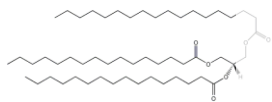
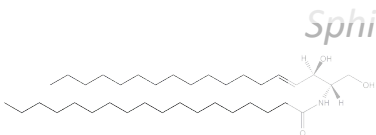


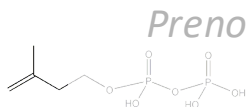
Glycerophospholipids



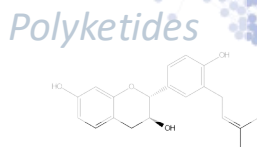
Glycerolipids



Sphingolipids



Prenol lipids



Polyketides

Sterol lipids



Computational Tools for Tandem Mass Spectrometry

Scott Walmsley, Ph.D.

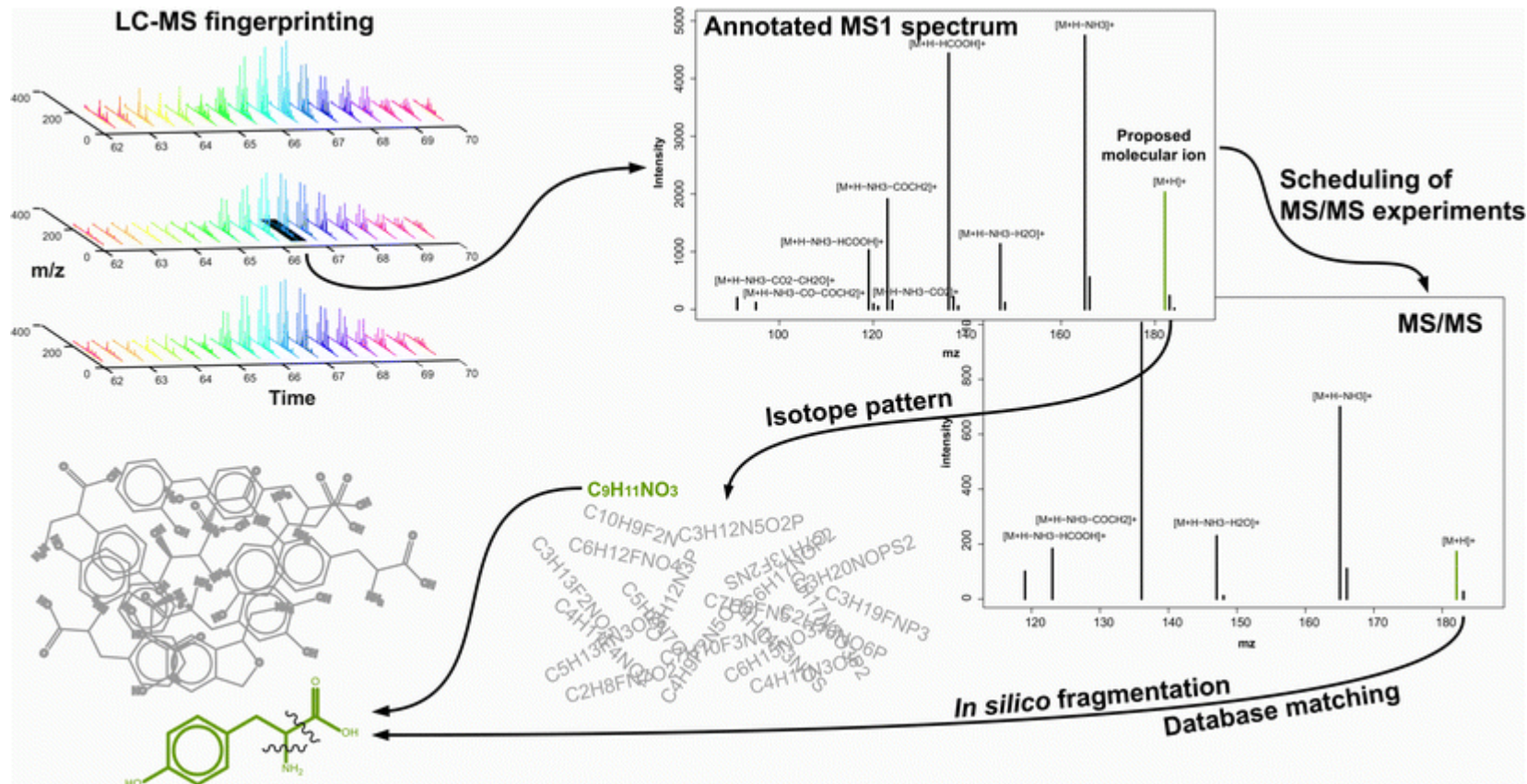
Mass Spectrometry Facility

DOPS-Skaggs SOP

UC Anschutz

Tandem Mass Spectrometry

- Peptide MS/MS
 - Proteomics
- Small Molecule MS/MS
 - Metabolomics, environmental, toxicology, etc...



“Shotgun” proteomics

1: Meissner F, Mann M. Quantitative shotgun proteomics: considerations for a high-quality workflow in immunology. *Nat Immunol.* 2014 Feb;15(2):112-7. doi: 10.1038/ni.2781. PubMed PMID: 24448568.

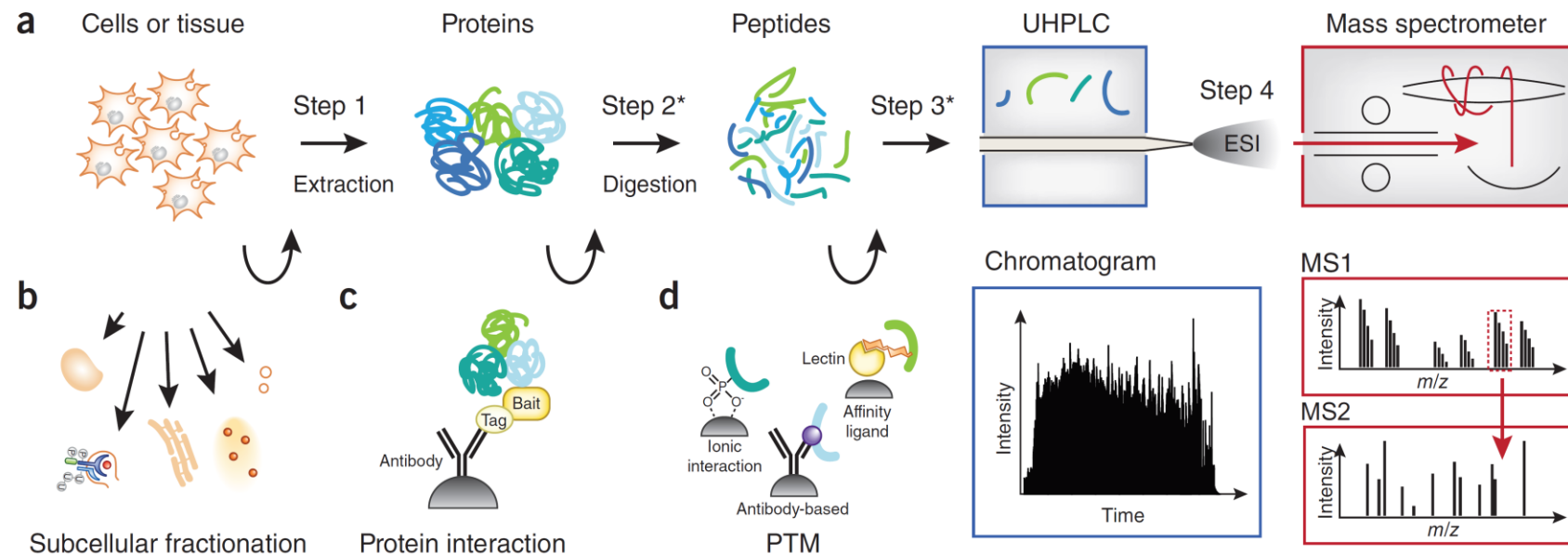
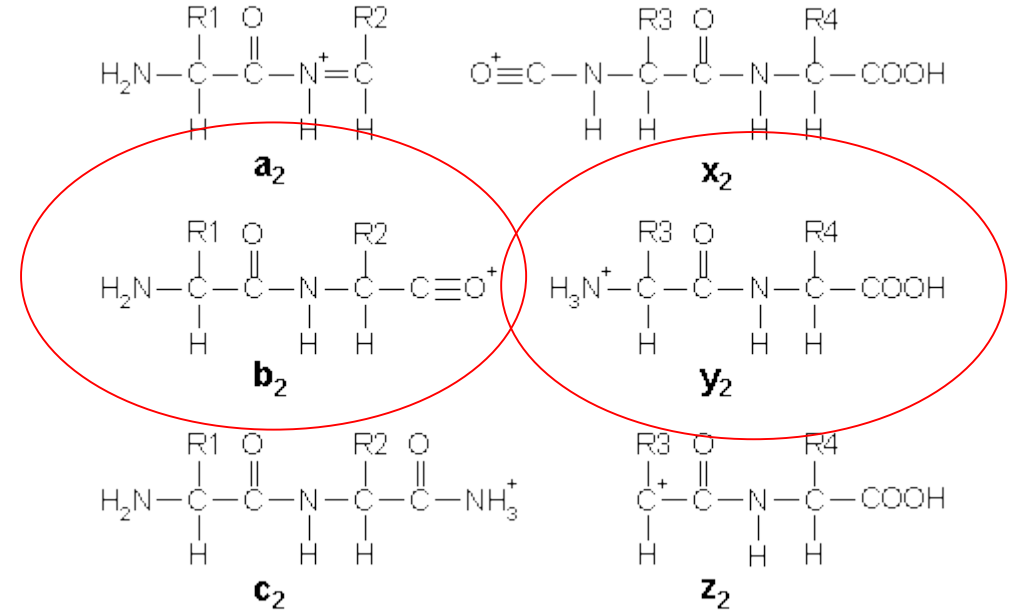
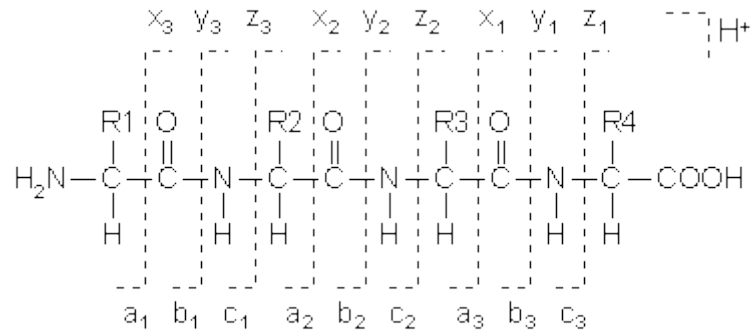
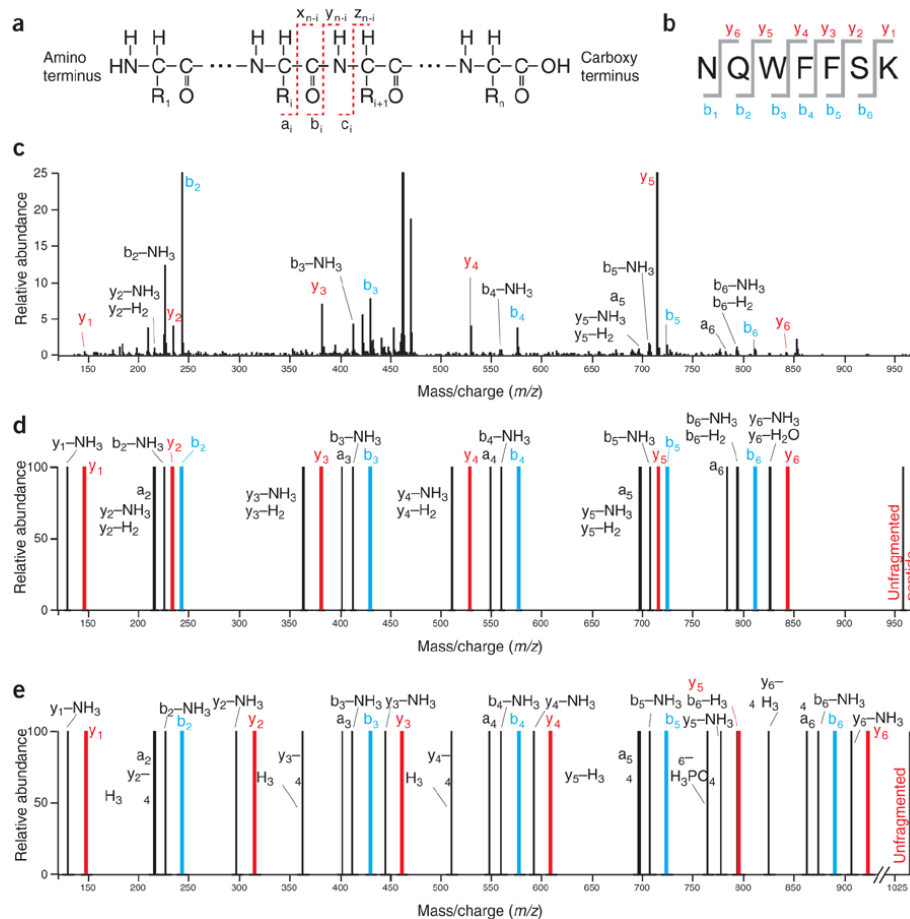


