

# Mass spectrometry-based proteomics

Lukas Käll

# Overview

## Mass spectrometry

- Ionization sources
- MS Technologies
- Fragmentation

- Statistics
- Post-translational Modifications (PTMs)

## Proteomics

- Why?
- Dynamic range
- Complexity
- Limitations

## Protein inference

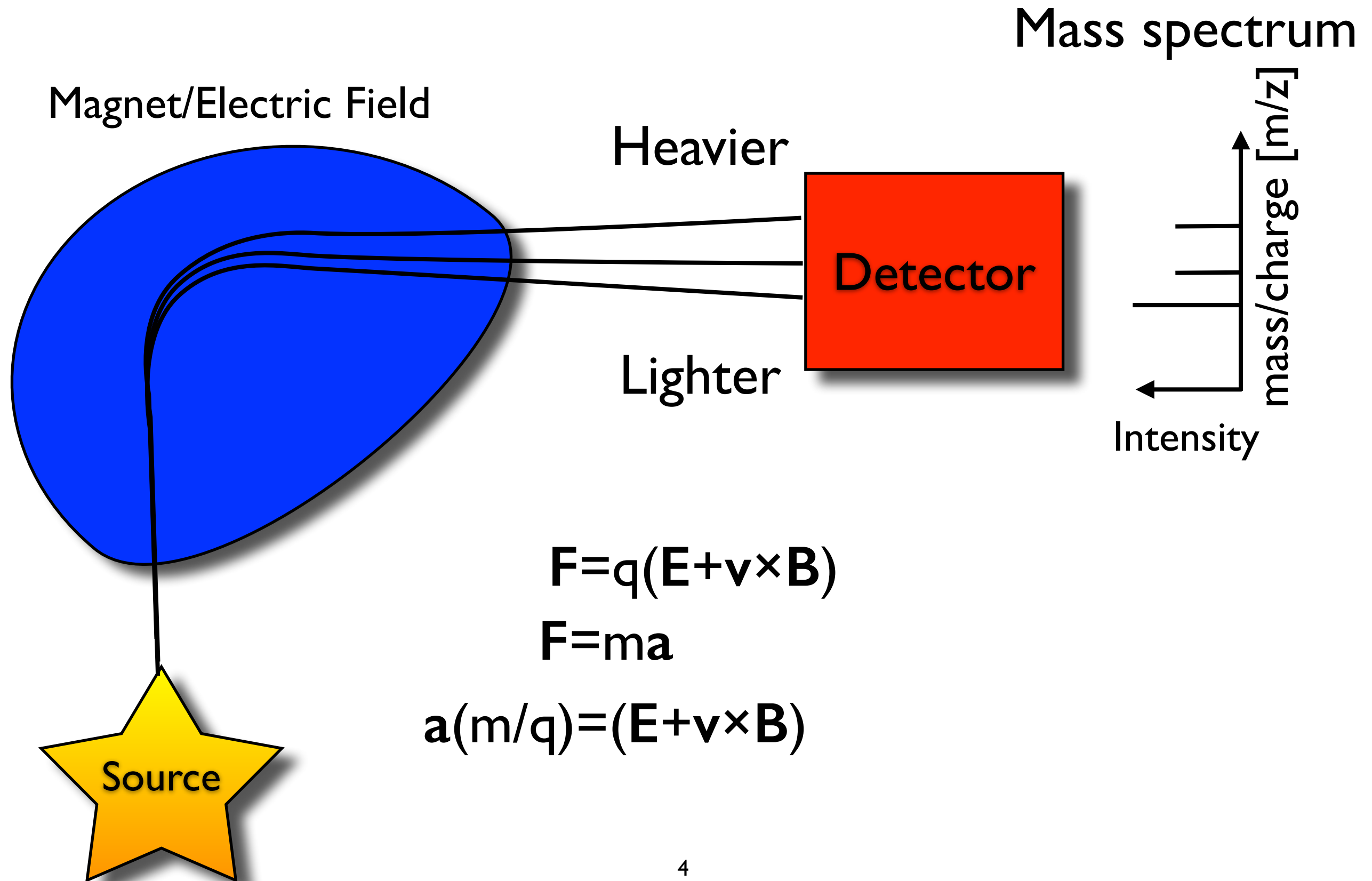
## Peptide Identification

- Search engines



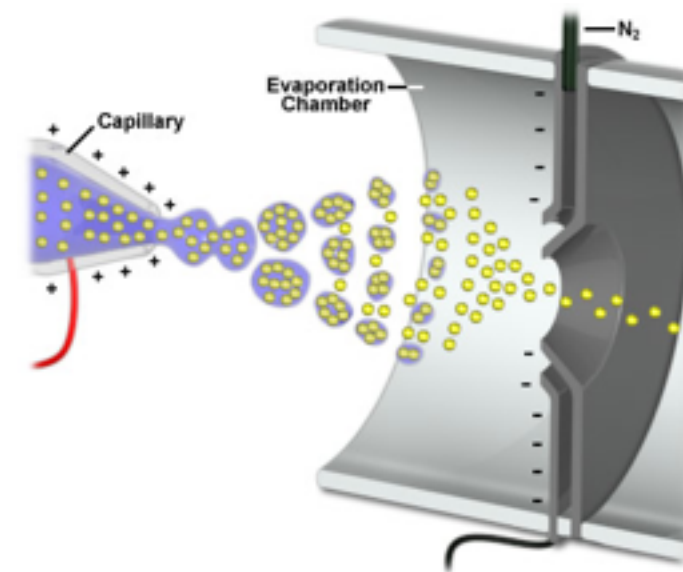
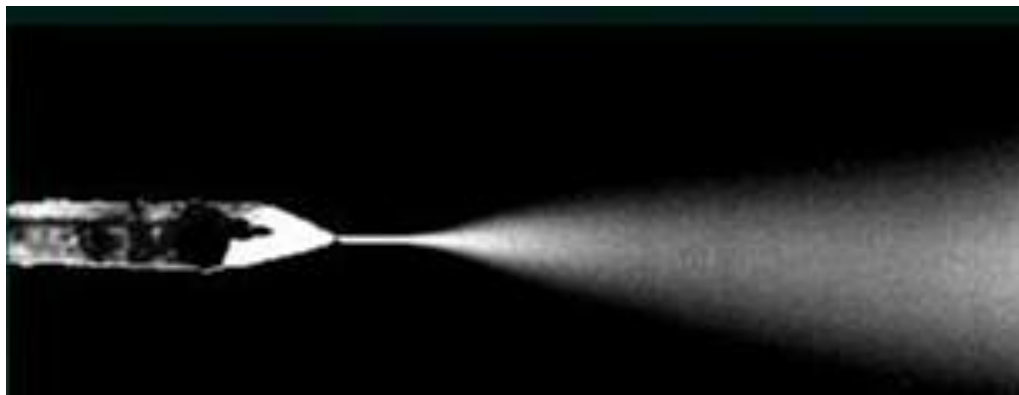
# Mass spectrometry

# Mass spectrometry



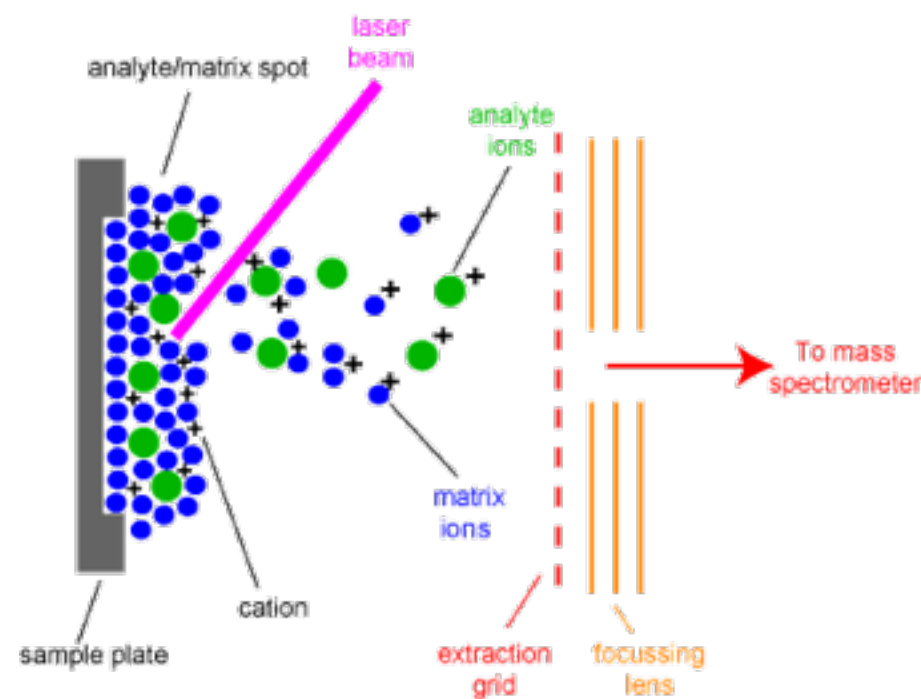
# Ion-Sources

## Electro spray ionization



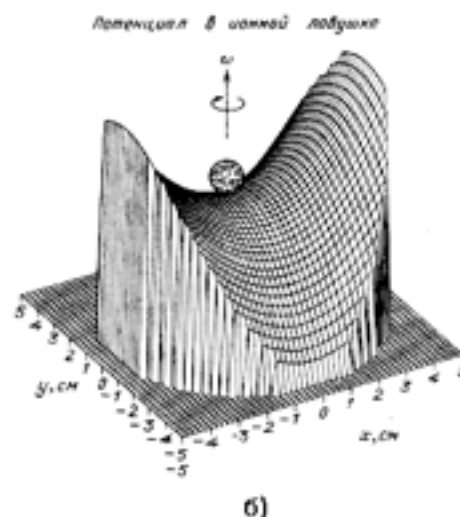
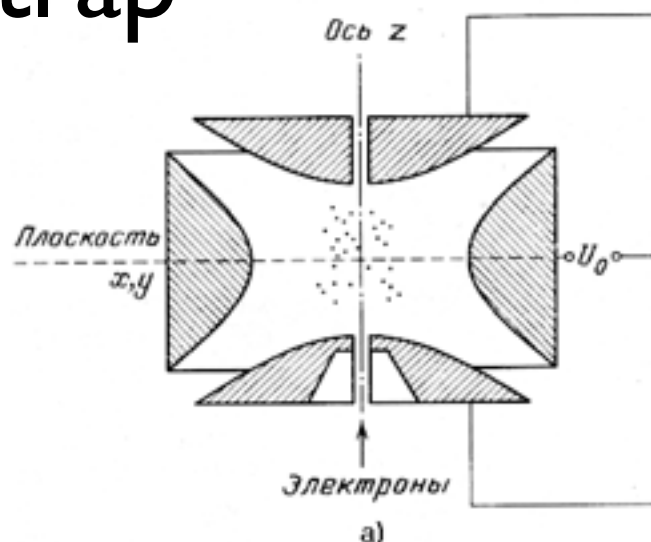
Nobel Prize 2002  
John Fenn

## MALDI



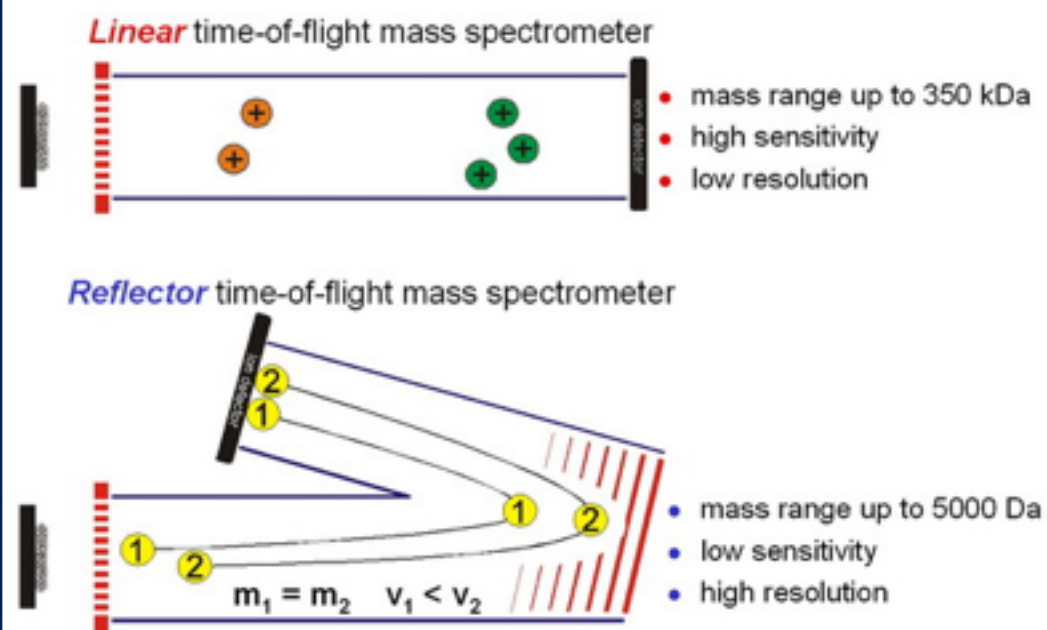
# Separation

## Ion trap

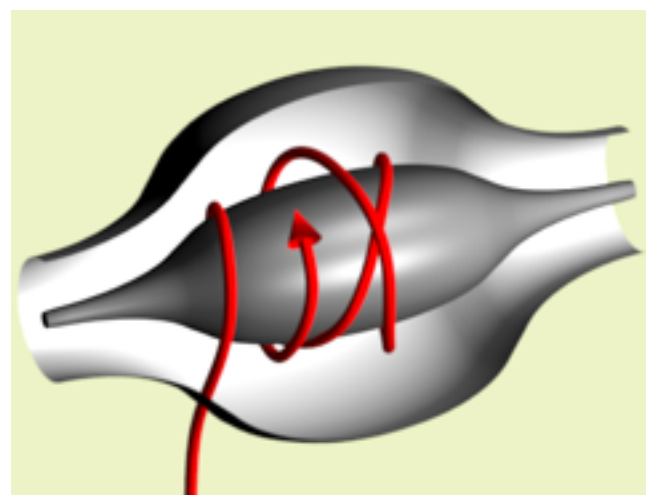


- Ions kept in place by AC current generating an electric field
- The ions will leave the trap as we increase the voltage of the control in order of their  $m/z$  (lowest first)

## Time of Flight (TOF)



## OrbiTrap

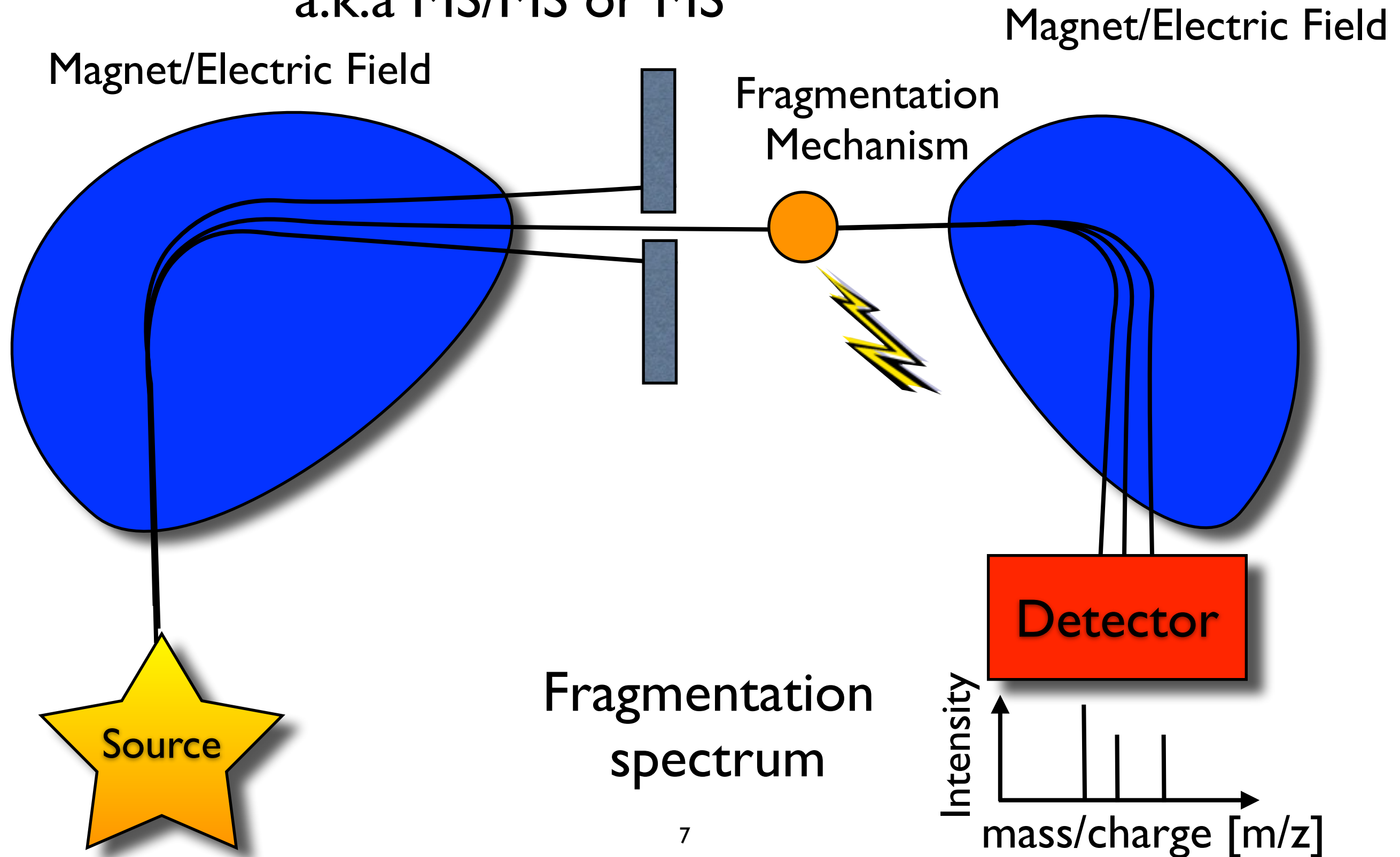


- Electrostatic attraction to the inner electrode is balanced by centrifugal forces.
- Detect the ions by their movement along the axis of the electrode.
- High mass accuracy (1–2 ppm)

<http://www.youtube.com/watch?v=KjUQYuy3msA>

# Tandem mass spectrometry

a.k.a MS/MS or MS<sup>2</sup>



# What do they look like?



Time of Flight



Orbitrap



# Why do we need proteomics?

- We already know the sequence of the genome!
- We know how to measure transcription!
- What else do we need to know?

