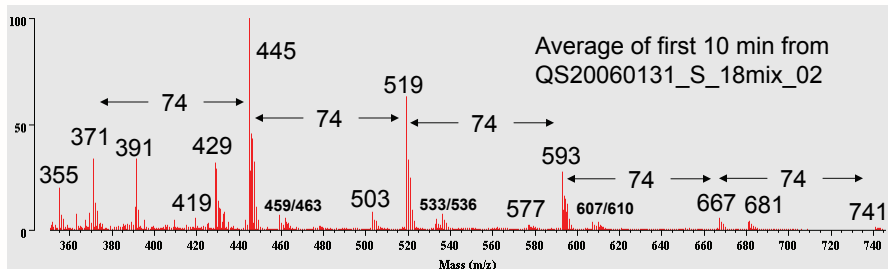
Background Ions

A 3D plot of the raw data clearly shows that there are no peptides in the first 12 min and the last 19 min.

Background ions can be seen as horizontal lines on this plot, the most intense being at m/z 445.

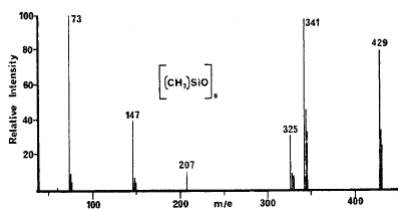
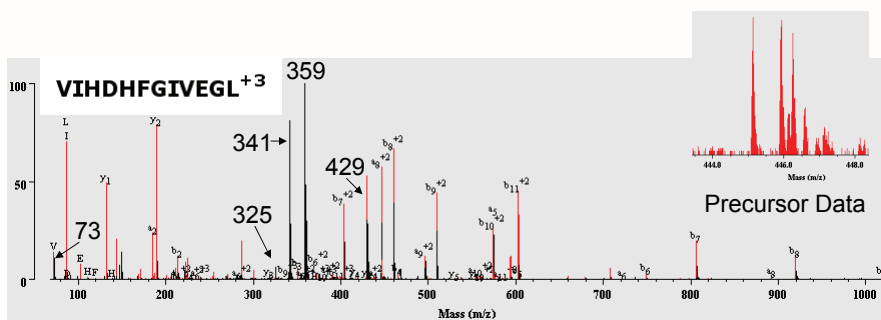
Protein Prospector can produce these plots in collaboration with MATLAB.

Background Ions



- 355, 371 – D5 decamethylcyclopentasiloxane
- 429, 445, 463 – D6 dodecamethylcyclohexasiloxane
- 503, 519, 536 – D7 tetradecamethylcycloheptasiloxane
- 577, 593, 610 – D8 hexadecamethylcyclooctasiloxane
- 651, 667 – D9 octadecamethylcyclononasiloxane
- 725, 741 – D10 eicosamethylcyclodecasiloxane
- Lower peak is loss of methane (-16 Da)
- Higher peak is an ammonium adduct (+17 Da)
- 391 – Dioctyl phthalate

MSMS Spectrum Contaminated with D6



There was some background component in at least 84 spectra.

Is it worth filtering out these masses from the peaklist?

Non-Specific and Missed Cleavages

• Prospector now has a comprehensive range of options for searching for enzyme non-specific cleavage peptides.

- No enzyme
- Non-Specific at N termini
- Non-Specific at C termini } **Unique to Protein Prospector**
- Non-Specific at 1 termini (N or C but not both)
- Non-Specific at 2 termini (like No Enzyme but considers missed cleavages)

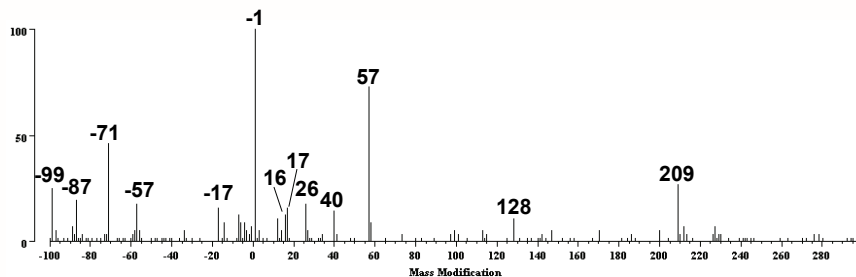
• This sample had a larger than usual number of peptide with non-specific cleavages but very few with missed cleavages.

- Product of microwave-assisted digestion.