Figure 2 Shotgun proteomics workflow. (a) The generic workflow of modern LC-MS-based proteomics consists of four steps. In step 1, proteins are extracted from tissues, body fluids, cells or subcellular compartments. In step 2, proteins are proteolytically digested. In step 3, peptides are separated by UHPLC. In step 4, peptides are ionized by electrospray (electrospray ionization, ESI), and their masses and fragment masses are acquired in a mass spectrometer. In the workflow described here, LC-MS is performed with an UHPLC system coupled online to a Q Exactive (Thermo Fisher Scientific). (b–d) Variations of the workflow that include enrichment steps for proteins in subcellular compartments (b), for interacting proteins (c) and for peptides with PTMs (d). At steps 2 and 3 (*), additional fractionation of proteins or peptides is possible.

“Tryptic” peptides can be predictably fragmented in a Mass Spectrometer



Fragmentation by MSMS



Expected fragmentation patterns are matched to an experimental spectra.

WE call this the peptide spectrum match (PSM)

This is performed using a “search engine”

Search engines rank order the “best hit” followed by the next best matches

Tandem Mass Spectra

- 1000's to millions produced each file
- Peptide-spectra match (PSM) accomplished via “search engines”
- Examples of search engines that perform PSMs:
 - Theoretical spectral matching:
 - SEQUEST
 - MASCOT
 - X!Tandem
 - Experimental Library Matching:
 - NIST MS
 - SpectraST

Search engines

→ Match MSMS spectra to protein sequence

- X!TANDEM

<https://www.ncbi.nlm.nih.gov/pubmed/14976030>

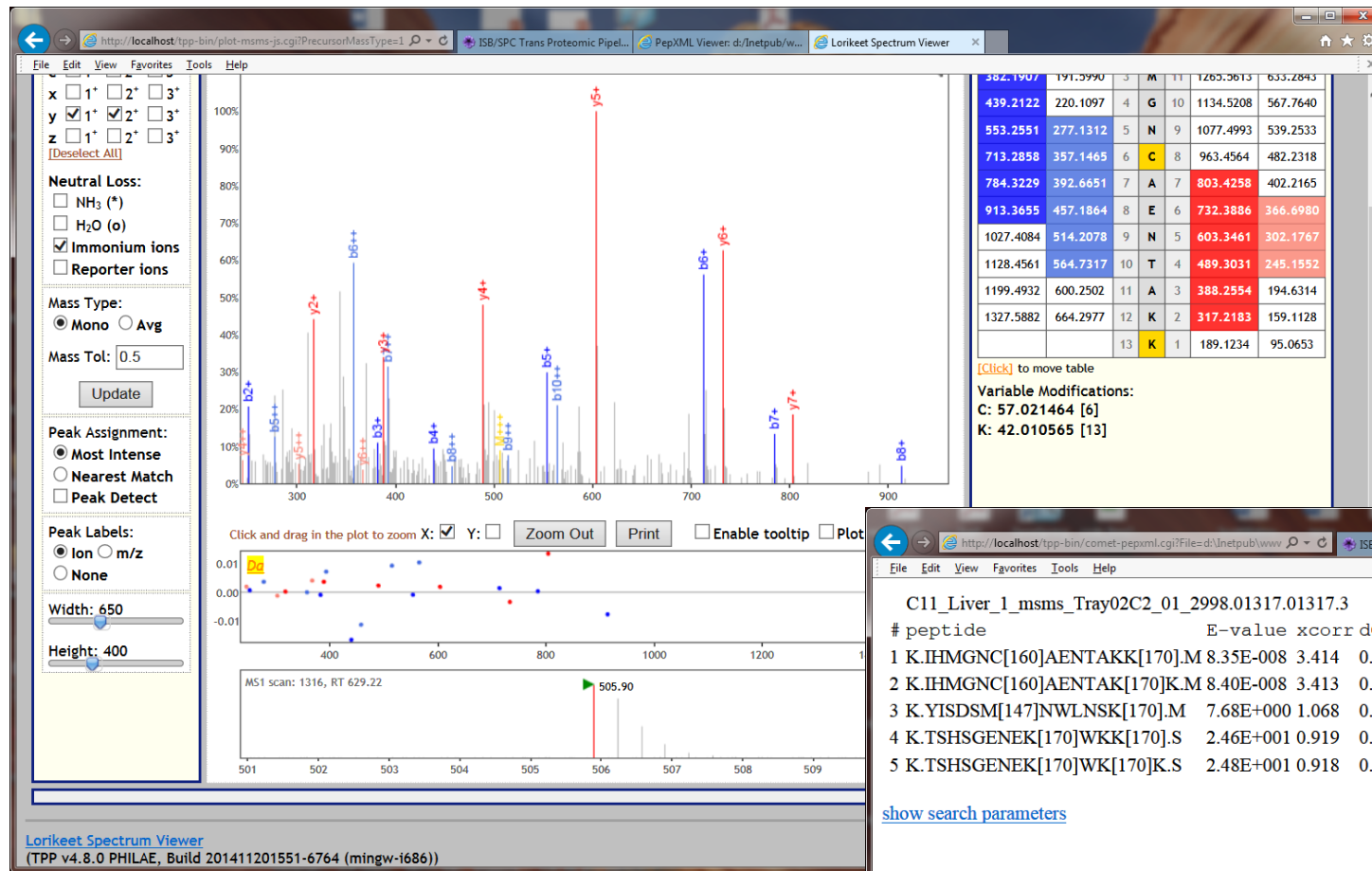
- COMET (SEQUEST)