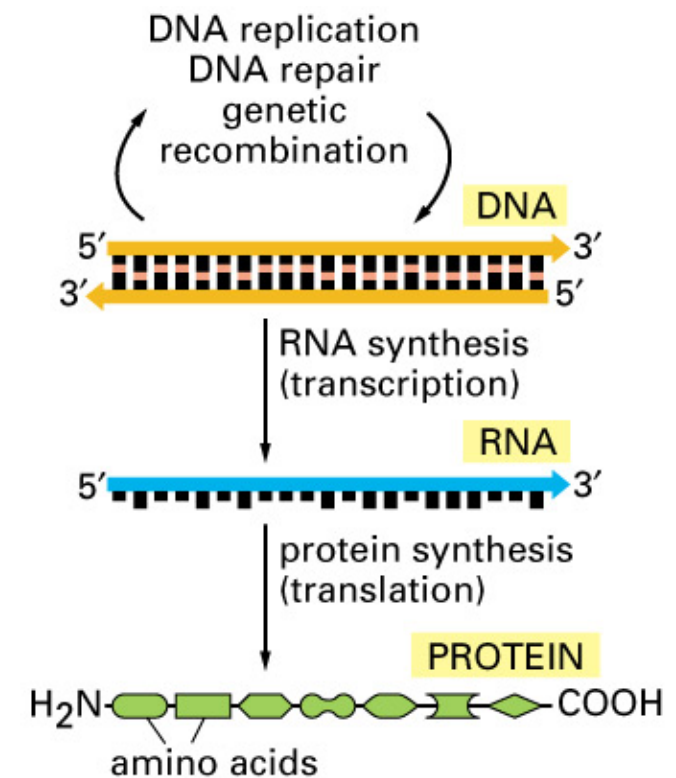


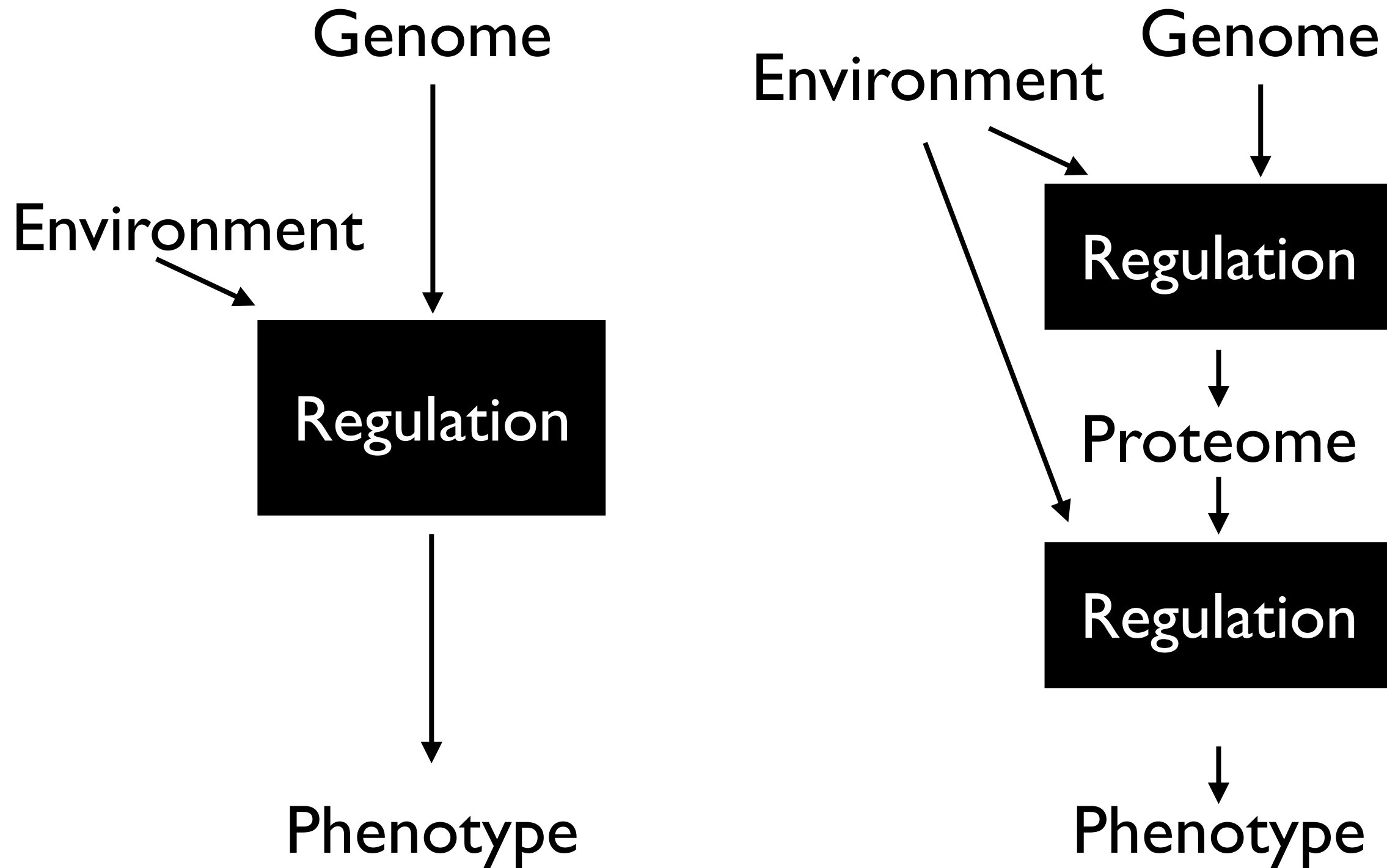
Figure 6-2. Molecular Biology of the Cell, 4th Edition.

# Same DNA, different configuration of proteins



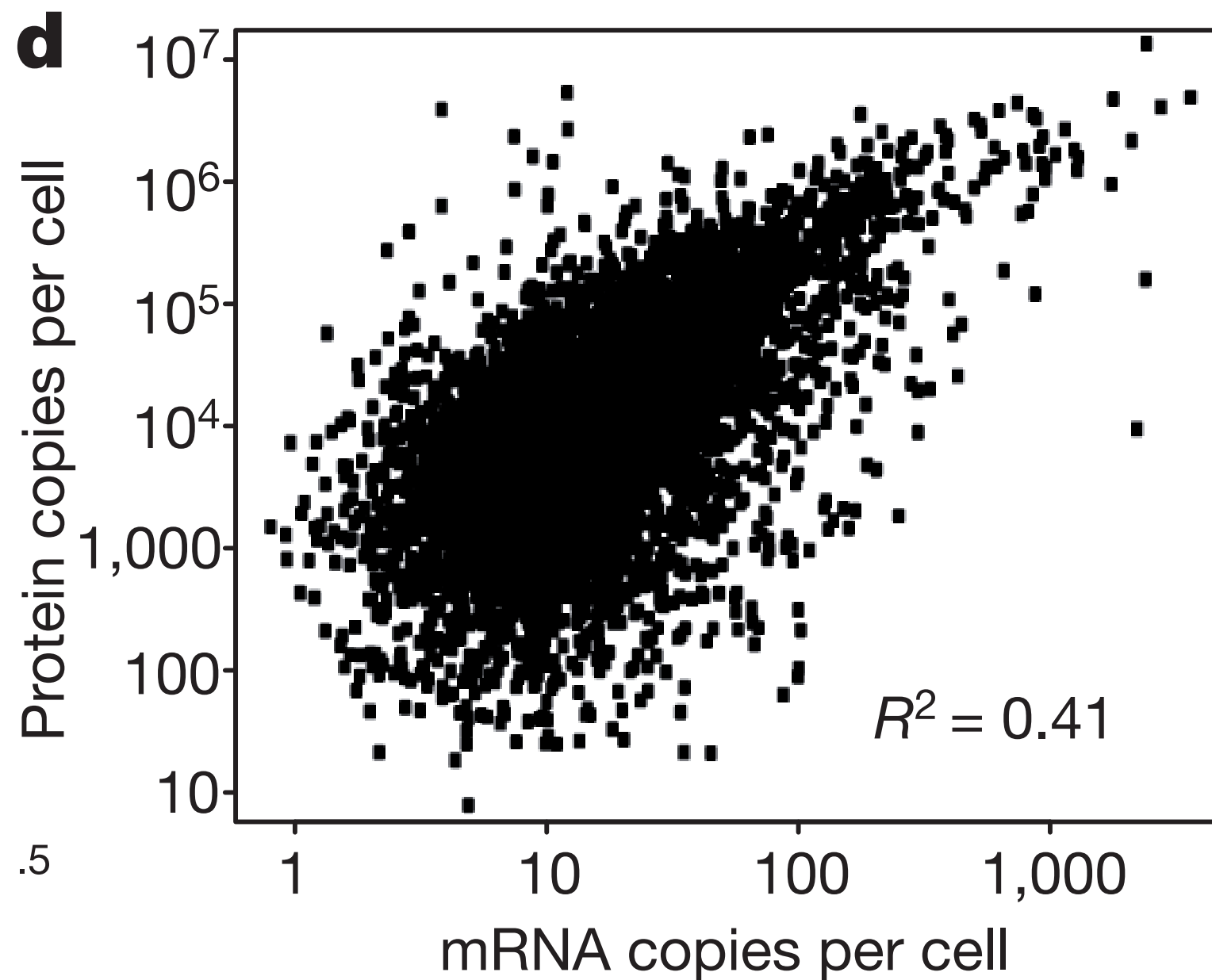
<https://youtu.be/jEtaqmW3ZK4>

# Proteins are closer to the Phenotype



# mRNA levels correlate only weakly with protein levels

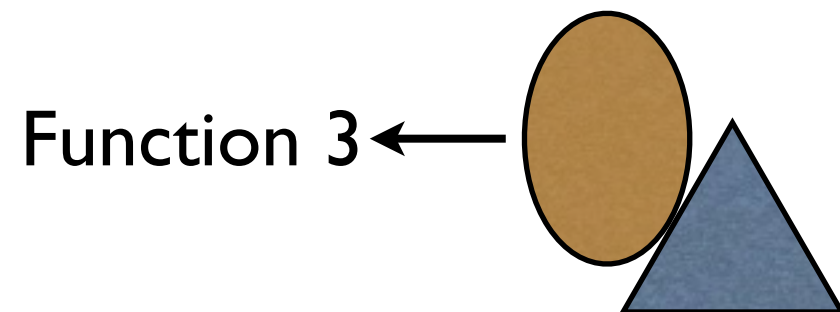
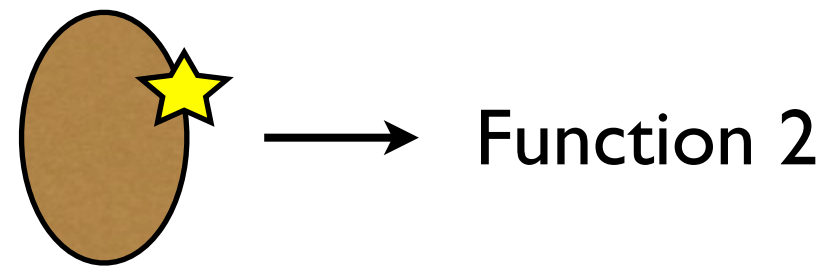
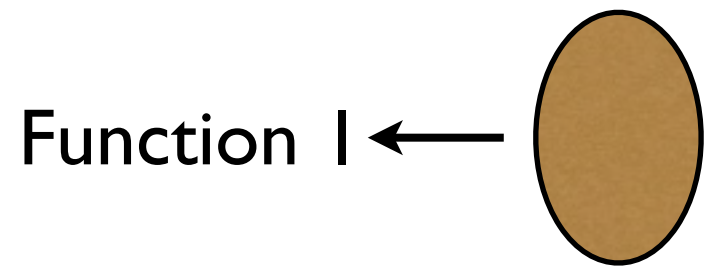
[Schwanhäusser *et al.* Nature 2011]



# One protein - Many functions

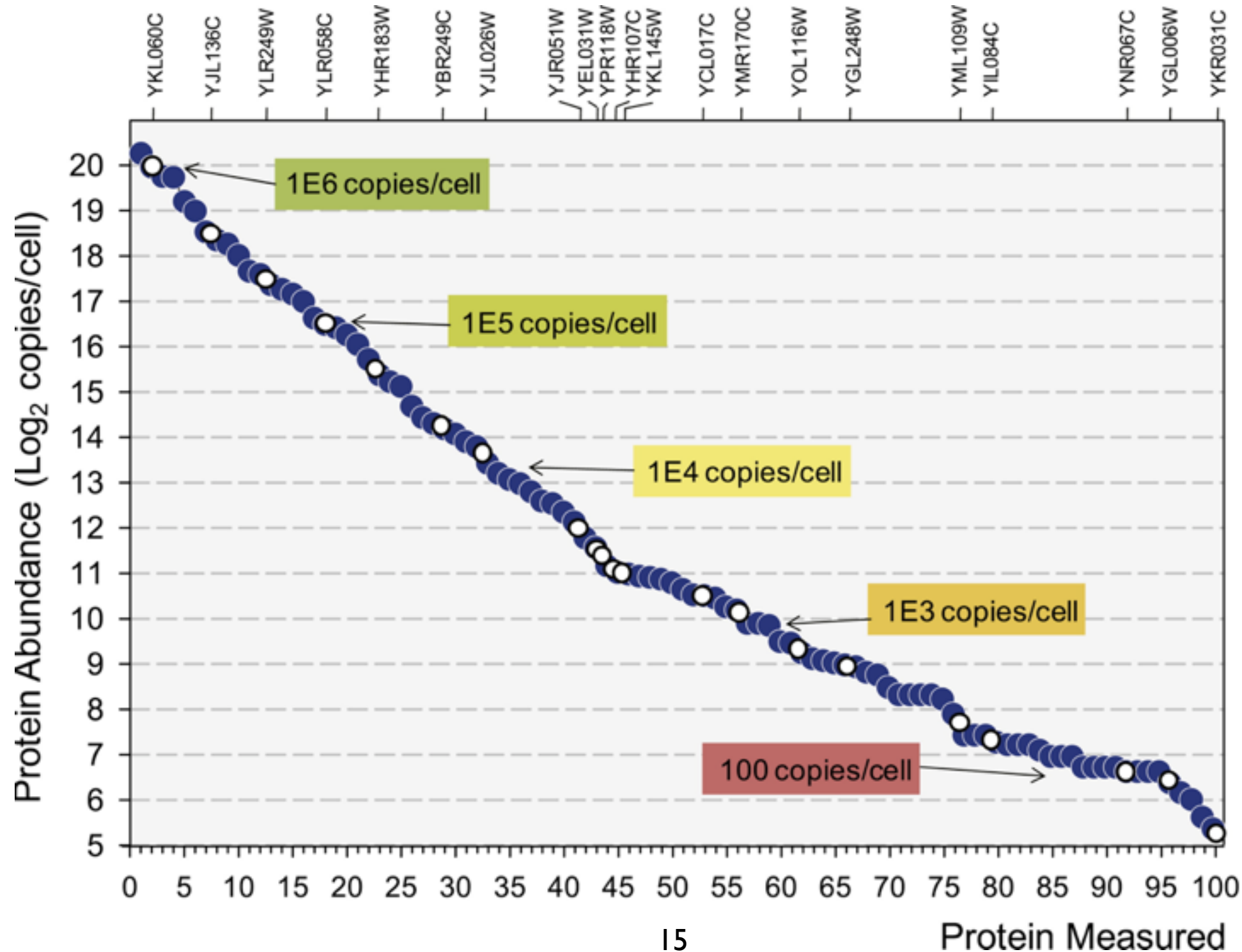
Proteins undergo post-translational modifications