Type	1 st Pep	2 nd Pep	Total	%
0	669	43	712	57
N	246	8	254	20.4
C	220	7	227	18.2
2	53	2	55	4.4

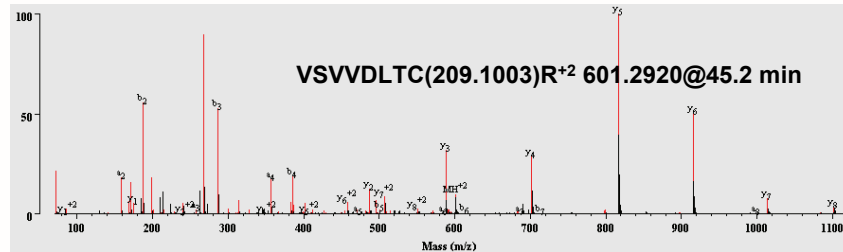
Sometimes a second co-eluting peptide could be identified in a single MSMS spectrum.

Mass Modifications

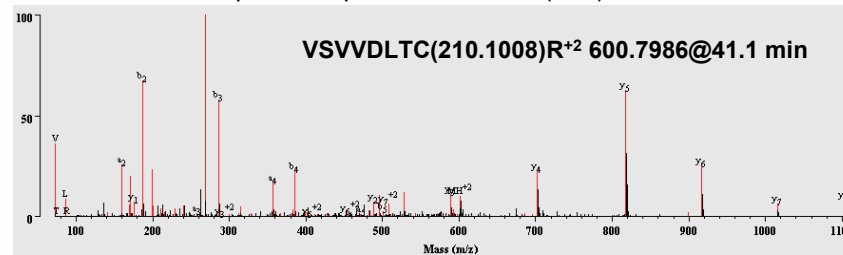


- Prospector can display a histogram showing the mass modifications found.
- Results can be sorted by mass modification.
 - If a modification is observed several times it is more likely to be real.
 - Helps identify novel modifications such as C(209).
 - This modification is formed by the alkylation of the Cys-SH with dithiothreitol, then alkylation of DTT on its other -SH group by iodoacetamide (Chalkley et al 2008).

Unexpected Modifications - C(209) and C(210)



An example of a spectra with the C(209) modification.

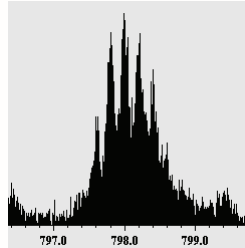


C(210) was assigned as Cys(Carboxymethyl-DTT)-derivative

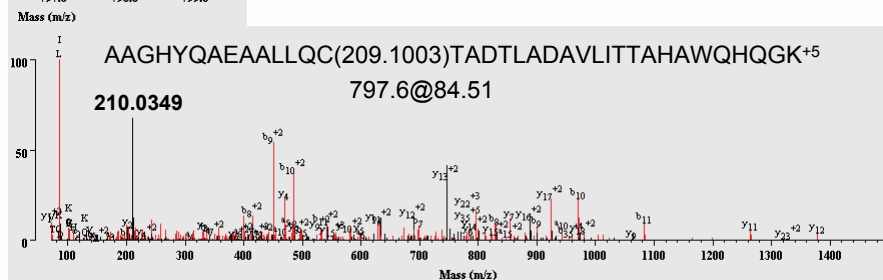
List of Modifications Found

- C(Carbamidomethyl), C(Dehydro)
- C(209), C(210) – see above
- C(12) – tentatively assigned as carbon incorporation similar to P(12)
- C(32) – tentatively assigned as persulfide
- C(56) – maybe Acrolein 56
- -17Da - ring formation from C(Carbamidomethyl) at N-terminus
- E(128) – possible insertion of Lys residue (see below)
- H(Methyl) - Tele-methylhistidine in Alpha skeletal muscle actin - listed in Swiss Prot
- K(Methyl), K(Dimethyl) - In trypsin
- K(Delta:H(6)C(6)O(1)) – maybe Acrolein 94
- M(Met-loss+Acetyl), M (Met-loss) – both from N-terminal processing
- M (Oxidation)
- N(Asn->Succinimide), Asn(Deamidated)
- P(-30) - Maybe Pro->Pyrrolidinone
- P(12) – See Nielsen 2006
- Q(Gln->pyro-Glu)
- +17 Neutral loss - Ammonium adduct
- +32 Neutral loss
- +41 Neutral loss – Acetonitrile adduct
- W(12) – tentatively assigned as carbon incorporation similar to P(12)
- W(Oxidation)