<https://link.springer.com/article/10.1007%2Fs13361-015-1179-x>

- MASCOT (Mowse)

<https://www.ncbi.nlm.nih.gov/pubmed/15335725?dopt=AbstractPlus>

MSMS Search Result (Comet)



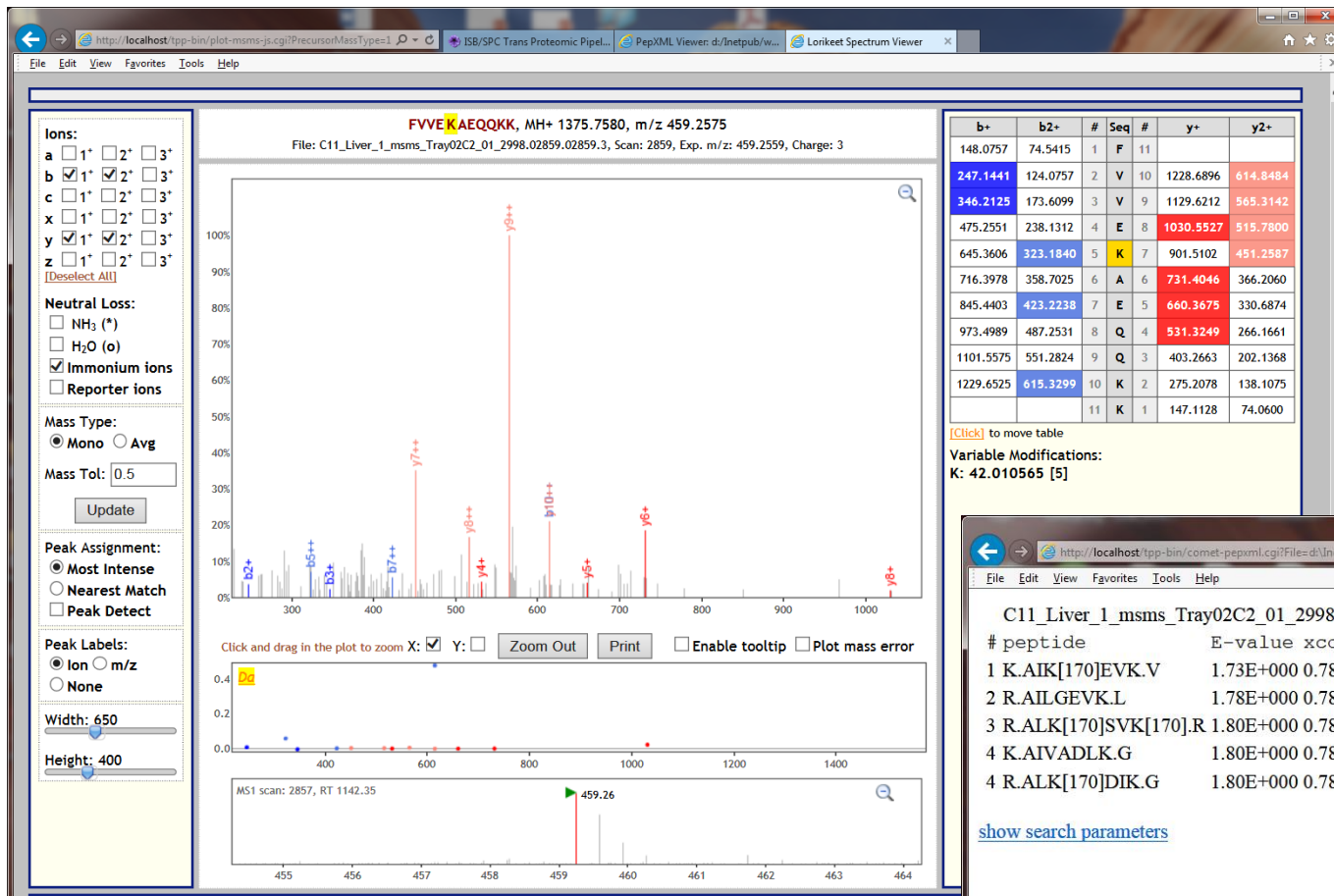
PSM for a good match

C11_Liver_1_msms_Tray02C2_01_2998.01317.01317.3											mass 1514.690054	
#	peptide	E-value	xcorr	dCn	Sp	rankSp	neutral	mass	ions	protein		
1	K.IHMGNC[160]AENTAKK[170].M	8.35E-008	3.414	0.687	537.0	1	1514.696983	17/48	sp Q8QZT1 THIL_MOUSE	+2		
2	K.IHMGNC[160]AENTAK[170].K	8.40E-008	3.413	0.687	537.0	1	1514.696983	17/48	sp Q8QZT1 THIL_MOUSE	+2		
3	K.YISDSM[147]NWLNSK[170].M	7.68E+000	1.068	0.731	29.3	2	1514.671145	5/44	sp P48722 HS74L_MOUSE	+3		
4	K.TSHSGENEK[170]WKK[170].S	2.46E+001	0.919	0.732	23.1	5	1513.716122	4/44	tr Q68FG1 Q68FG1_MOUSE	+3		
5	K.TSHSGENEK[170]WK[170].K	2.48E+001	0.918	0.732	23.1	5	1513.716122	4/44	tr Q68FG1 Q68FG1_MOUSE	+3		

[show search parameters](#)

MSMS Search Result #2 (COMET)

PSM for a bad match



Browser window: <http://localhost:7777/comet-pepxml.cgi?File=d:\inetpub\www\comet-pep.xml>

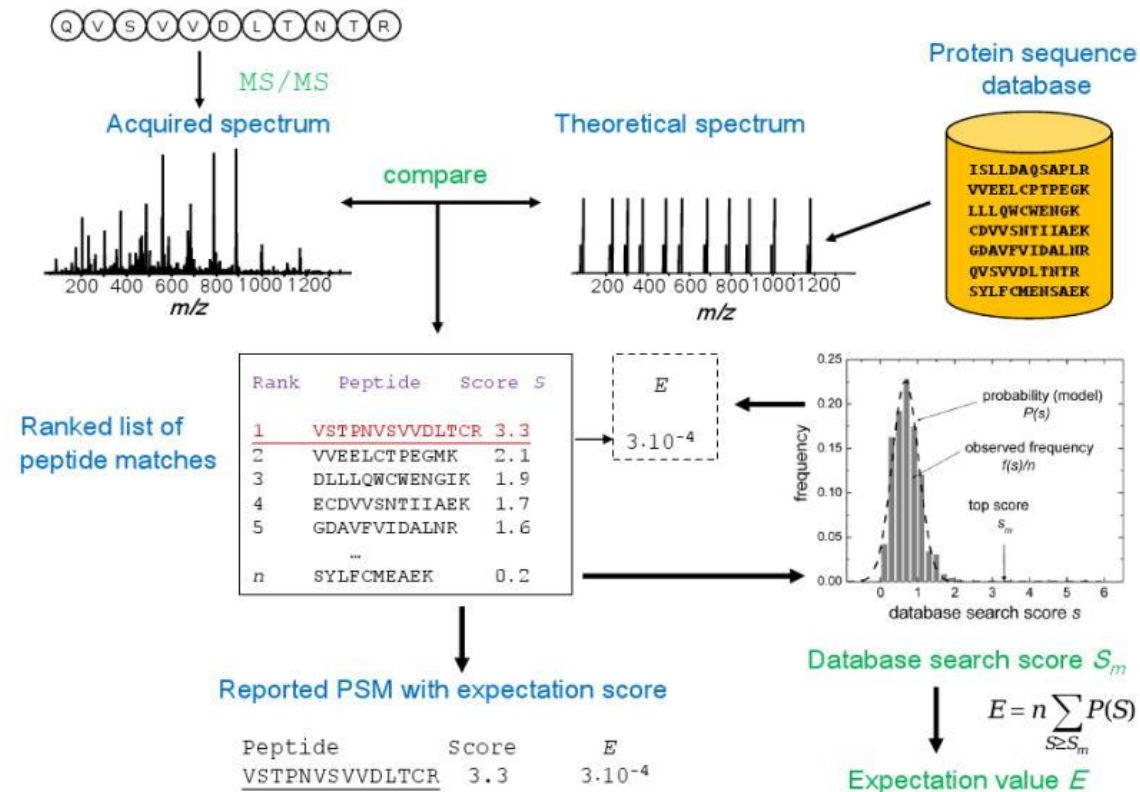
File: C11_Liver_1_msms_Tray02C2_01_2998.02343.02343.2 mass 728.440037

#	peptide	E-value	xcorr	dCn	Sp	rankSp	neutral	mass	ions	protein
1	K.AIK[170]EVK.V	1.73E+000	0.789	0.003	101.9	1	728.443241	4/10	DECOY_1_tr_52196	
2	R.AILGEVK.L	1.78E+000	0.787	0.004	79.4	2	728.443240	4/12	sp Q9D8N0 EF1G_MOUSE +3	
3	R.ALK[170]SVK[170].R	1.80E+000	0.786	0.004	101.9	1	728.443241	4/10	DECOY_0_sp_19583 +3	
4	K.AIVADLK.G	1.80E+000	0.786	0.004	79.4	2	728.443240	4/12	DECOY_0_sp_13685	
4	R.ALK[170]DIK.G	1.80E+000	0.786	1.000	101.9	1	728.443241	4/10	sp Q6ZWQ0 SYNE2_MOUSE +6	

[show search parameters](#)

Peptide Identification from MS/MS spectra

1: Nesvizhskii AI. A survey of computational methods and error rate estimation procedures for peptide and protein identification in shotgun proteomics. J Proteomics. 2010 Oct 10;73(11):2092-123. doi: 10.1016/j.jprot.2010.08.009. Epub 2010 Sep 8. Review. PubMed PMID: 20816881; PubMed Central PMCID: PMC2956504.