Proteins may have different functions in different complexes or compartments



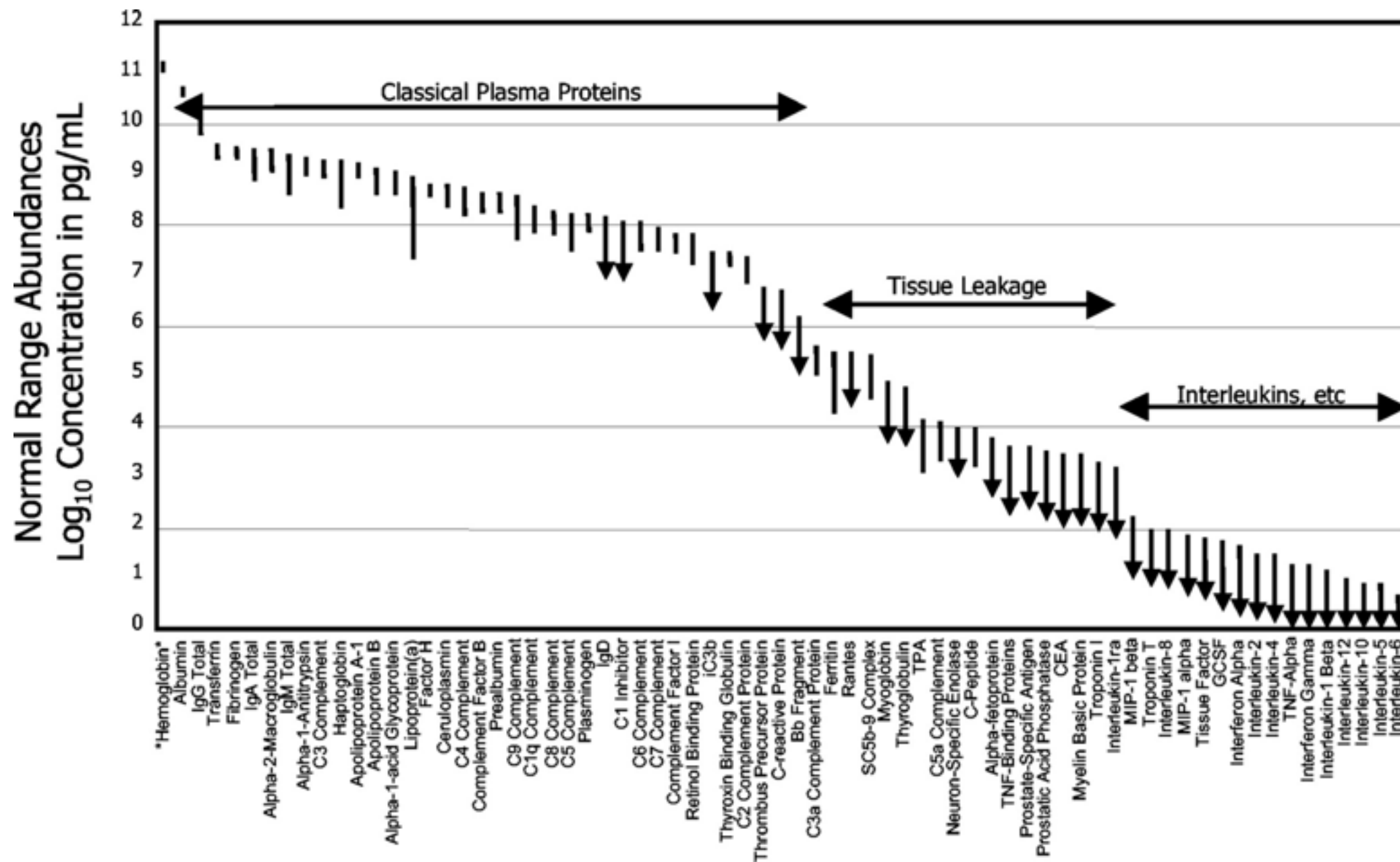
Why don't everybody  
do proteomics?

# Proteins concentration in yeast range $>4$ orders of magnitude



[Picotti et al Cell 2009]

# Protein concentration in blood plasma range >10 orders of magnitude



Difference between earth diameter and 1 mm is about 10 orders of magnitude

[Andersson & Andersson, MCP 2002]



# How to define a protein?

## *Definition*

## *Occurrence in Human*

Protein coding ORF

21,257 [Ensembl]

Splice variant

148,792 [Ensembl]

Protein species

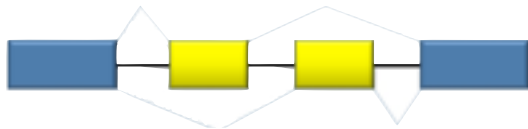
$> 10^6$

Cell specific  
protein species

$> 10^7$

PTMs

Sequence  
rearrangements,  
mutations

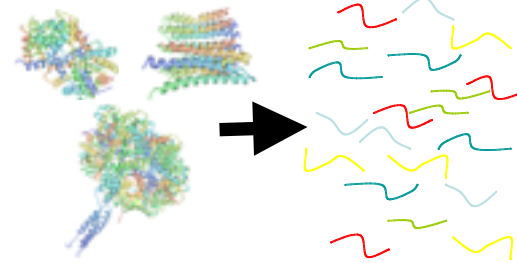


# Proteomic techniques

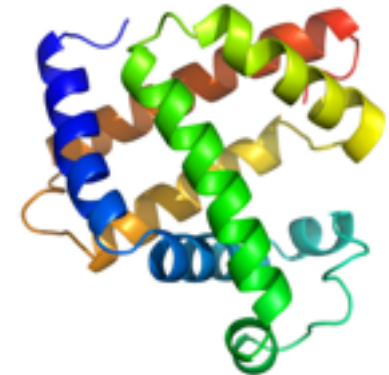


Mass spectrometry-based approaches

e.g.



Shotgun proteomics

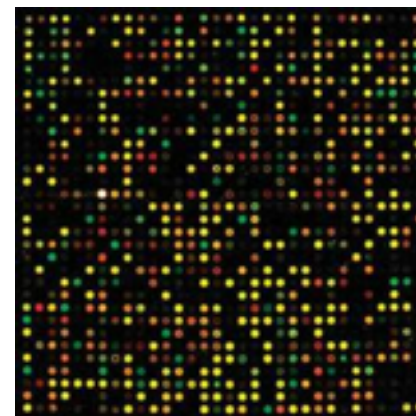


Top-down proteomics

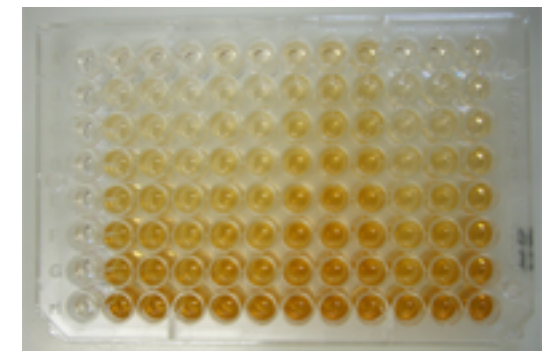


Antibody-based approaches

e.g.



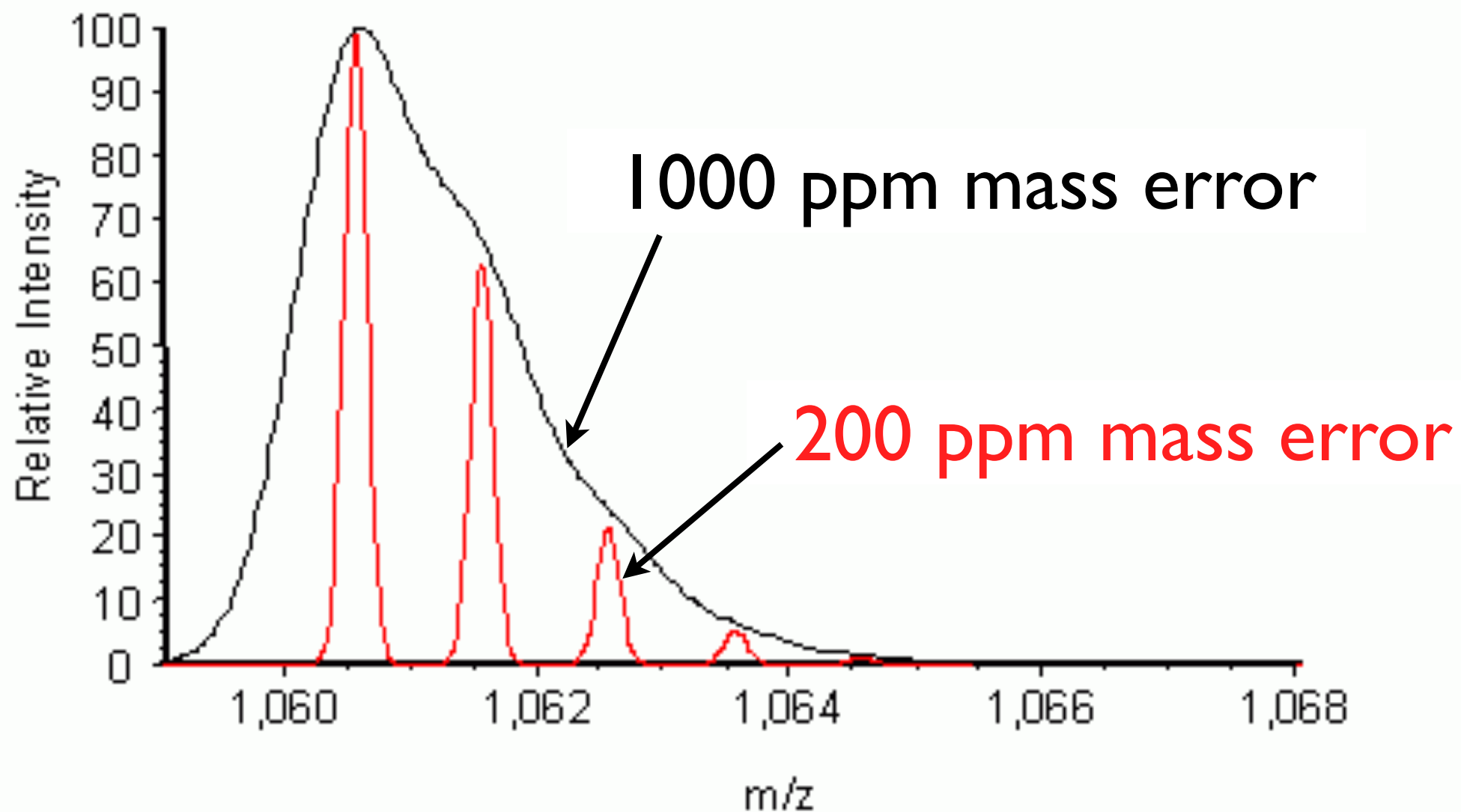
Protein Arrays



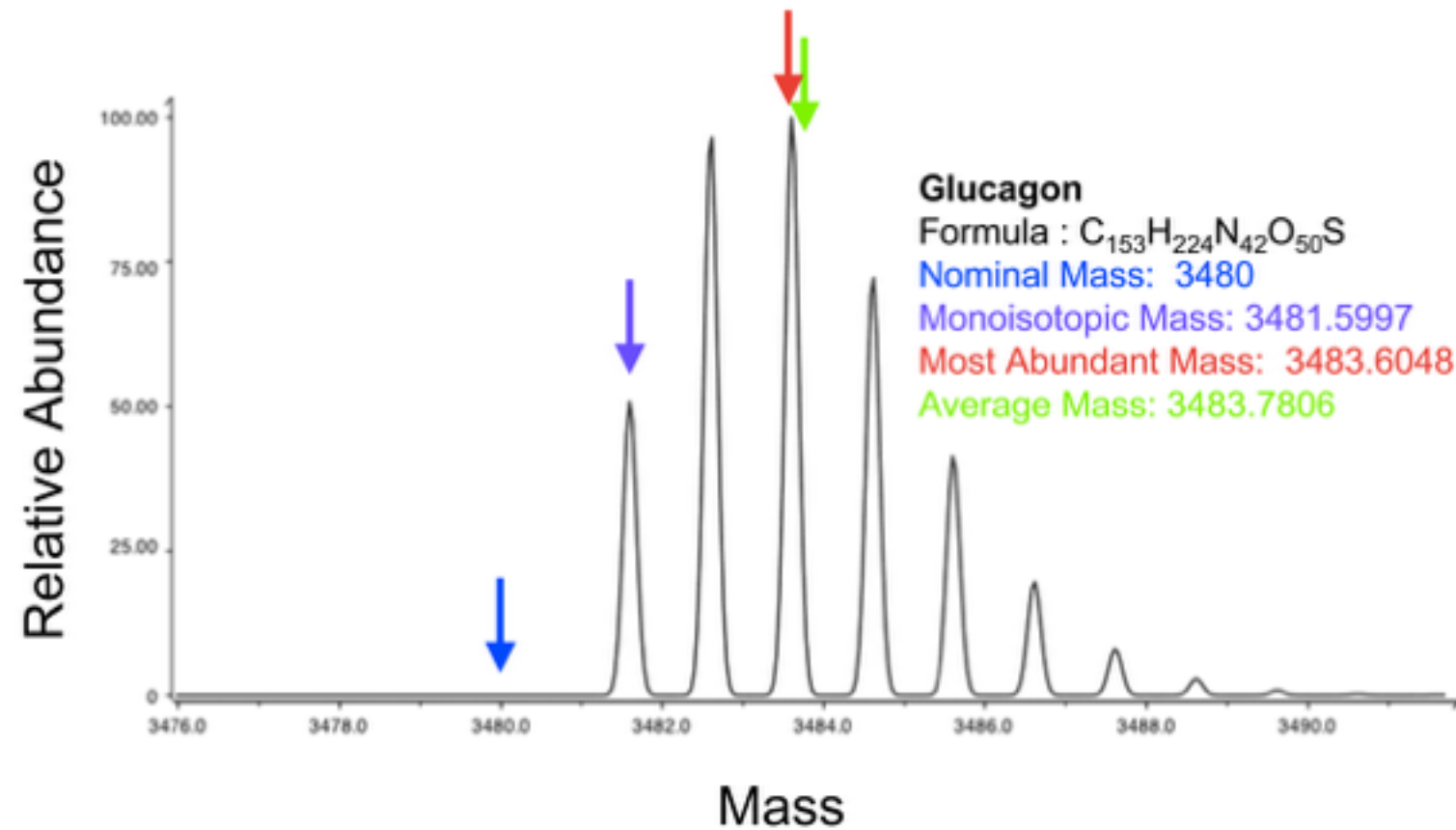
Enzyme-linked immunosorbent assay (ELISA)