Incorrect Charge /Monoisotopic Mass

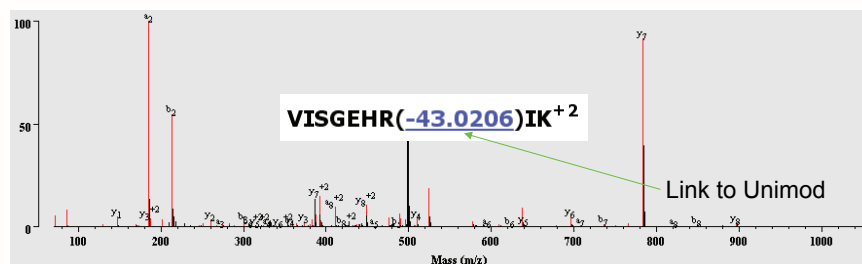


This 5+ peak had the charge and monoisotopic mass incorrectly assigned in the peak list. The charge and precursor m/z were determined by examining the precursor scan and the hit was found by using MS-Tag from File, manually specifying the charge.



Homologous Peptides

- Several spectra are of homologs of peptides found in the proteins identified in this dataset:
- Initially found by either:
 - mass modification search
 - homology tick box (amino acid substitution) option.
- Can be checked by searching a large database (NCBI/UniprotKB) using MS-Pattern or MS-Homology.



Homologous Peptides

PSI-MS Name	Interim name	Description	Monoisotopic mass	Average mass	Composition
	Asp->Ala	Asp->Ala substitution	-43.989829	-44.0095	C(-1) O(-2)
Arg->GluSA	Argglutamicealde	Arginine oxidation to glutamic semialdehyde	-43.053433	-43.0711	H(-5) C(-1) N(-3) O
	Arg->Ile	Arg->Ile or Arg->Leu substitution	-43.017047	-43.0280	H(-1) N(-3)
	Val->Gly	Val->Gly substitution	-42.046950	-42.0797	H(-6) C(-3)
Arg->Orn	Arg2Orn	Ornithine from Arginine	-42.021798	-42.0400	H(-2) C(-1) N(-2)

Hits from Unimod for -43 Da

MS-Pattern

Database: [UnProt20071010] DNA_Frame Translation: [3] Species: [All] Output Type: [HTML] Hits to file: [] Name: [lastes] Digest: [No enzyme]

Pre-Search Parameters

Start Search

Sample ID (comment): []

Reg Expression (Use CAPS): [VISGEHLIK] Maximum Reported Hits: [200] Pre Search Only: []

Max. # of Mismatched AA's: [0]

MS-Pattern Search Results

Search completed. 11 sec elapsed. 0 sec remaining.

[+] Parameters

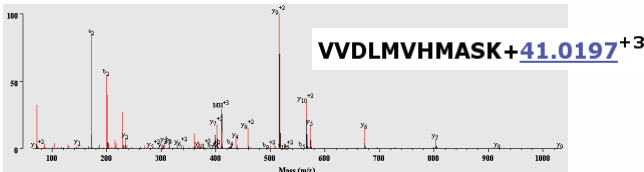
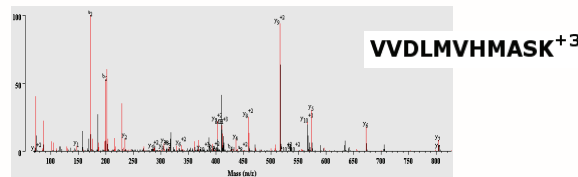
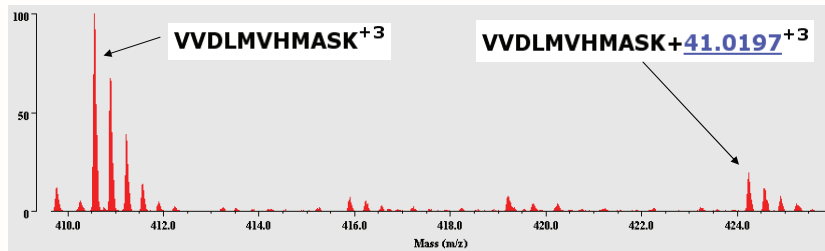
[+] Pre Search Results

MS-Pattern search selects 1 entry.