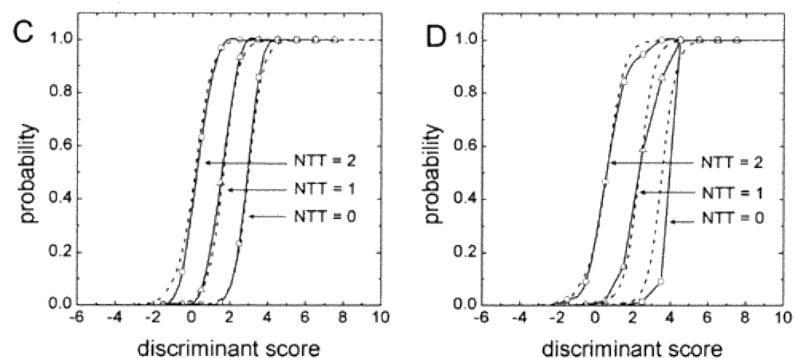
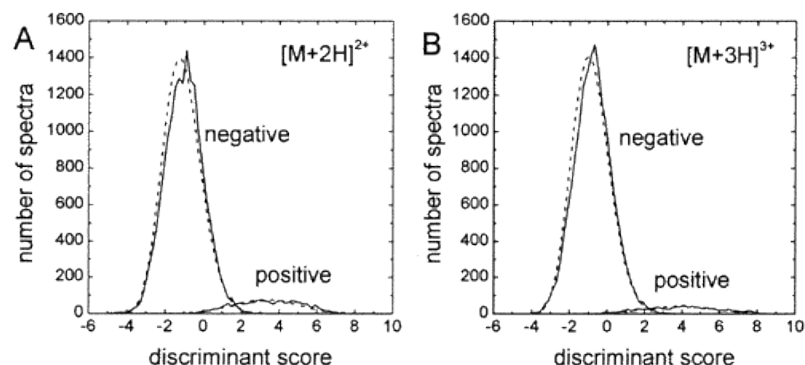


Problems with PSMs

- Different search engines give different results
- Different scoring methods
- Need to 'unify' results
- ➔ answer:
 - Modelling of distribution of best hits factoring other effects
 - NTT: number of tryptic termini
 - NMC: number of missed cleavages
 - Mass accuracy
 - Mixture modelling and Bayes

Peptide Validation from MSMS Search Engine

1: Keller A, Nesvizhskii AI, Kolker E, Aebersold R. Empirical statistical model to estimate the accuracy of peptide identifications made by MS/MS and database search. *Anal Chem.* 2002 Oct 15;74(20):5383-92. PubMed PMID: 12403597.



$$p(+|F,NTT) = \frac{p(F|+)p(NTT|+)p(+)}{p(F|+)p(NTT|+)p(+)+p(F|-)p(NTT|-)p(-)} \quad (8)$$

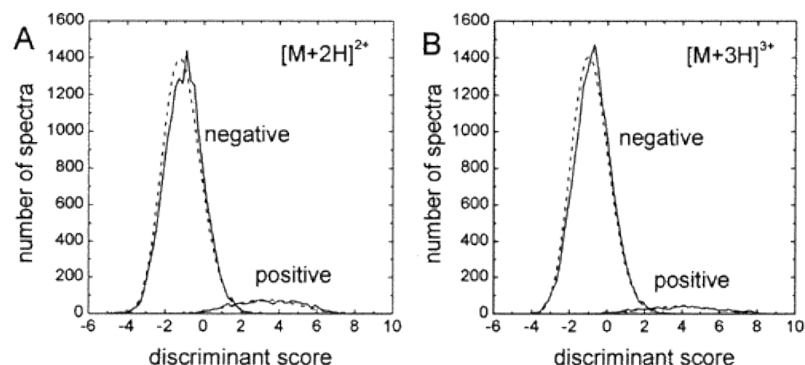
Models optimum discriminant score

Parameters include mass, # amino acids, search engine scores

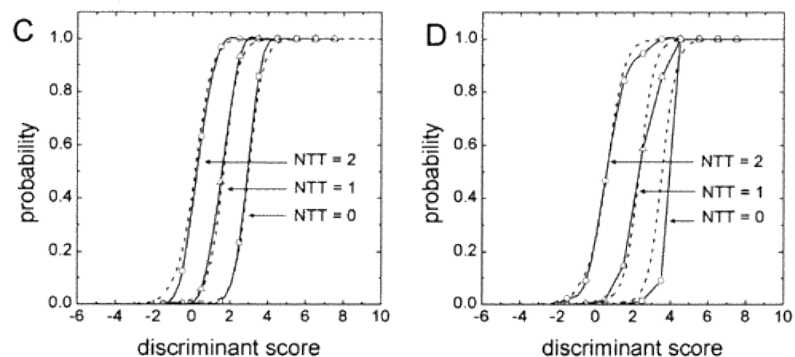
Separate models for charge state

Peptide Validation from MSMS Search Engine

1: Keller A, Nesvizhskii AI, Kolker E, Aebersold R. Empirical statistical model to estimate the accuracy of peptide identifications made by MS/MS and database search. Anal Chem. 2002 Oct 15;74(20):5383-92. PubMed PMID: 12403597.



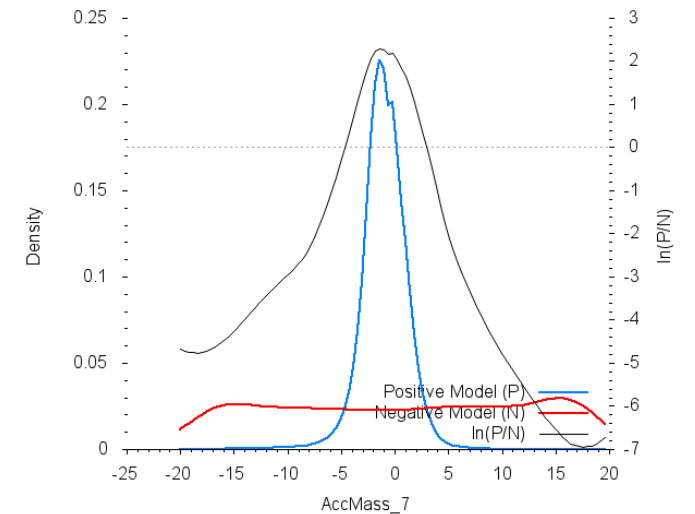
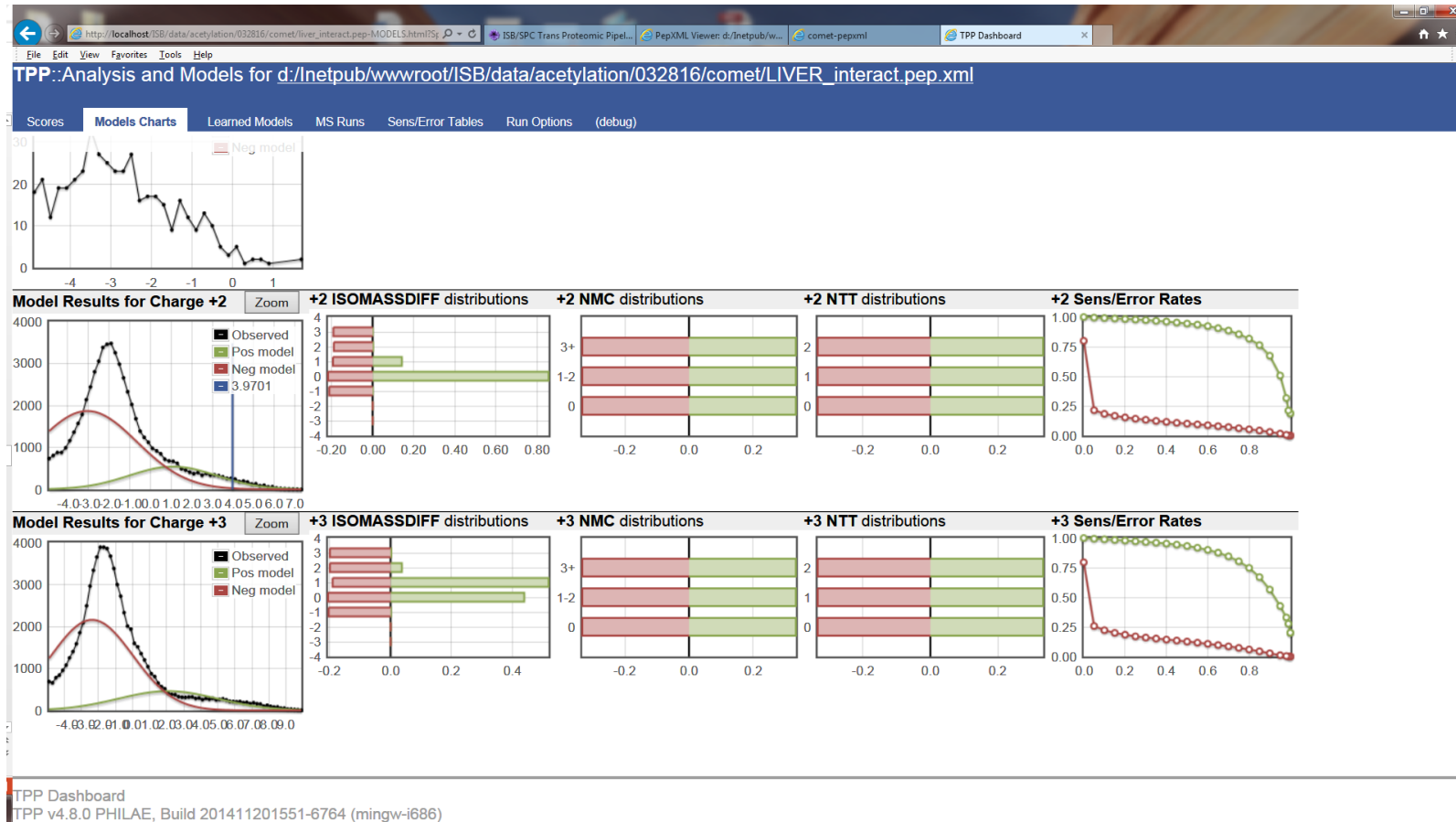
$$p(+|F,NTT) = \frac{p(F|+)p(NTT|+)p(+)}{p(F|+)p(NTT|+)p(+)+p(F|-)p(NTT|-)p(-)} \quad (8)$$