# Resolving isotopes

Braykinin=RPPGFSPFR



# Different definitions of Masses



Mass Type	The sum of the molecules' atoms':
Nominal Mass	Integer mass of the most abundant isotope
Monoisotopic Mass	Masses of unbound, ground-state, isotope <i>The preferred measure for high-resolution MS</i>
Average Mass	Average Mass given the isotopes and their natural abundance <i>The preferred measure for low-resolution MS</i>

# Mass of Elements

Element	Average mass	Monoisotopic Mass
C	12.0107	12
N	14.00674	14.00307
H	1.00794	1.00782
O	15.9994	15.99492
S	32.066	31.97207
Pepide ALLETYCATPAKSE, C <sub>65</sub> H <sub>105</sub> N <sub>15</sub> O <sub>23</sub> S	1496.682	1495.72281

**monoisotopic mass** is the mass of the principal (most abundant) isotope.

**average mass** is the average mass of all isotopes, normalized for natural abundance.

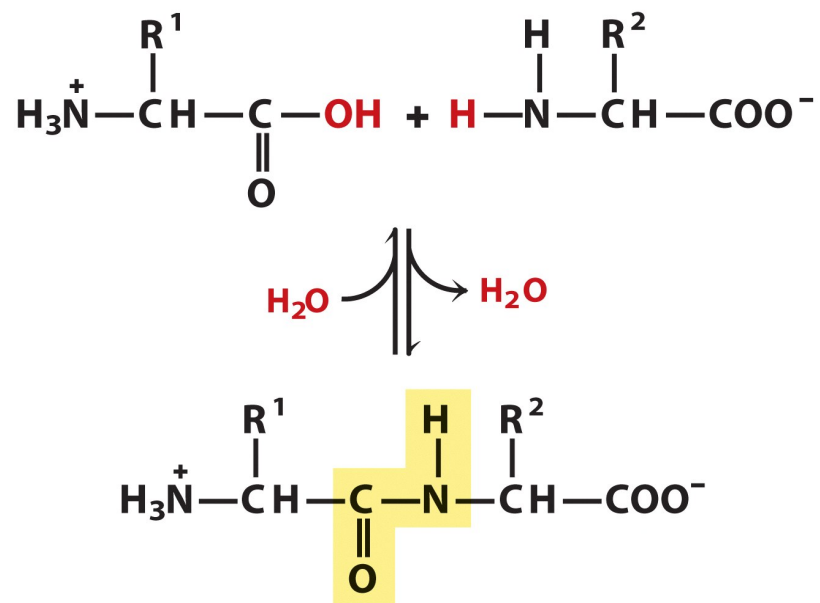
# Residue mass of amino acids

Amino Acid ↕	Short ↕	Abbrev. ↕	Formula ↕	Mon. Mass§ (Da) ↕	Avg. Mass (Da) ↕
Alanine	A	Ala	C <sub>3</sub> H <sub>5</sub> NO	71.03711	71.0788
Cysteine	C	Cys	C <sub>3</sub> H <sub>5</sub> NOS	103.00919	103.1388
Aspartic acid	D	Asp	C <sub>4</sub> H <sub>5</sub> NO <sub>3</sub>	115.02694	115.0886
Glutamic acid	E	Glu	C <sub>5</sub> H <sub>7</sub> NO <sub>3</sub>	129.04259	129.1155
Phenylalanine	F	Phe	C <sub>9</sub> H <sub>9</sub> NO	147.06841	147.1766
Glycine	G	Gly	C <sub>2</sub> H <sub>3</sub> NO	57.02146	57.0519
Histidine	H	His	C <sub>6</sub> H <sub>7</sub> N <sub>3</sub> O	137.05891	137.1411
Isoleucine	I	Ile	C <sub>6</sub> H <sub>11</sub> NO	113.08406	113.1594
Lysine	K	Lys	C <sub>6</sub> H <sub>12</sub> N <sub>2</sub> O	128.09496	128.1741
Leucine	L	Leu	C <sub>6</sub> H <sub>11</sub> NO	113.08406	113.1594
Methionine	M	Met	C <sub>5</sub> H <sub>9</sub> NOS	131.04049	131.1986
Asparagine	N	Asn	C <sub>4</sub> H <sub>6</sub> N <sub>2</sub> O <sub>2</sub>	114.04293	114.1039
Pyrrolysine	O	Pyl	C <sub>10</sub> H <sub>11</sub> N <sub>3</sub> O <sub>3</sub>	255.15820	255.2470

etc...

The free form of the amino acids are a the equivalent of a water molecule heavier (~18 Da) than its residue mass

# The mass of a peptide

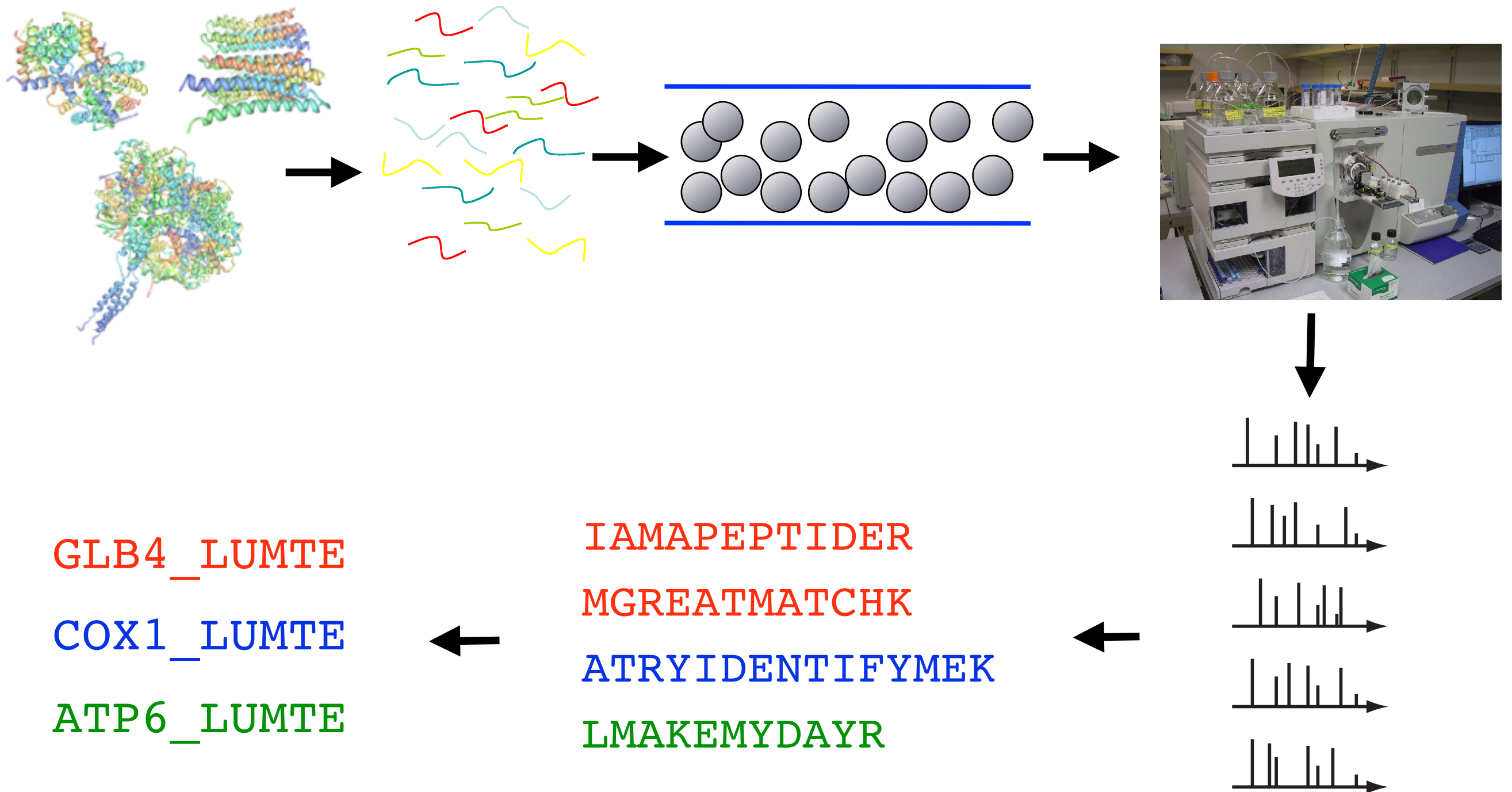


The mass  $m(p)$  of peptide  $p$  can be calculated as the residue mass of its constituent amino acids,  $a_1...a_n$ , and the mass of a water molecule