Matching Sequence Peptide M+H MS-Digest Index # Protein MW (Da)/pI Accession # Species Protein Name
 (R)VISGEHLIK(A) 995.5884 3439688 58493/6.3 Q208A7 BACLI Thermotolerant alpha-amylose

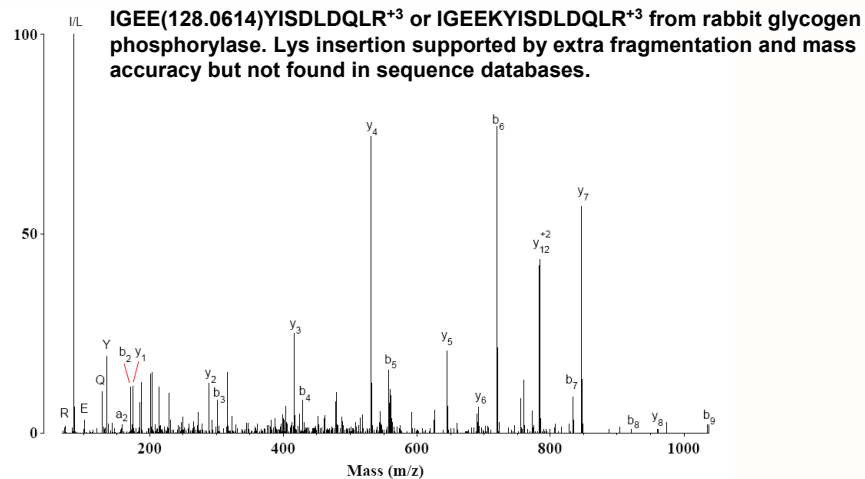
MS-Pattern identifies VISGEHLIK from Q208A7 from UniprotKB as a possible hit.

A Neutral Loss - Acetonitrile



Inserted Amino Acid

The dataset contains one example of what is thought to be an inserted amino acid. This was found by the mass modification search.



IGEEKYISDL DQLR is a better match than IGEEKEYISDL DQLR

Conclusions

- Programs in Protein Prospector allowed interpretation of many more (the majority) of the spectra:
- Regular database search identifies about 20% of spectra.
- Allowing for defined variable modifications and non-tryptic cleavages, Sequest identified about 25% of spectra (Klimek et al.)
- Improved Protein Prospector (with mass modification searching) along with manual inspection of the matches identified 64% of the spectra.
- If the spectra at the start and end of the run, where there are no peptides, are excluded 77% are interpreted

This software is freely available for use on the web at:

<http://prospector2.ucsf.edu>

Acknowledgements

Funding: NIH NCRG grant 001614 and the Vincent Coates Foundation.

References

- 1) Klimek, J., Eddes, J. S., Hohmann, L., Jackson, J., Peterson, A., Letarte, S., Gafken, P. R., Katz, J. E., Mallick, P., Lee, H., Schmidt, A., Ossola, R., Eng, J. K., Aebersold, R., and Martin, D. B., The Standard Protein Mix Database: A Diverse Data Set To Assist in the Production of Improved Peptide and Protein Identification Software Tools. *J Proteome Res*, 2007.
- 2) Baker, P. R., Medzihradzky, K. F., and Burlingame, A. L. Improved Methods for Comprehensive Sample Analysis Using Protein Prospector. *Proceedings of the 54th ASMS Conference*. 2006. Seattle, WA.
- 3) Schlosser, A. and Volkmer-Engert, R. Volatile polydimethylcyclsiloxanes in the ambient laboratory air identified as source of extreme background signals in nanoelectrospray mass spectrometry. *J. Mass Spectrom.*, **38**: 523-525 (2003)
- 4) Pickering, G. R., Olliff, C. J. and Rutt, K. J. The mass spectrometric behaviour of dimethylcyclsiloxanes. *Organic Mass Spectrometry.*, **10**: 1035-1045 (1975)
- 5) Chalkley, R. J., Baker, P. R., Medzihradzky, K. F., Lynn, A. J. and Burlingame, A. L. In-depth Analysis of Tandem Mass Spectrometry Data from Disparate Instrument Types. *Submitted* (2008)
- 6) Baker P. R., Chalkley R. J., Medzihradzky K. F. and Burlingame A. L. Discovery of Unanticipated Modifications using Protein Prospector, *Proceedings of the 55th ASMS Conference*. 2007 Indianapolis, IA
- 7) Nielsen, M. L., Savitski, M. M. and Zubarev, R. A. Discovery and Unexpected Modifications by Modificomb: the Case of Glu-Methylation in Human Cells. *Proceedings of the 54th ASMS Conference*. 2006. Seattle, WA.