Goal: ease the burden of manually validating 100,000's spectra

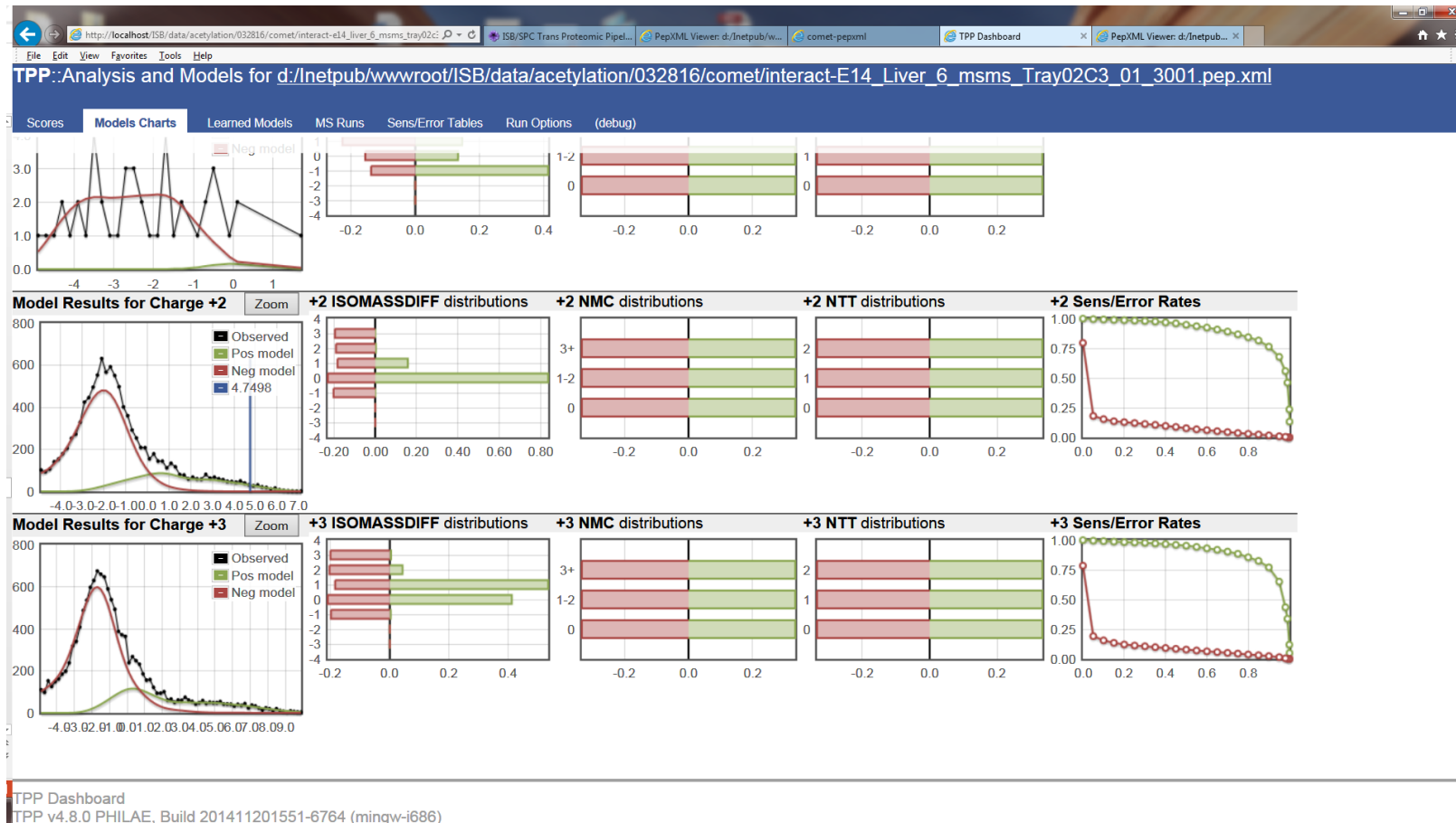
Typical Modelling Results

GOOD model “fit”?



Modelling mass accuracy

Typical Modelling Results



Familiarization with the parameters for modelling is always beneficial.

Error rate estimation of PSMs

Choi H, Ghosh D, Nesvizhskii AI. Statistical validation of peptide identifications in large-scale proteomics using the target-decoy database search strategy and flexible mixture modeling. J Proteome Res. 2008 Jan;7(1):286-92. Epub 2007 Dec 14. PubMed PMID: 18078310.

Decoy databases

- **“Synthetic” sequences not native to the proteome.**
- **Typically “reversed” and added to the search database**
- Lower quality spectra likely to randomly match to “false” sequence.
- Allows an estimate of error rate of PSM

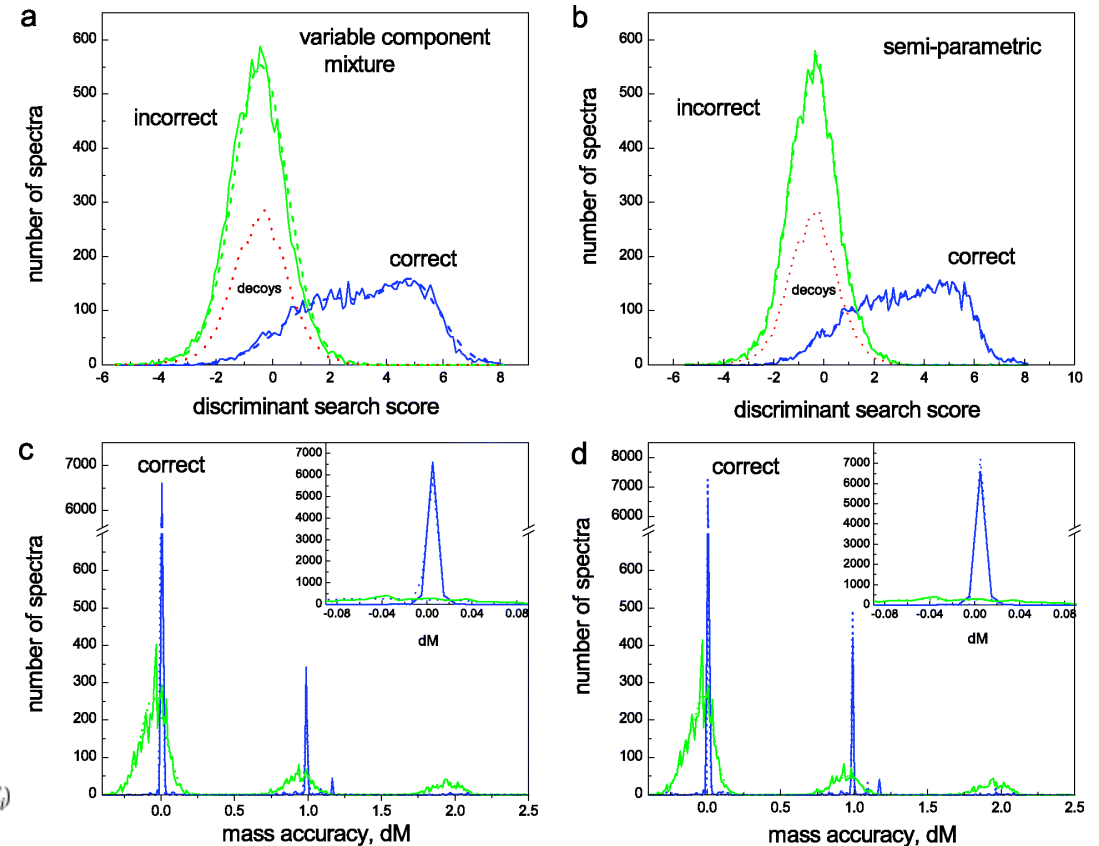
FDR:

		The MS/MS spectrum comes from a peptide sequence in the database	
		True	False
Search reports a match to the correct sequence	True	True positive	False positive
	False	False negative	True negative

www.matrixscience.com

Or model the distributions:

$$\begin{aligned}
 P(+|S_p, E) &= \int_{\Theta} P(+|\Theta, S_p, E) dF(\Theta|S_p, E) \\
 &= \sum_d p(d) \int_{\Theta^{(d)}} P(+|\Theta^{(d)}, S_p, E) dF(\Theta^{(d)}|S_p, E) \\
 &\approx \frac{1}{I} \sum_{k=1}^I 1(+|\Theta_k, S_p, E)
 \end{aligned}$$



Error rate estimation of PSMs

1: Käll L, Storey JD, MacCoss MJ, Noble WS. Posterior error probabilities and false discovery rates: two sides of the same coin. J Proteome Res. 2008 Jan;7(1):40-4. Epub 2007 Dec 4. Review. PubMed PMID: 18052118.