$$m(p) = m(H_2O) + \sum_{i=1...n} m(a_i)$$

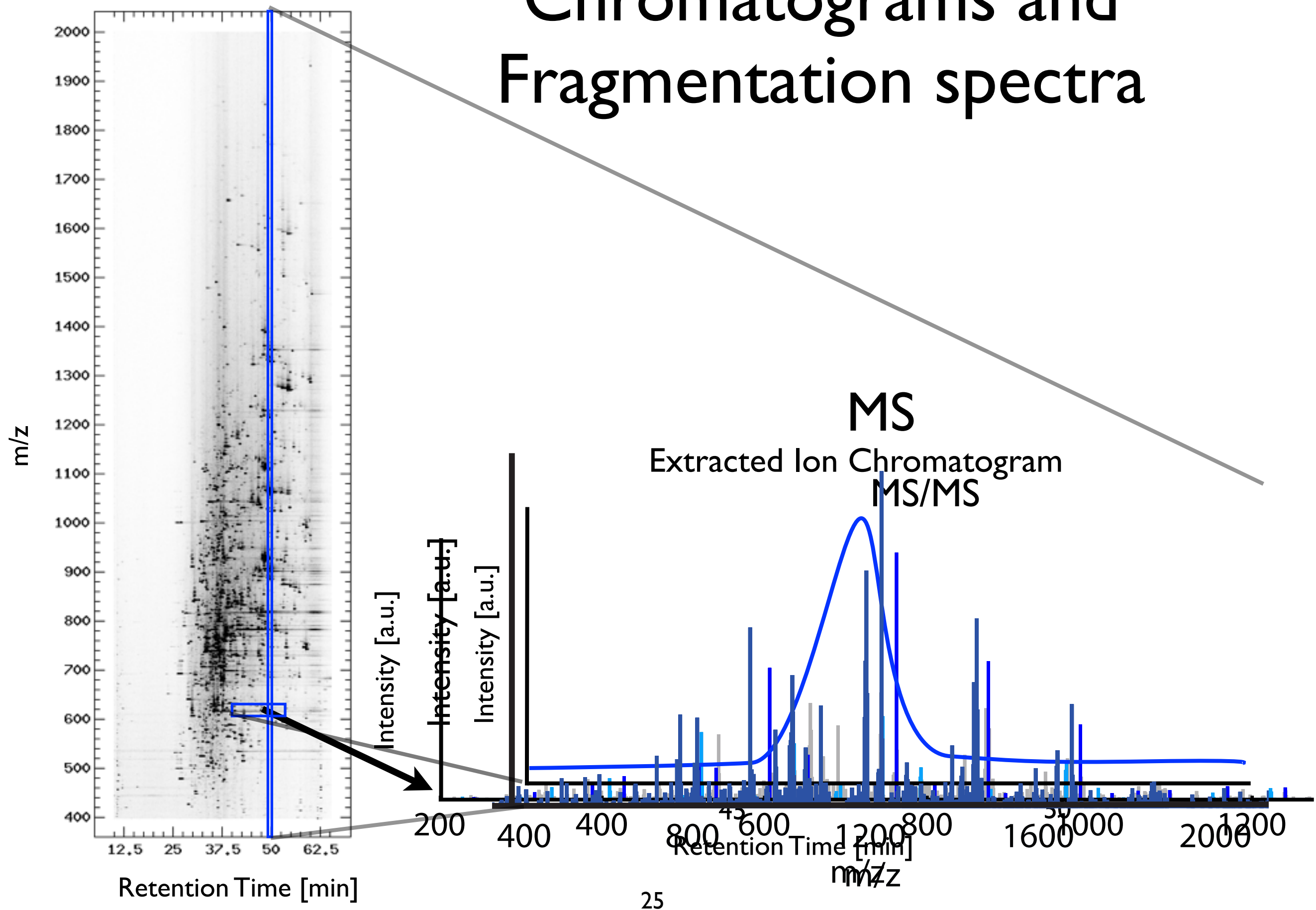


# Shotgun proteomics



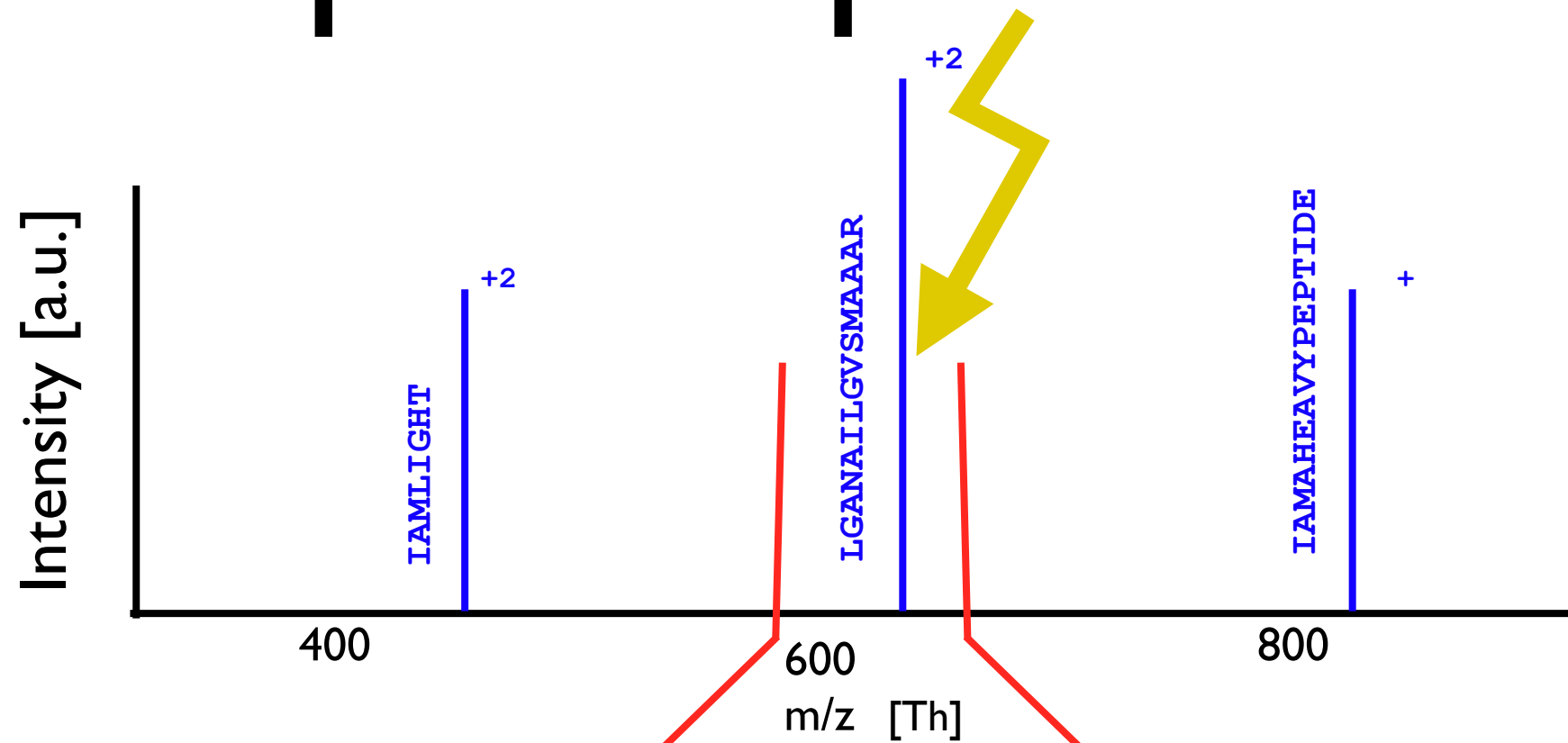


# Chromatograms and Fragmentation spectra

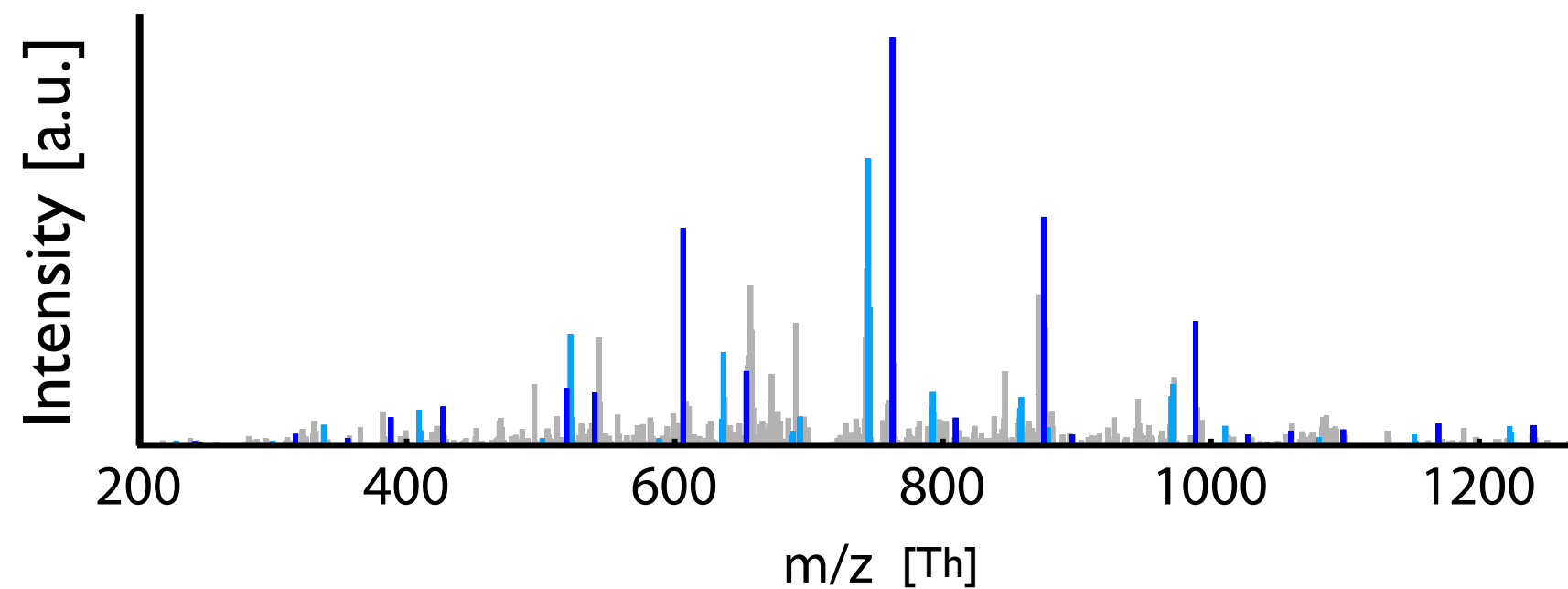


# Peptide spectra

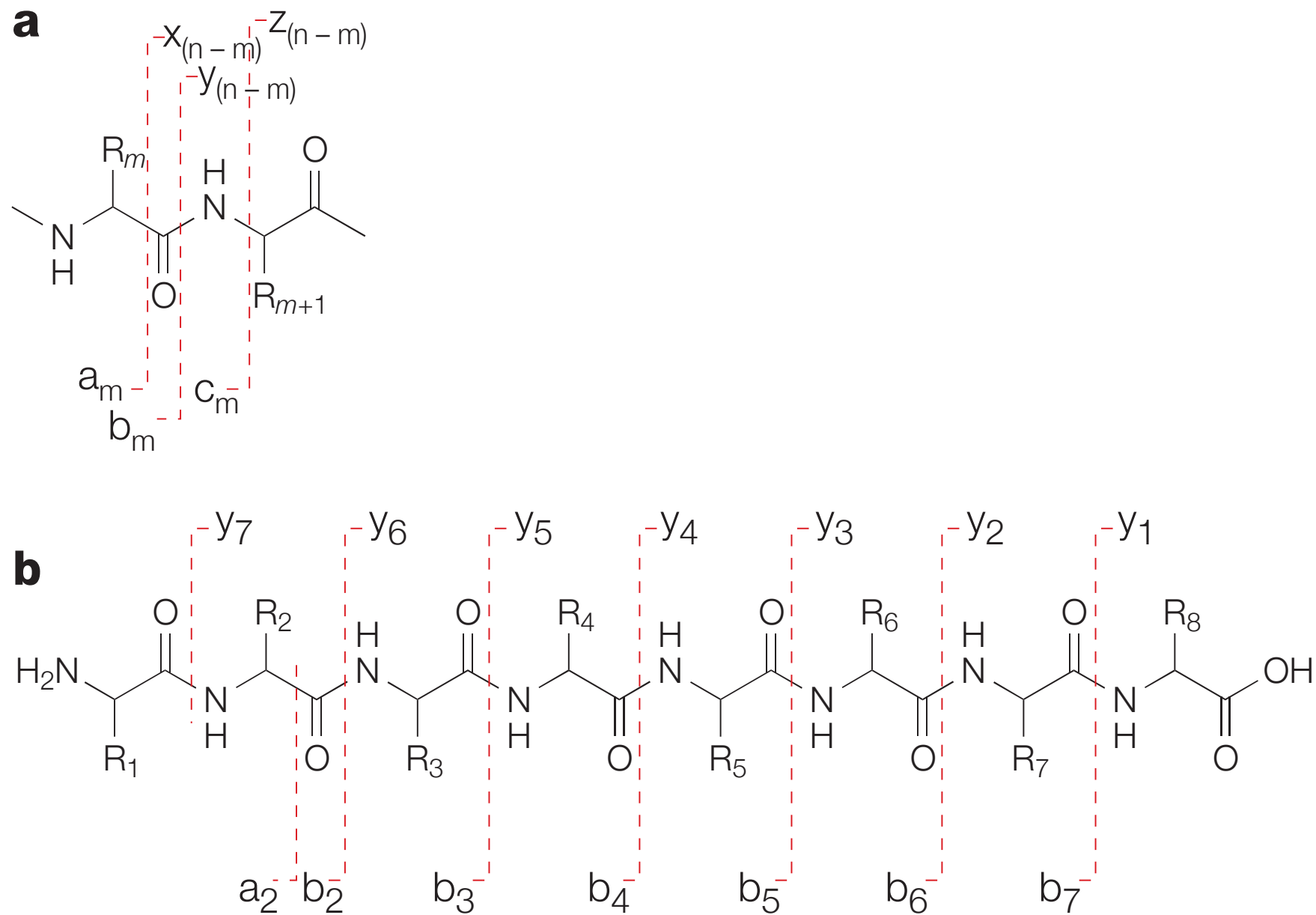
MS<sup>1</sup>



MS<sup>2</sup>



# Peptide Fragmentation

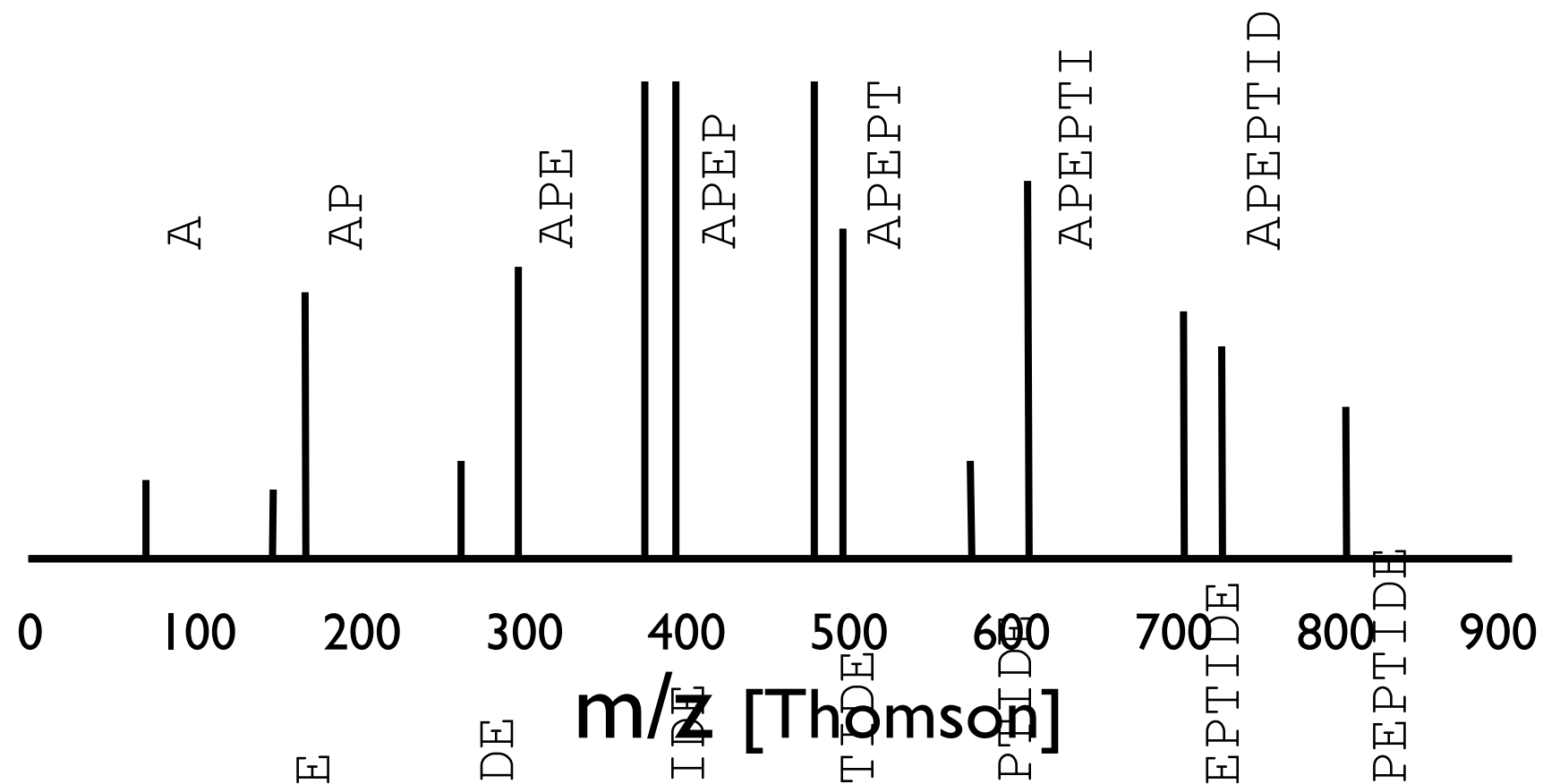


# Fragmentation Spectrum

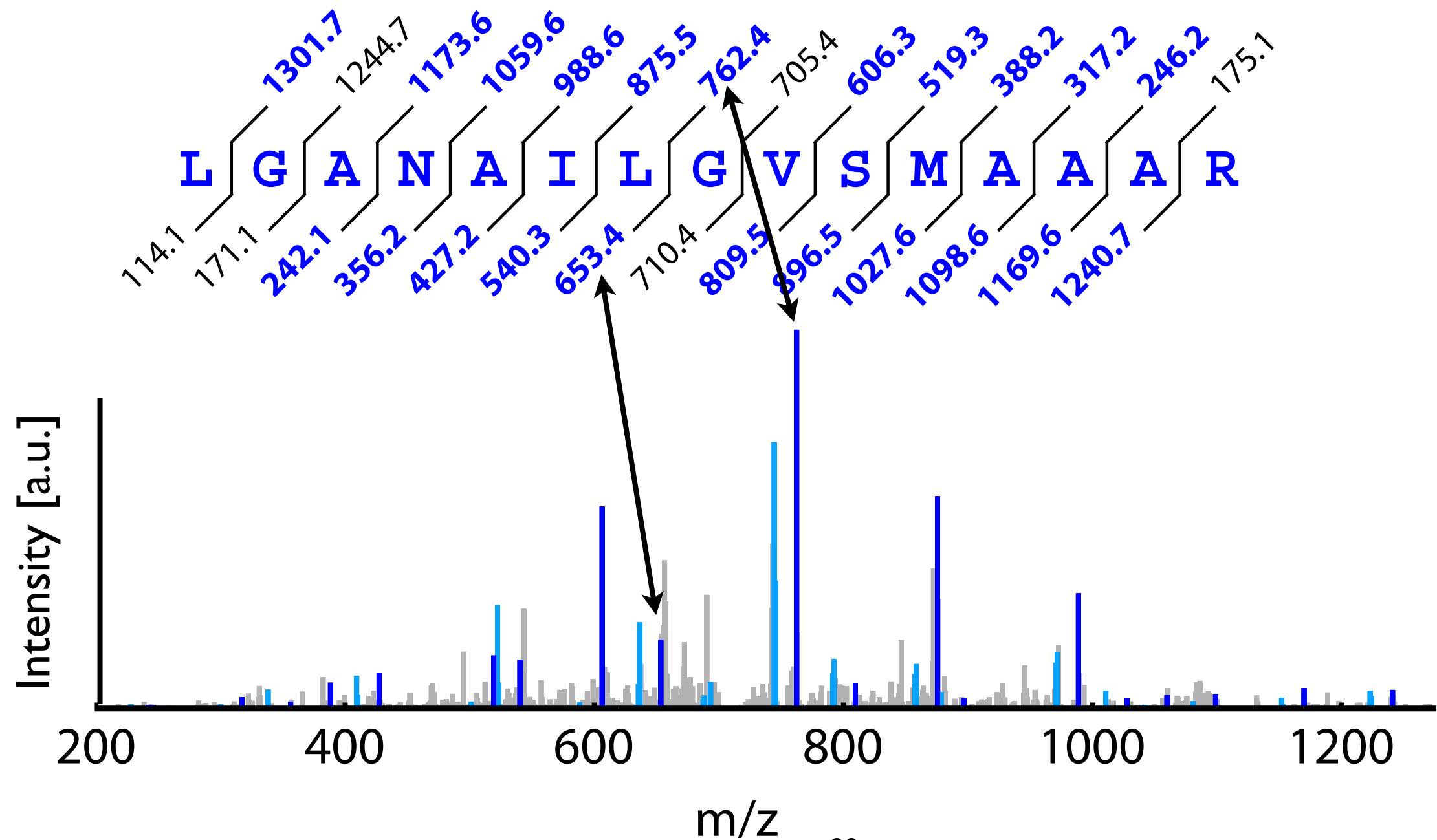
A|P|E|P|T|I|D|E

b:

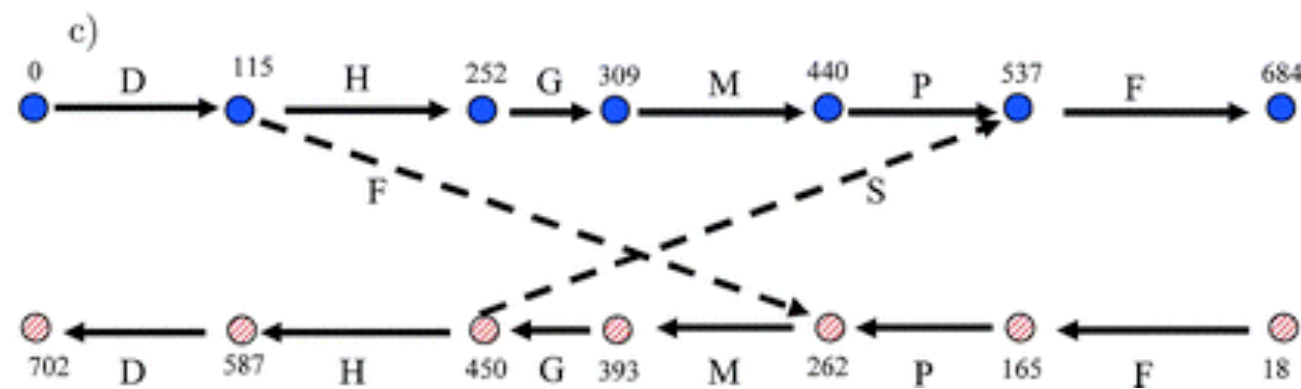
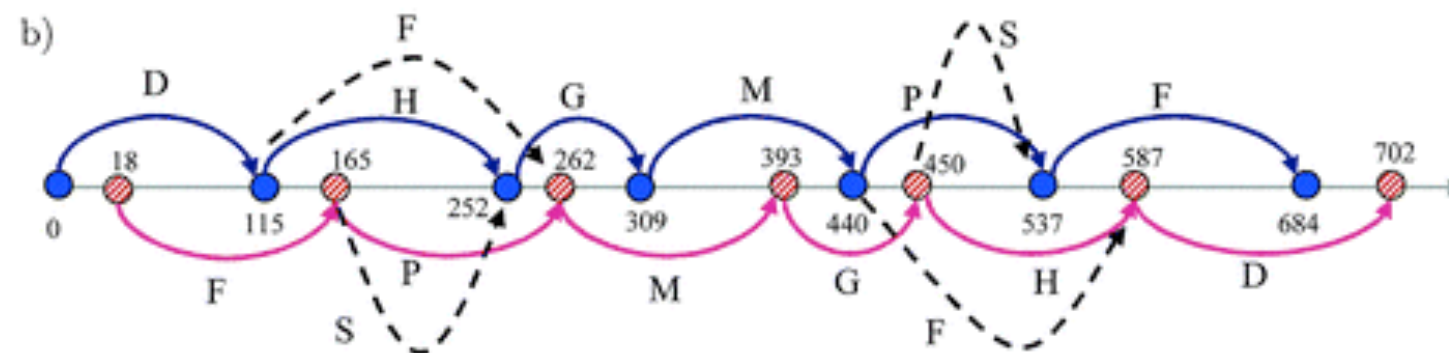
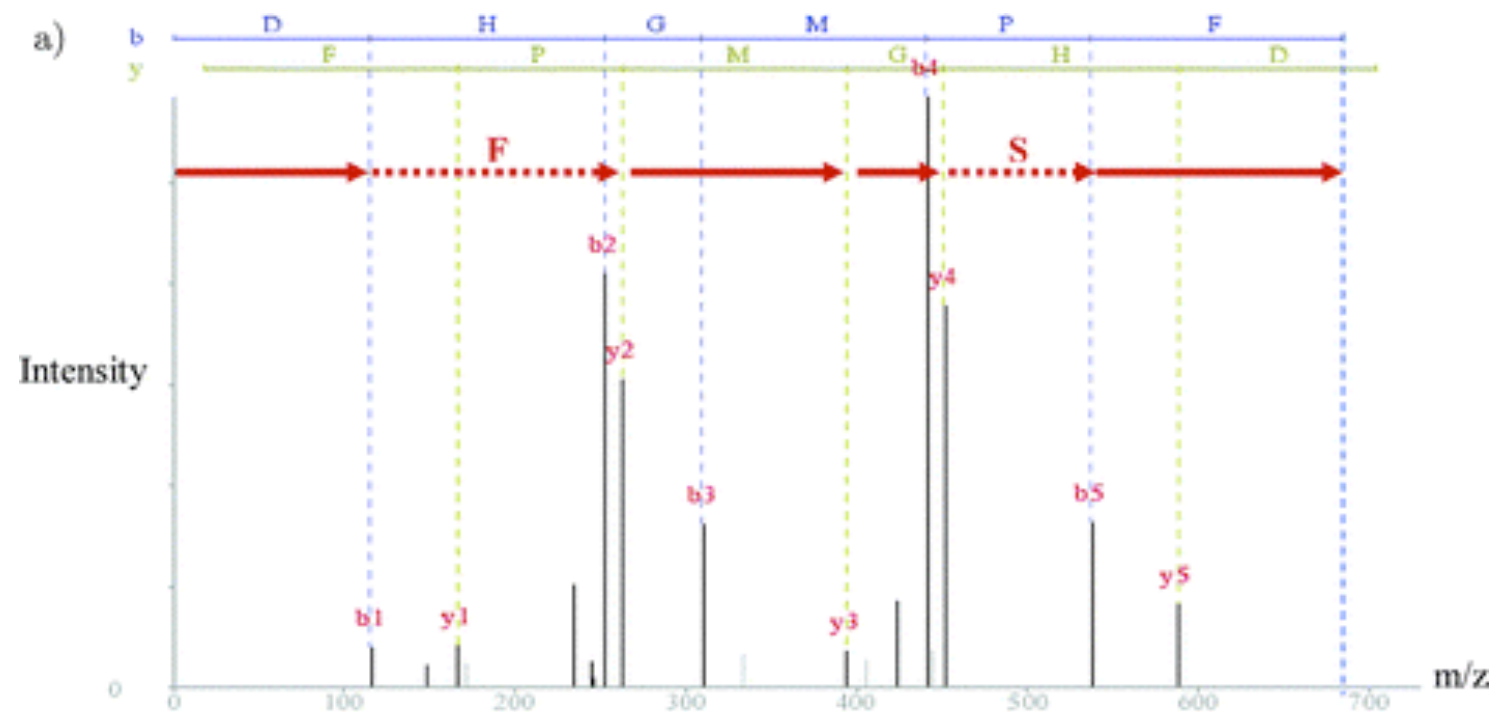
y:



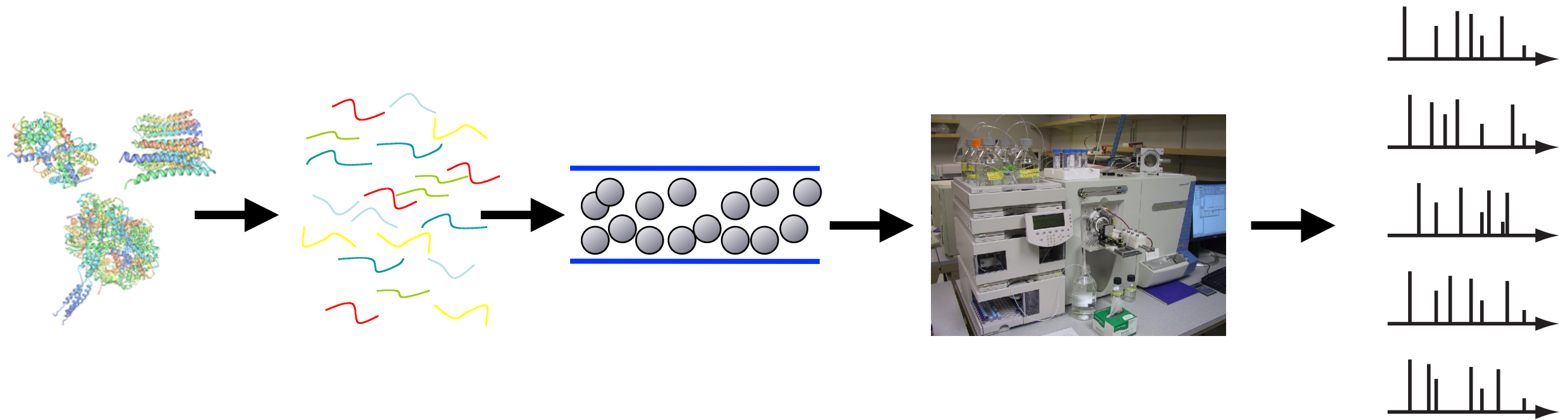
# Peptide fragmentation spectrum



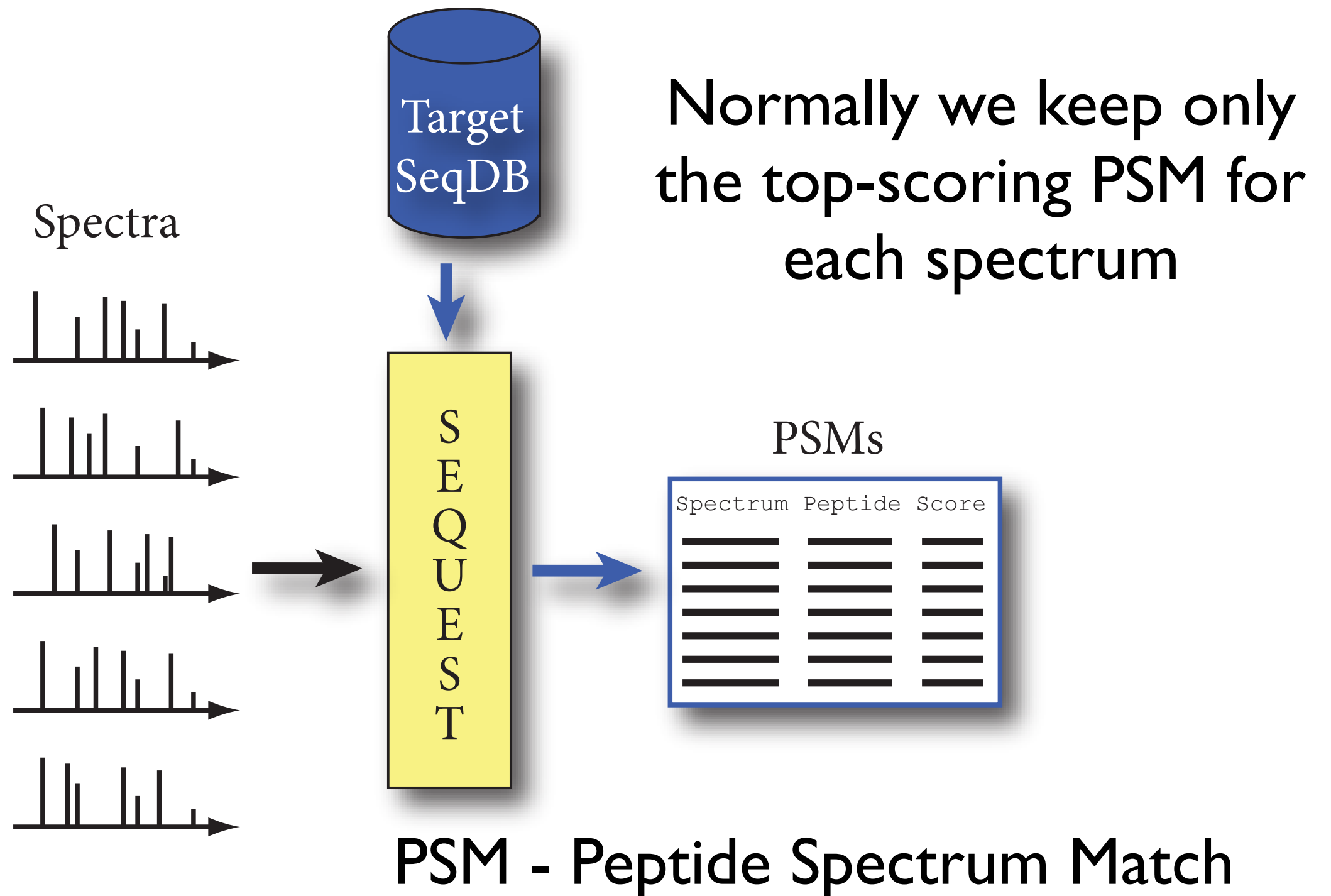
# de Novo sequencing



# Peptide Identification



# Peptide identification





# Four popular search engines

- SEQUEST (Scripps, Thermo Fisher Scientific)  
<http://fields.scripps.edu/sequest>
- MASCOT (Matrix Science)  
<http://www.matrixscience.com>
- X! Tandem (The Global Proteome Machine Organization)  
<http://www.thegpm.org/TANDEM>
- MS-GFDB  
<http://proteomics.ucsd.edu/Software/MSGFDB.html>

# Sequest

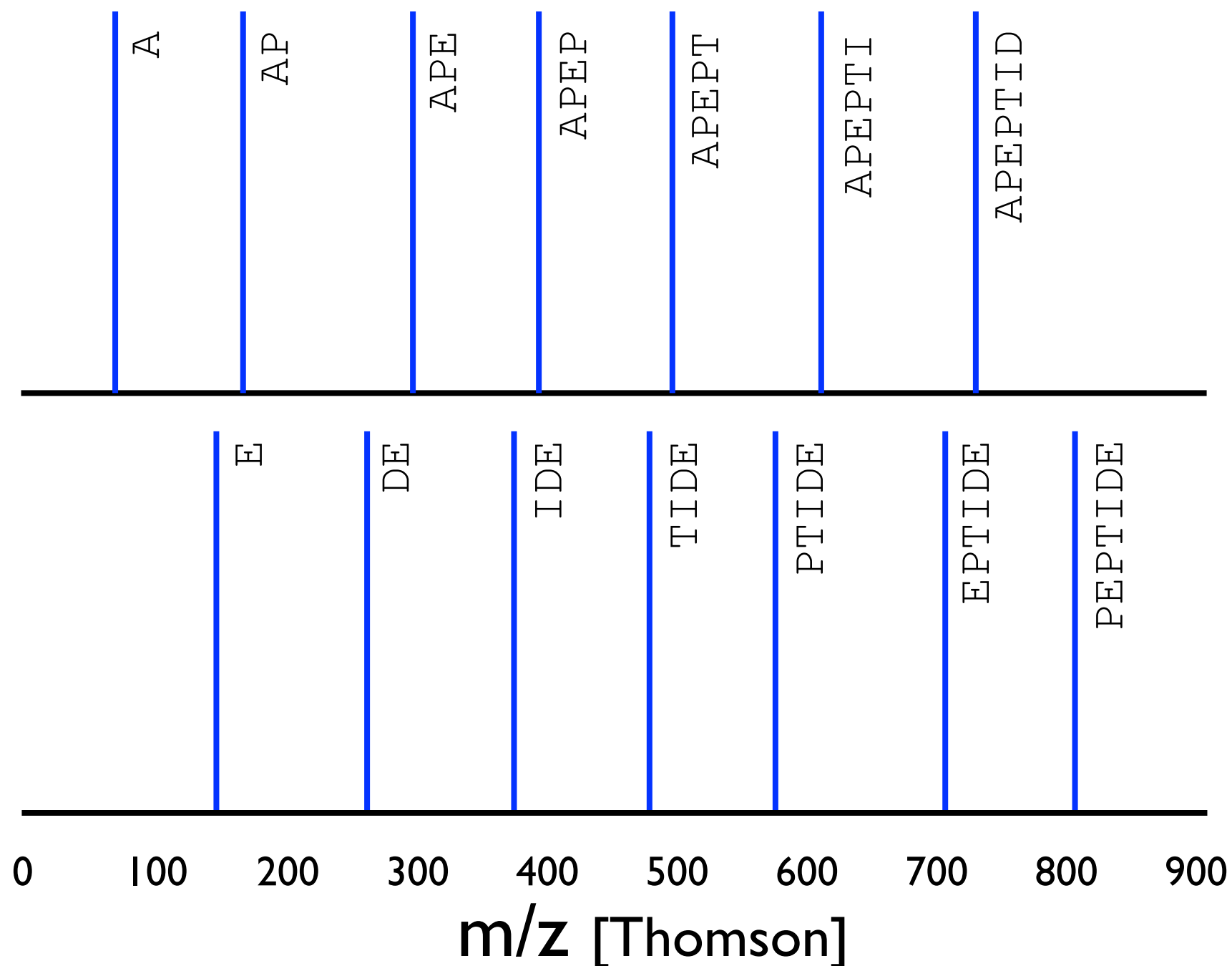
- First published automated spectral search engine
- Published but patented algorithm [Eng *et al.* JASMS 1994]
- For each spectrum  $x$  the top 500 candidate peptides are selected by a fast calculated preliminary score  $Sp$ .
- The theoretical spectra  $y$  are calculated for each of these top candidates and a background normalized cross correlation score,  $Xcorr$ , is calculated
$$R_i = \sum_{j=1}^n x_j y_{j+i}$$
$$X = R_0 - \frac{1}{151} \sum_{i=-75}^{75} R_i$$
- A score  $\Delta Cn$  is provided which gives the relative difference between the first and second best  $Xcorr$
- Re-implementations free for Academic users: Crux and Tide

# Theoretical Spectrum of a peptide

A|P|E|P|T|I|D|E

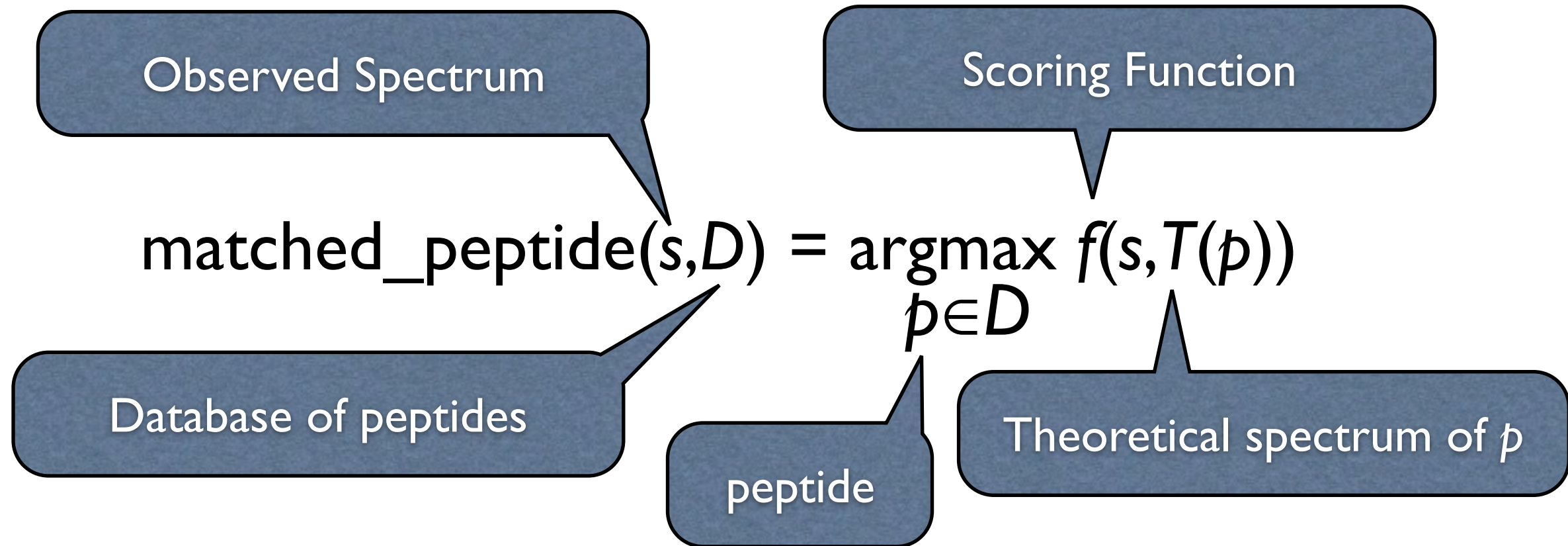
b:

y:



# Search engine

SEQUENT:



other:

$$\text{matched\_peptide}(s, D) = \underset{p \in D}{\operatorname{argmax}} f(s, p)$$

# Mascot



- Probably the most spread commercial spectral search engine
- Unpublished scoring function (Trade secret),
- Reports Rank, score and E-value for each PSM
- Predicts a *homology threshold* from database size and instrument accuracy which each PSM should pass
- Provides a fancy web report

# X! Tandem



- Open source, published algorithm  
[Craig & Beavis Rapid Commun. Mass Spectrom 2003]
- Scoring function, HyperScore, is build around the hypergeometric distribution (of number of matched b- and y-ions)
- Provides hyperscore and E-value for each PSM
- Relatively fast

$$\text{HyperScore} = \left( \sum_{i=0}^n I_i * P_i \right) * N_b! * N_y!$$

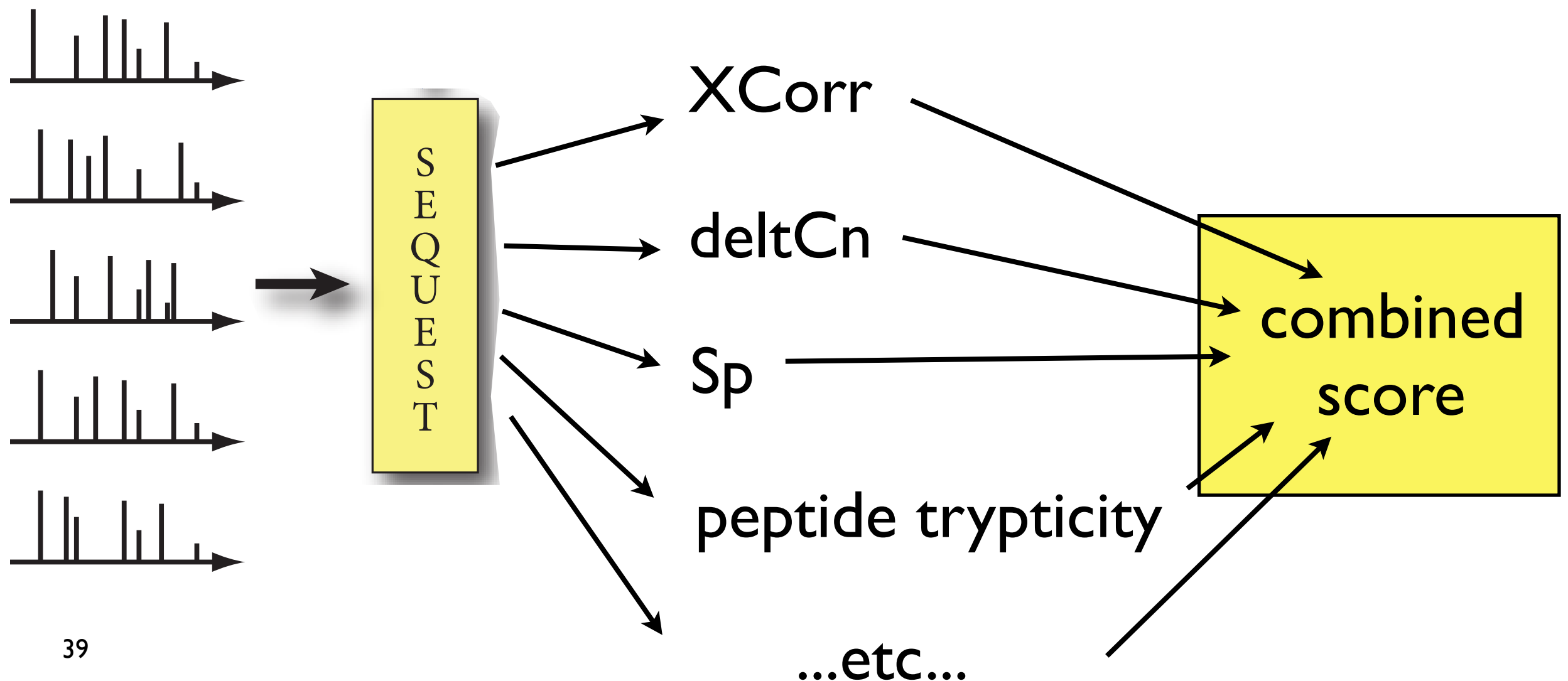
Intensity

Present 0/1

# Post Processors

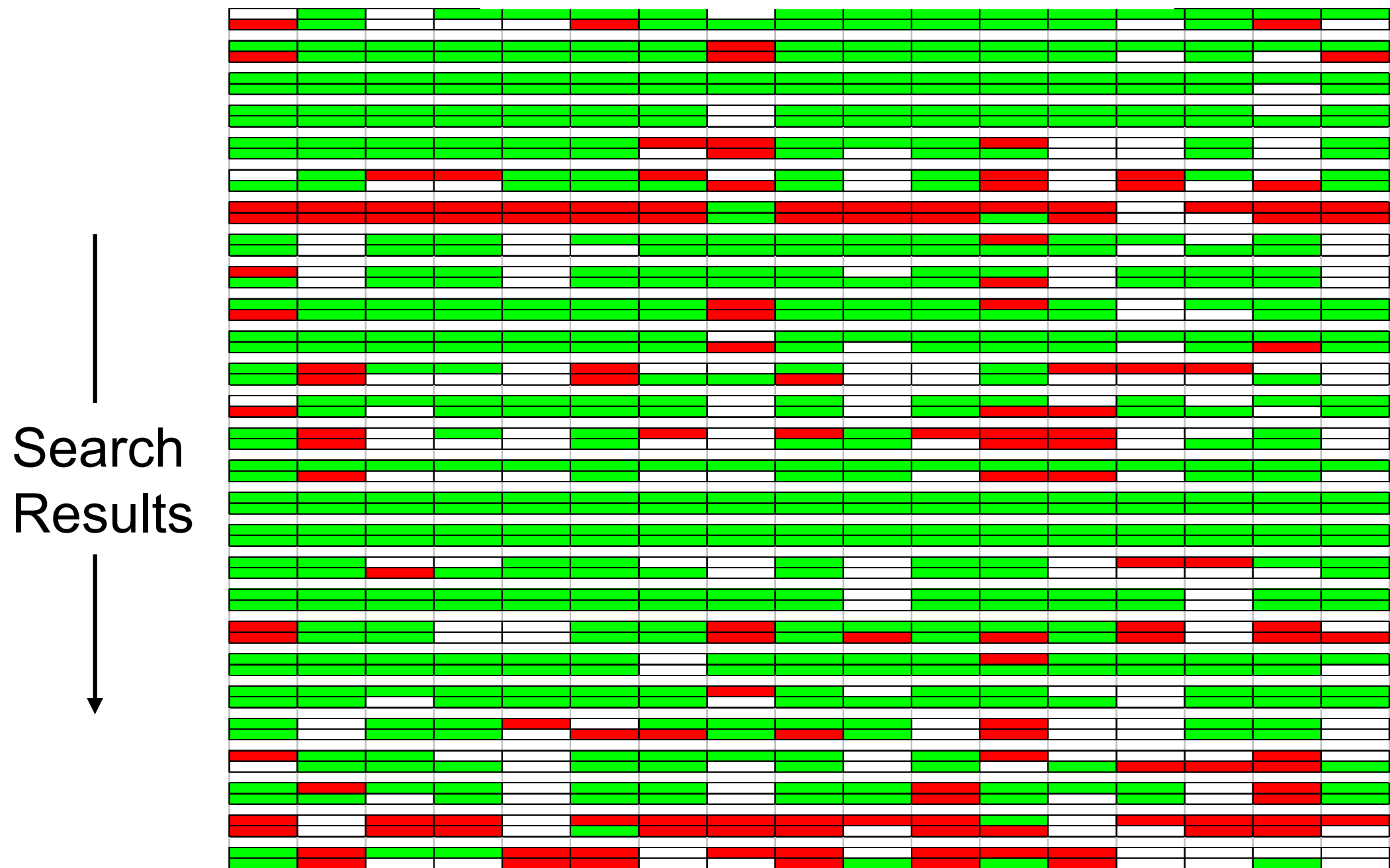
Combinations of scores are known to give better yield than individual scores. Two examples are:

- PeptideProphet (LDA) [Keller *et al.* 2002 Anal Chem]
- Percolator (semi-supervised SVM) [Käll *et al.* Nat Methods 2007]



# (Un)reliability of manual validation

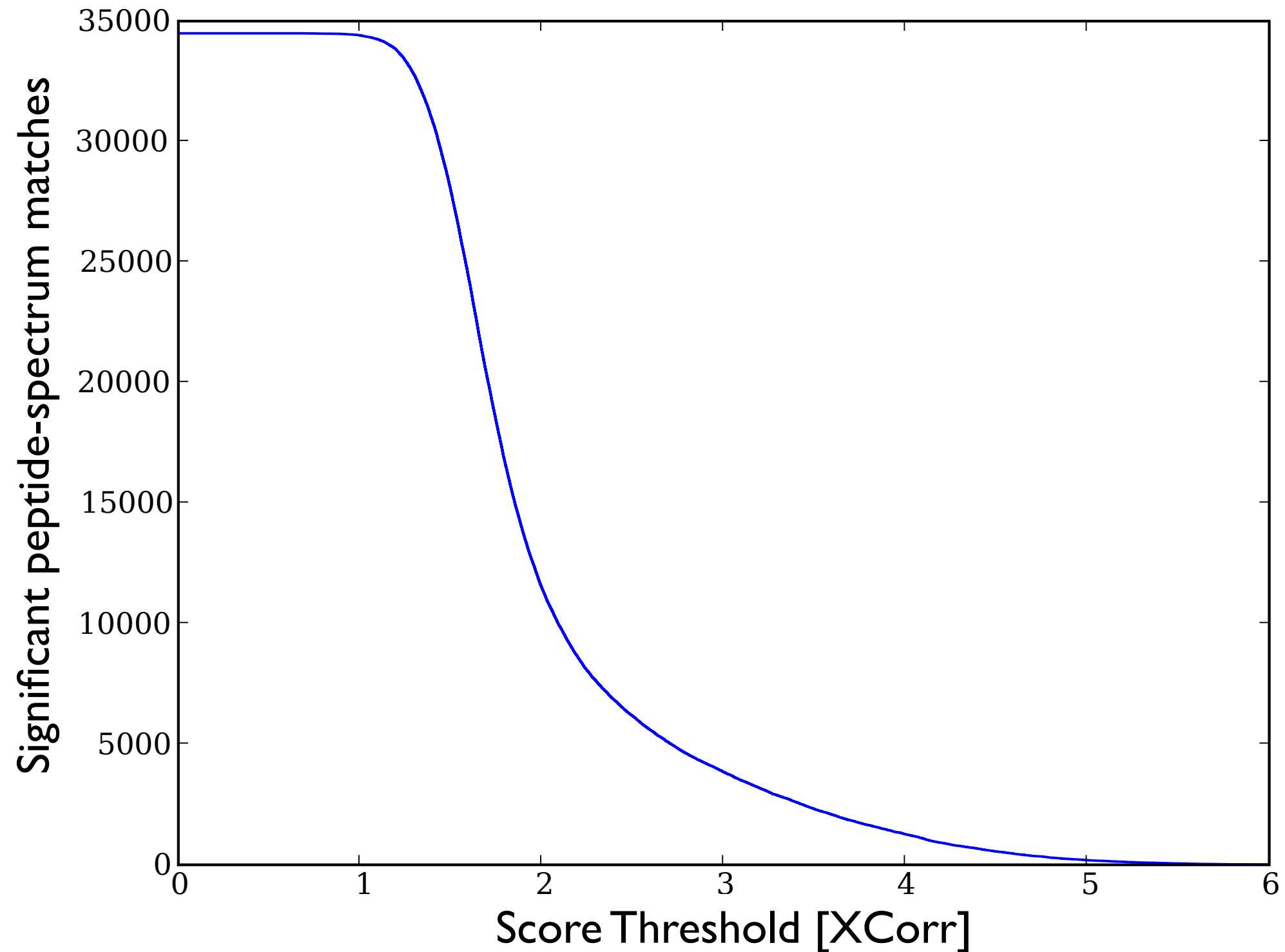
Manual Authenticators →



Correct Validation    Incorrect Validation    Validation Withheld



# Score thresholds



# correct/incorrect target PSMs

score	type
7.5	correct
7.2	correct
6.9	correct
6.8	correct
6.7	incorrect
6.5	correct
6.4	correct
6.4	correct
6.3	incorrect
6.1	correct
6	incorrect
5.9	correct
5.7	incorrect
...	...