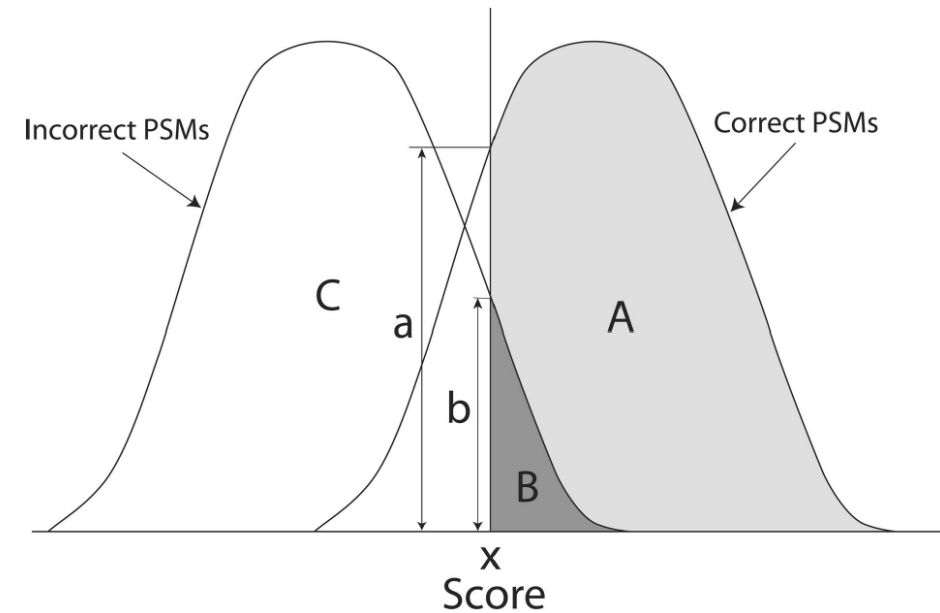
Decoy databases

- Lower quality spectra likely to randomly match to “false” sequence.
- Allows an estimate of error rate of PSM

FDR:

		The MS/MS spectrum comes from a peptide sequence in the database	
		True	False
Search reports a match to the correct sequence	True	True positive	False positive
	False	False negative	True negative

www.matrixscience.com



$$\text{FDR} = B/(A + B)$$

$$\text{PEP} = b/(a + b)$$

Figure 1. Two complementary methods for assessing statistical significance.

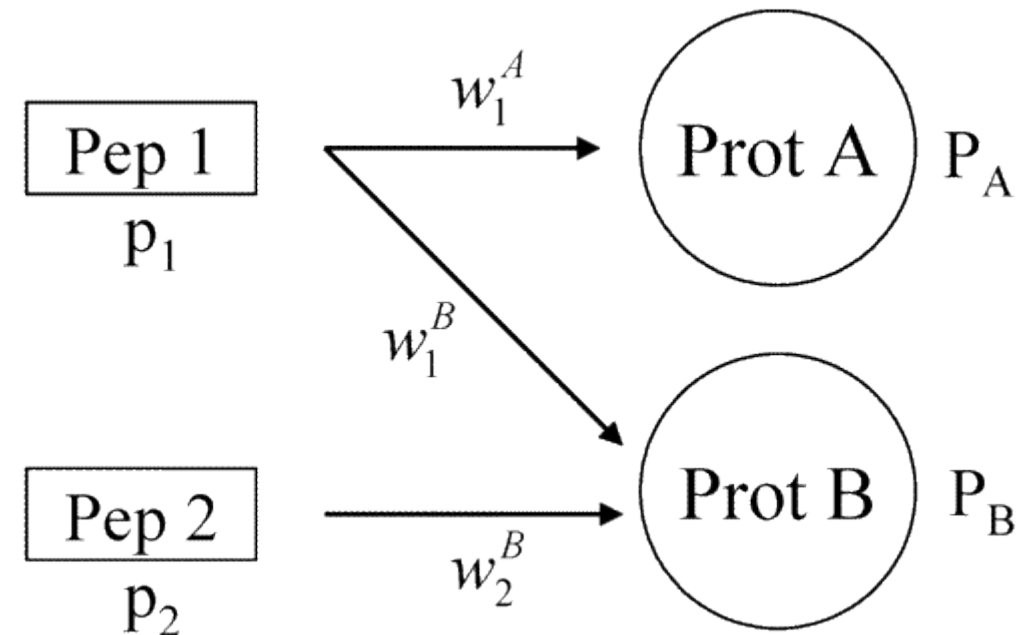
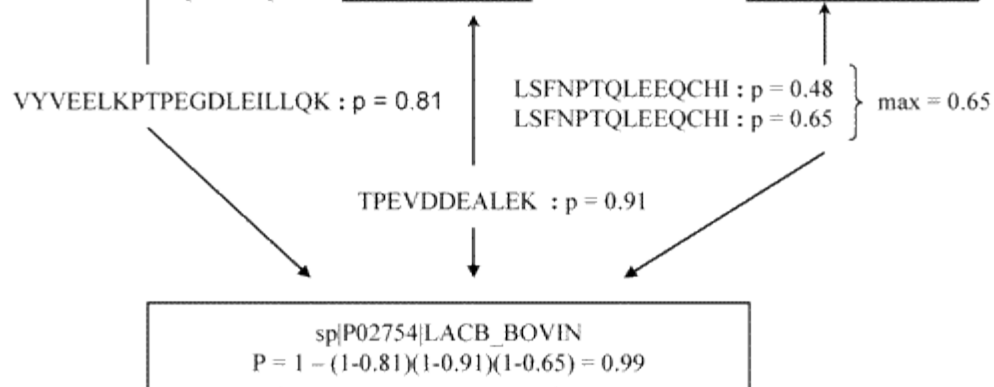
Problems with PSMs for protein inference

- Sequences are redundant
- Isoforms and peptide protein mapping
- Replicate PSMs acquired (multiple samples).....which peptide probability is best?
- What's an appropriate threshold?
- ANSWER→
 - NSP...Number sibling peptides model..is the sequence unique?
 - Best observable peptide
 - Protein probabilities
 - “decoy” databases

Statistical Model for Protein Inference

>sp|P02754|LACB_BOVIN BETA-LACTOGLOBULIN PRECURSOR (BETA-LG)
(ALLERGEN BOS D 5) - Bos taurus (Bovine).

MKCLLLALALTCGAQALIVTQTMKGLDIQKVAGTWYSLAMAASDISLLDAQSAPLRVYV
EELKPTPEGDLEILLQKWENGEC AQKKIIAEKTKIPAVFKIDALNENKLVLDTDYKKYLL
FCMENSAEPEQSLACQCLVRTTPEVDDEALEKFDKALKALPMHIRLSFNPTQLEEQCHI



$$p(+|D, \text{NSP}) = \frac{p(+|D)p(\text{NSP}|+)}{p(+|D)p(\text{NSP}|+) + p(-|D)p(\text{NSP}|-)} \quad (5)$$

Protein Identification from Peptides

Shared peptides

ProtXML Viewer results for interact-e14_liver_6_msms_tray02c3_01_3001.prot.xml

ProtXML Viewer :: d:/Inetpub/wwwroot/TSB/data/acetylation/032816/comet/interact-e14_liv
http://localhost/tpp-bin/ProtXMLViewer.pl?file=d:/Inetpub/wwwroot/TSB/data/acetylation/

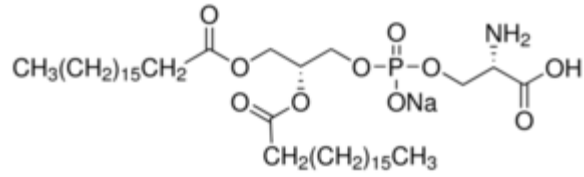
#	Main Entry Accession	NTT	Probability	# Tot PSMs	% Coverage / Weight	% Spectrum ids / NSP
94	Group # 1	1.0000	<div></div>			
94a	sp A2AS89 SPEB_MOUSE Agmatinase, mitochondrial OS=Mus musculus GN=Agmat PE=1	1.0000	<div></div>	14	8.9%	0.25%
+3	RSVDEGLLDISK ₁₇₀ R	<div></div>	0.9973	2	1.00	3.92
+2	SVDEGLLDISK ₁₇₀ R	<div></div>	0.9998	2	1.00	3.91
+3	SVDEGLLDISK ₁₇₀ R	<div></div>	0.9829	2	1.00	3.98
+2	VVLAEDCWMK ₁₇₀ SLVPLMAEVR	<div></div>	0.9992	5	1.00	3.91
+3	VVLAEDCWMK ₁₇₀ SLVPLMAEVR	<div></div>	0.9982	3	1.00	3.92
94b	tr Q8R0Z1 Q8R0Z1_MOUSE Agmat protein (Fragment) OS=Mus musculus GN=Agmat PE=	0.0000	<div></div>	8	0%	0%
+2	VVLAEDCWMK ₁₇₀ SLVPLMAEVR	<div></div>	0.9987	5	0.00	0.99
+3	VVLAEDCWMK ₁₇₀ SLVPLMAEVR	<div></div>	0.9969	3	0.00	1.00

ProtXML Viewer
TPP v4.8.0 PHILAE, Build 201411201551-6764 (mingw-i686)

Small Molecule MS/MS

- Unknown predictable patterns for electrospray ionization- MS/MS (ESI-MS/MS)
- Empirical knowledge versus predictive methods
 - Empirical patterns, 1970's, 1990's
 - More accurate
 - Dependent on identifying and storing validated spectra
 - Predictive: 2010's

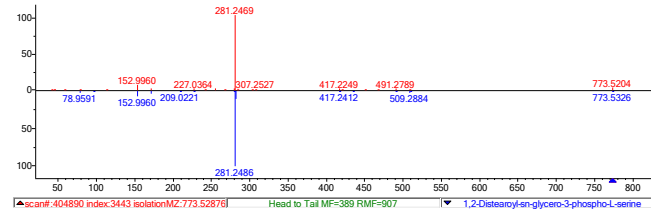
MS² at work:



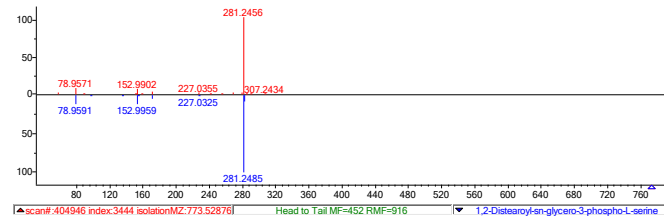
PS(18:0/18:0)

1,2-Distearoyl-sn-glycero-3-phospho-L-serine

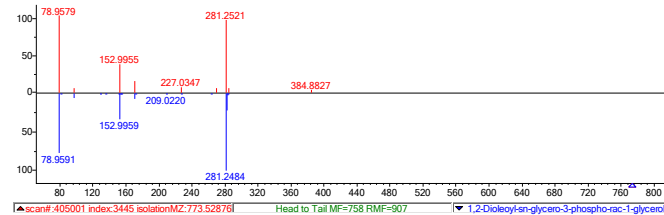
C₄₂H₈₂NO₁₀P



40V [M-H]- : 773.52876



60V [M-H-NH₃]- : 791.567635



80V [M-H]- : 773.52876

Early work

Optimization and testing of mass spectral library search algorithms for compound identification, Stein and Scott (1994)

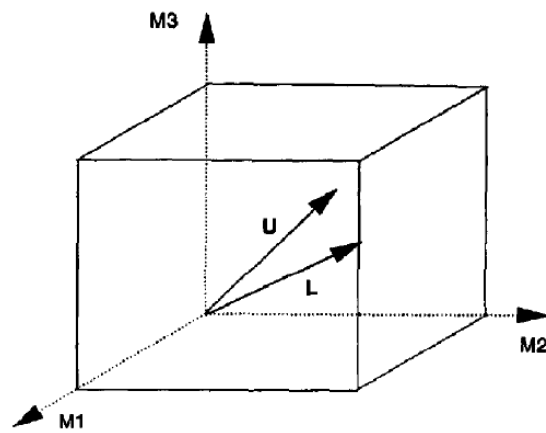
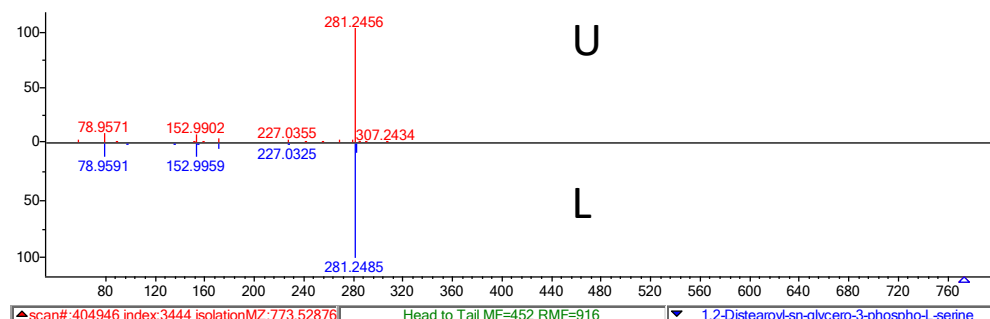


Figure 1. Vector representation of a hypothetical three-peak unknown (U) and library (L) mass spectrum in three-dimensional space (peaks have mass M1, M2, and M3).

Table 1. Search algorithms investigated^a

Euclidean distance

$$\left(1 + \frac{\sum (w_L - w_U)^2}{\sum w_U^2}\right)^{-1}$$

Absolute Value Distance

$$\left(1 + \frac{\sum |w_L - w_U|}{\sum w_U}\right)^{-1}$$

Hertz et al. [9]

$$\frac{\text{average of weighted peak intensity ratios}}{[1 + \text{fraction of unmatched intensities}]}$$

Dot-product (cosine), F_D

$$\frac{(\sum w_L w_U)^2}{\sum w_L^2 \sum w_U^2}$$

Probability-based matching (PBM) [5b, 5d, 10c]

Uses probability that, by chance, peaks match within an abundance window (W value) by using uniqueness values for mass (U value) and abundance (A) along with a variety of rules and correlation tables.

Composite:

$$\frac{N_U F_D + N_{L\&U} F_R}{N_U + N_{L\&U}}$$

F_D = Dot-Product Term Above

F_R = Ratio of Peak Pairs (below)

$W = [\text{Peak Intensity}]^n [\text{Mass}]^m = \text{Weighted Intensity}$

N = Number of peaks

$$F_R = \frac{1}{N_{L\&U}} \sum_i^{L\&U} \left(\frac{w_{L,i}}{w_{L,i-1}} \frac{w_{U,i-1}}{w_{U,i}} \right)^n$$

where $n = 1$ or -1 when the term in parentheses is less than or greater than unity, respectively

^a Subscripts: L = library, U = unknown, L & U = peak in both library and unknown spectrum.

Early work

Optimization and testing of mass spectral library search algorithms for compound identification, Stein and Scott (1994)

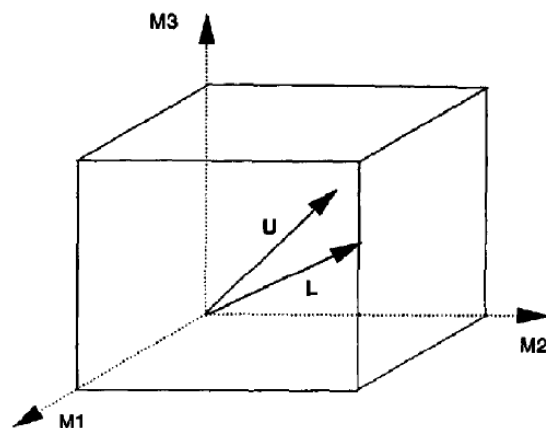


Figure 1. Vector representation of a hypothetical three-peak unknown (U) and library (L) mass spectrum in three-dimensional space (peaks have mass M1, M2, and M3).

Table 2. Performance of various algorithms

Algorithm	% Correct at Rank			Scaling /Comments	
	1	1-2	1-3	Mass Power	Intensity Power
Dot-product	72.9	85.9	90.8	1	0.5 ^a
	73.2	86.3	91.0	1	0.5 ^b
	72.8	85.9	90.8	1	0.5
Optimized	74.9	86.9	91.7	3	0.6
Euclidean distance	65.8	79.3	84.9	0	0.5
	69.9	82.9	88.2	1	0.5
Optimized	71.9	83.9	88.9	2	0.6
Absolute distance	61.4	74.9	81.2	0	1
	66.8	79.4	85.1	1	1
Optimized	67.9	80.3	85.5	2	0.9
PBM	57.1	71.5	78.5	"k value" ^c	
	64.0	77.7	84.3	Reliability ^d	
	64.7	78.4	84.8	Complete ^e	
Hertz et al.	59.9	73.9	81.1	0	See ref 9
Optimized	64.4	77.2	83.2	2	0.5
Composite	75.7	88.0	92.5	3, 0 ^f	0.5, 1 ^f

^a Used local and global "normalization" [6].

^b Used local "normalization" [6].

^c Based on "reverse-search" overall spectral match factor [10].

^d All recommended features except "quadratic scaling" [5b, d].

^e All recommended features [5b, d].

^f The first value is for the dot-product term and the second value is for the peak ratios term. The second and third power of the mass were equally effective for the first term.

Small Molecule MS/MS

- Still predominantly use DOT product for library- unknown matching.
 - Comprehensive limited by lack of library entries....eg not enough spectra.
- New predictive methods attempt to produce theoretical spectra from known structures
 - Then attempt to match experimental spectra to these theoretical entries.

Competitive fragmentation modeling of ESI-MS/MS spectra for putative metabolite identification, Allen et al (2015)

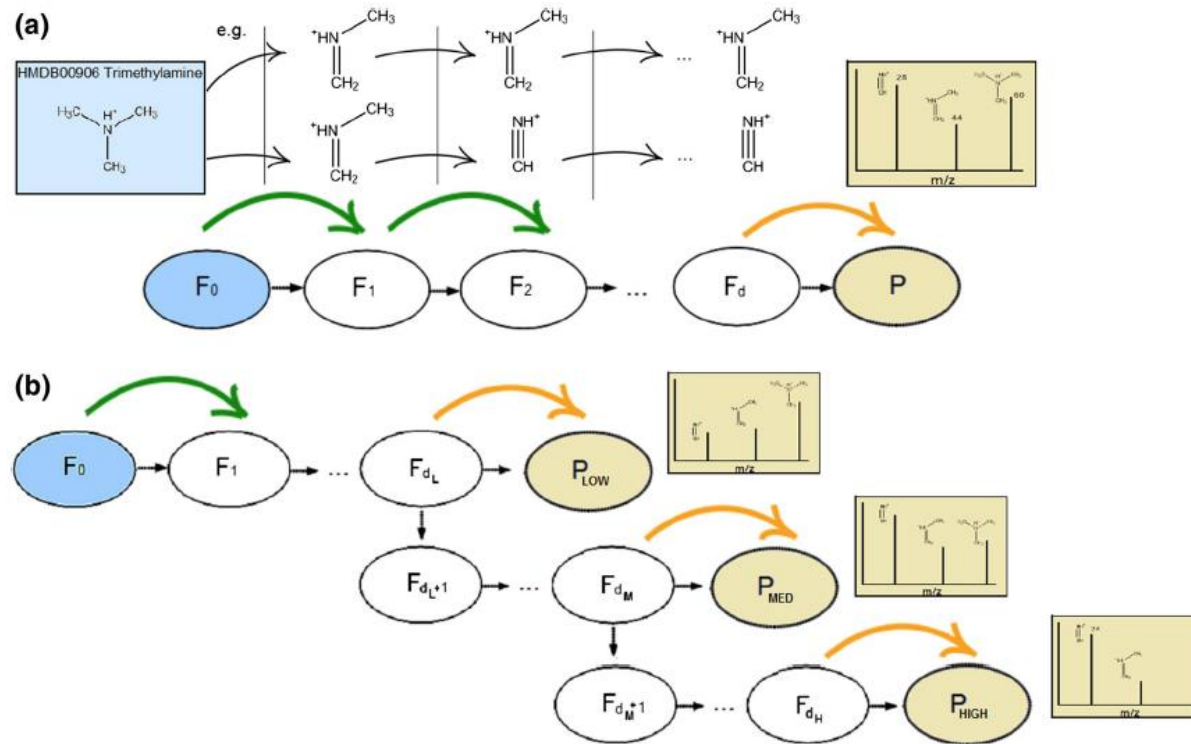
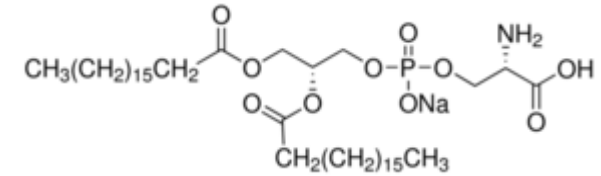


Fig. 1 **a** Single energy competitive fragmentation model (SE-CFM): a stochastic, Markov process of state transitions between charged fragments. **b** Combined energy competitive fragmentation model

(CE-CFM): an extension of SE-CFM that combines information from multiple collision energy spectra into one model

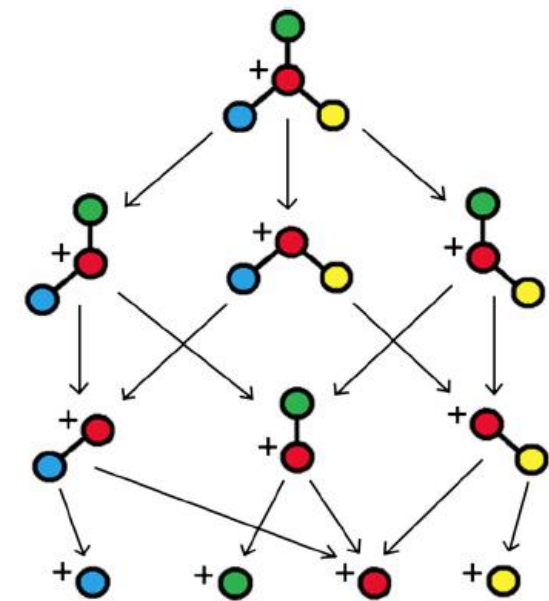


Fig. 2 An abstract example of a fragmentation graph, showing a directed acyclic graph of all possible ways in which a particular charged molecule may break to produce smaller charged fragments

Competitive fragmentation modeling of ESI-MS/MS spectra for putative metabolite identification, Allen et al (2015)

- Very complex problem
 - Competing moieties for fragmentation
 - Fragmentation tendency
 - Sequential transition state models

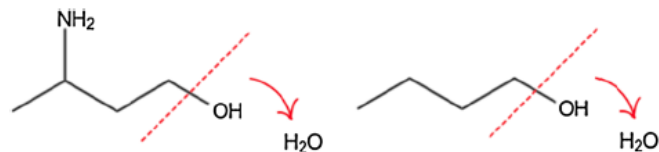


Fig. 3 Two similar breaks, both resulting in an H_2O neutral loss. The *right case* should be assigned a higher probability, as in the *left case*, the NH_3 is also likely to break away, reducing the probability of the H_2O loss

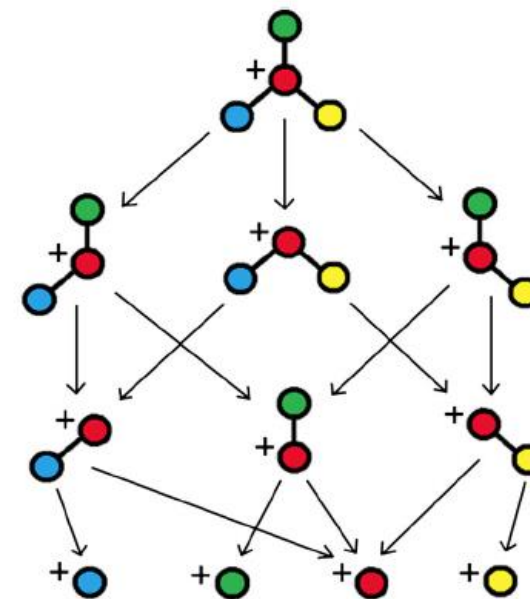
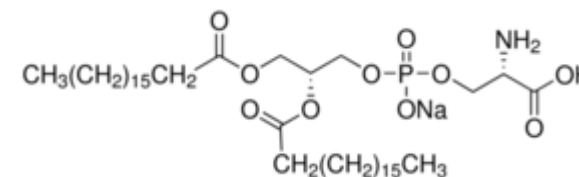
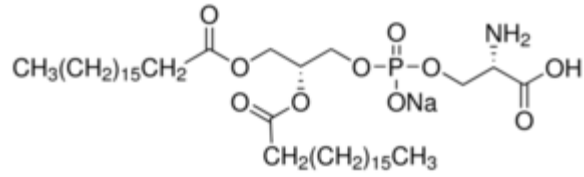


Fig. 2 An abstract example of a fragmentation graph, showing a directed acyclic graph of all possible ways in which a particular charged molecule may break to produce smaller charged fragments

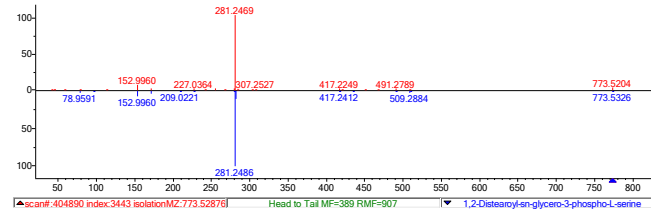
MS² at work:



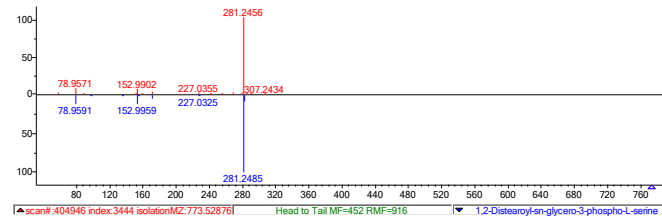
PS(18:0/18:0)

1,2-Distearoyl-sn-glycero-3-phospho-L-serine

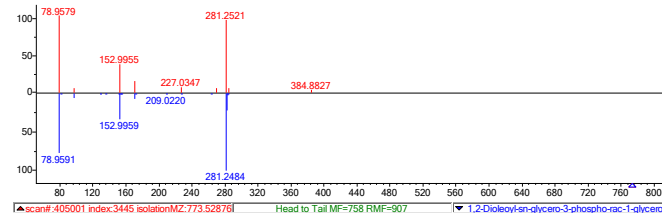
C₄₂H₈₂NO₁₀P



40V [M-H]⁻ : 773.52876



60V [M-H-NH₃]⁻ : 791.567635



80V [M-H]⁻ : 773.52876

Small Molecule MS/MS

- Old reliable methods remain in use
 - Labor intensive to produce search libraries
 - Not comprehensive enough
 - Not platform / laboratory specific
- New methods promising
 - Comprehensive / high specificity theoretical databases
 - Not truly benchmarked