$$\frac{2}{10}$$

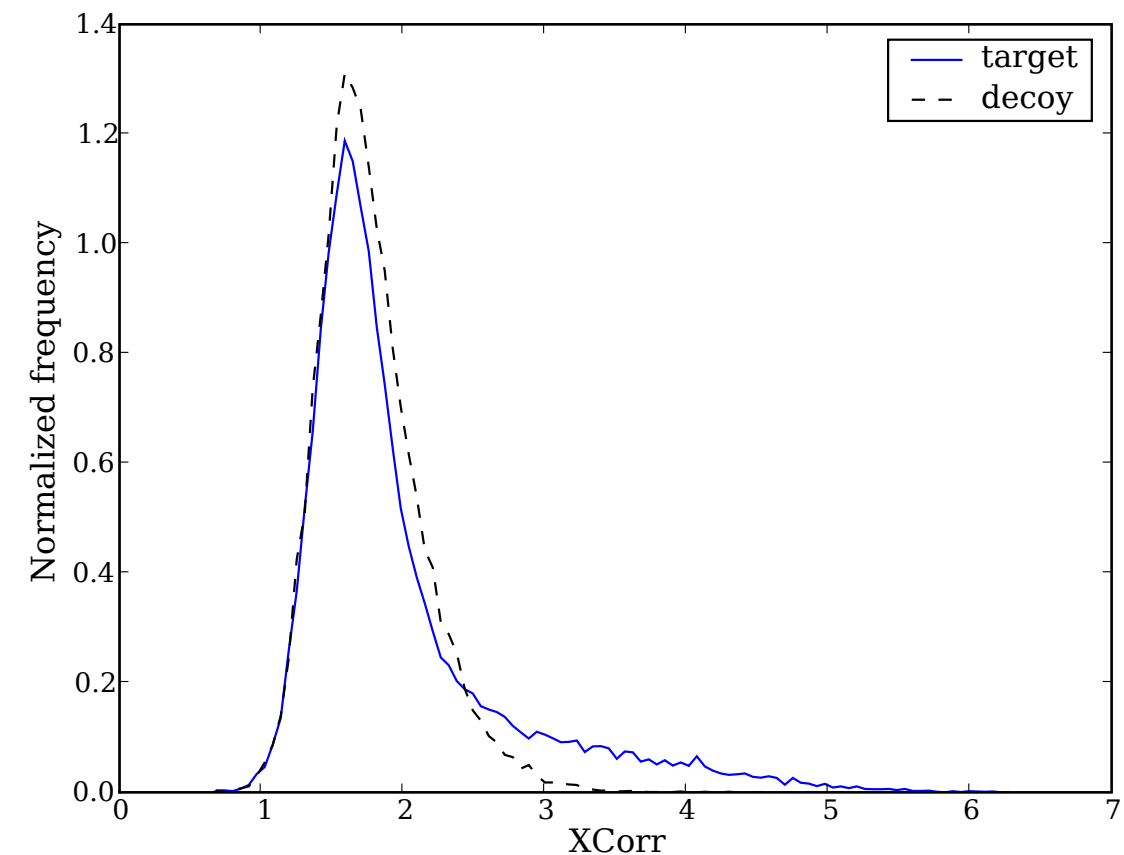
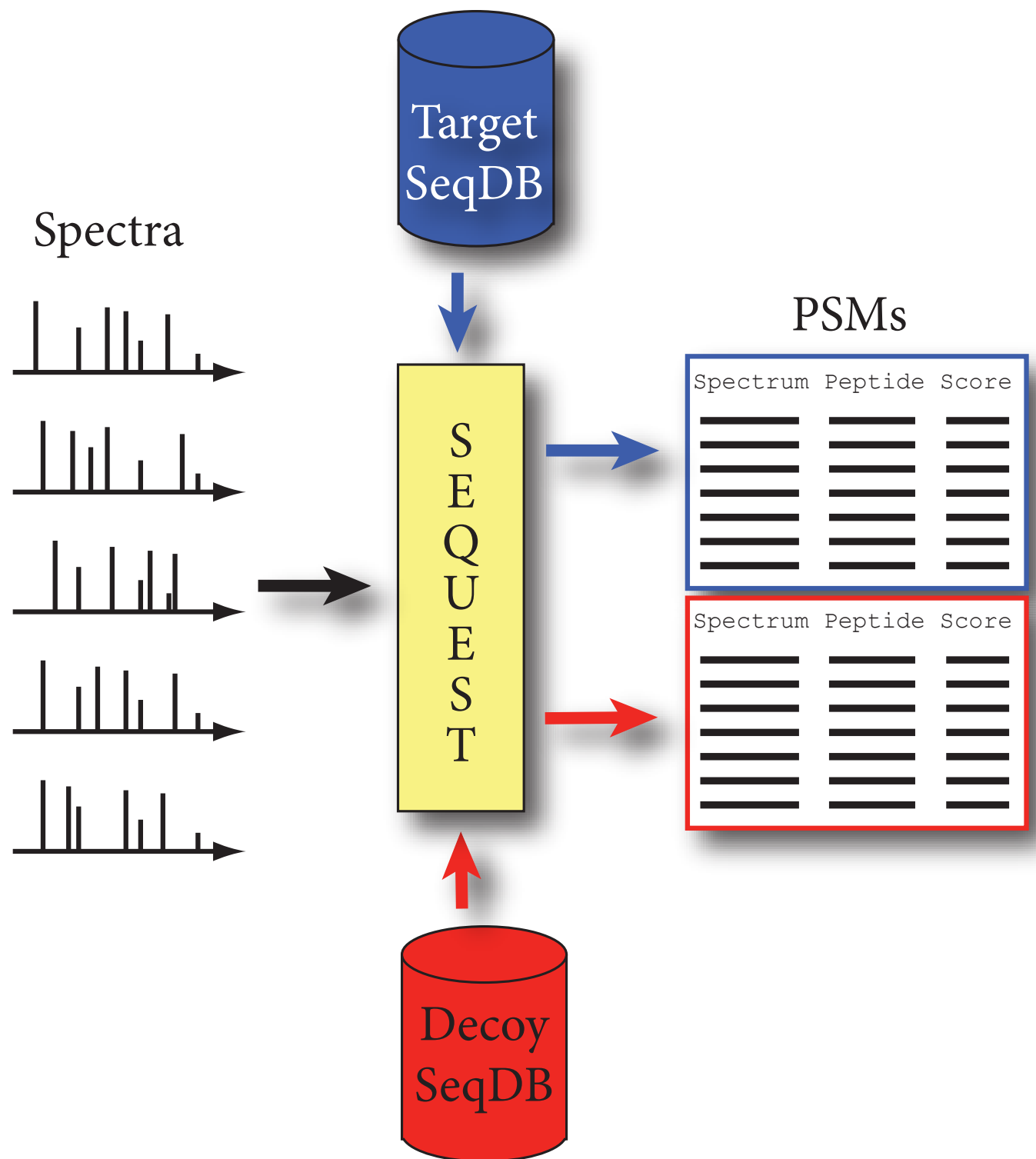
threshold

$FDR(x)$  is the expectation value of the fraction of PSMs above threshold  $x$  that are incorrect

# control for ...

- ... FDR or q value when you are interested in identifying a set of PSMs
- ... PEP when you are interested in assessing the quality of a particular PSM.
- ... p or E value in an experiment rendering one single spectrum.

# Target-decoy analysis



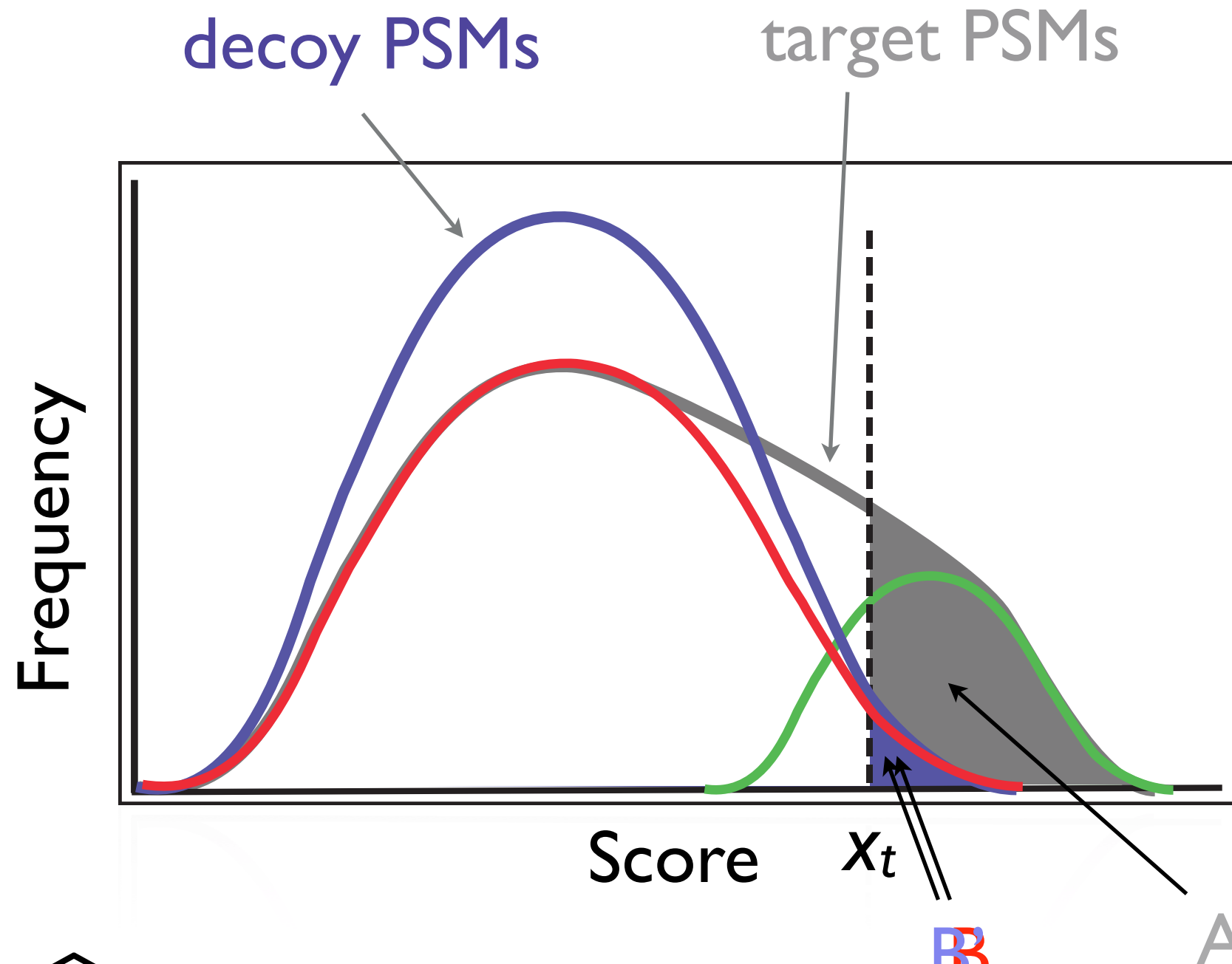
[Moore et al. JASMS 2002]

# Methods to generate decoy sequences

1. Shuffling sequences [Klammer *et al.* JPR 2006]
2. Markov models [Colinge *et al.* Proteomics 2003]
3. Reversing sequences [Moore *et al.* JASMS 2002]
4. Pseudo-Reversing sequences  
[Elias&Gygi *NMeth* 2007]

Its essential that the decoy PSMs are good proxies for incorrect target PSMs, which makes the first two methods less suitable

# Using decoy PSMs to estimate false discovery rate



$$\text{FDR}(x_t) = \frac{\Pr(x \geq x_t, H=0)}{\Pr(x \geq x_t)}$$

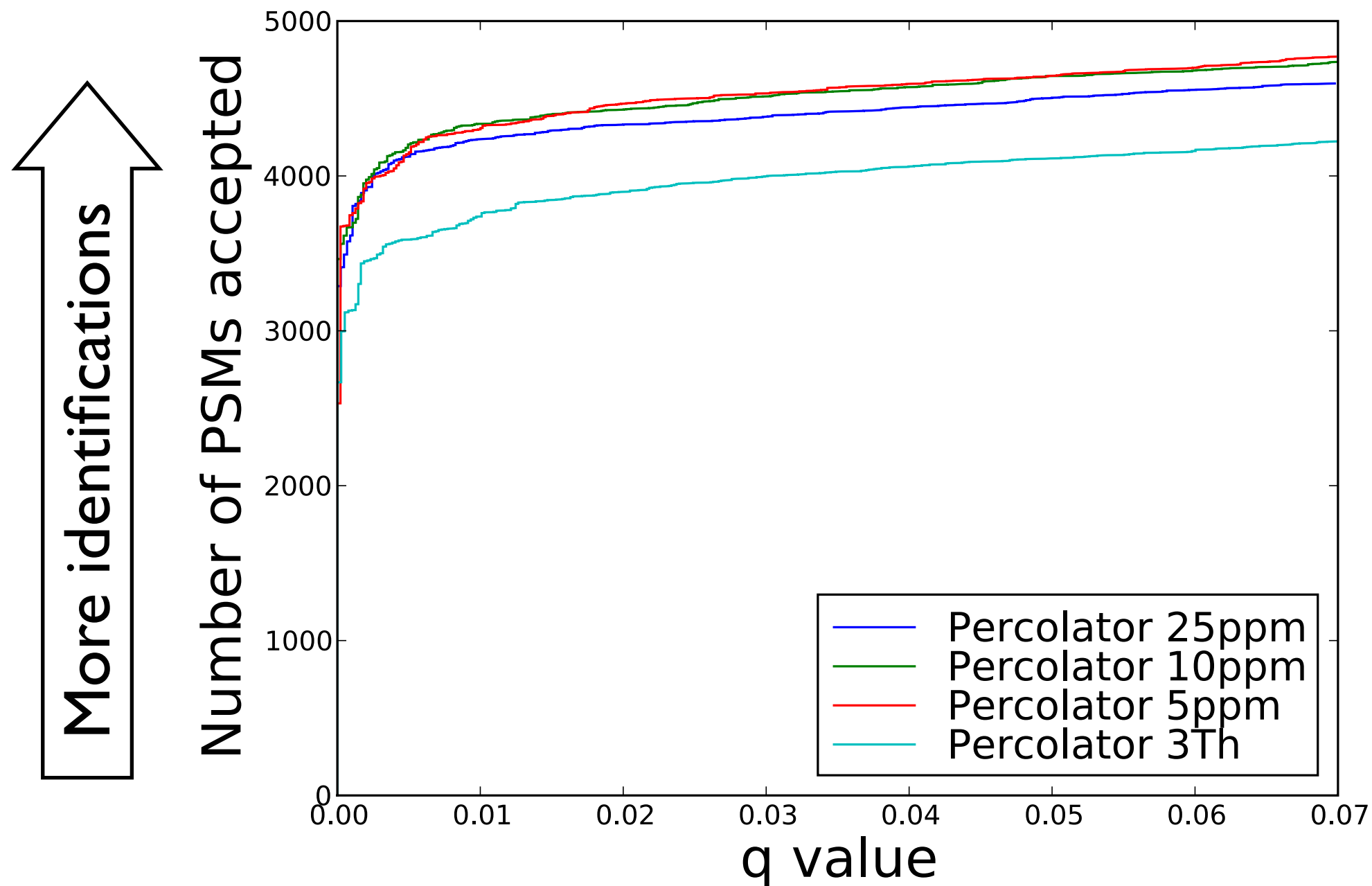
$$\widehat{\text{FDR}} = \frac{B}{A} = \frac{\widehat{\pi_0} B'}{A}$$

$$\widehat{q}(x_t) = \inf_{x \leq x_t} \{\widehat{\text{FDR}}(x)\}$$

[Käll et al. JPR 2008]

$\widehat{\pi_0}$  is the prior probability that a target PSM is incorrectly matched

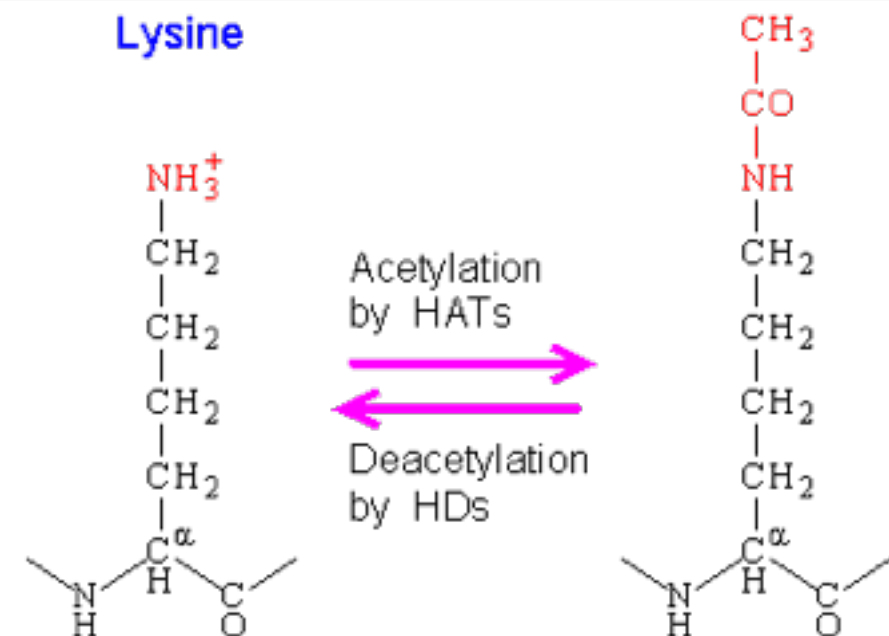
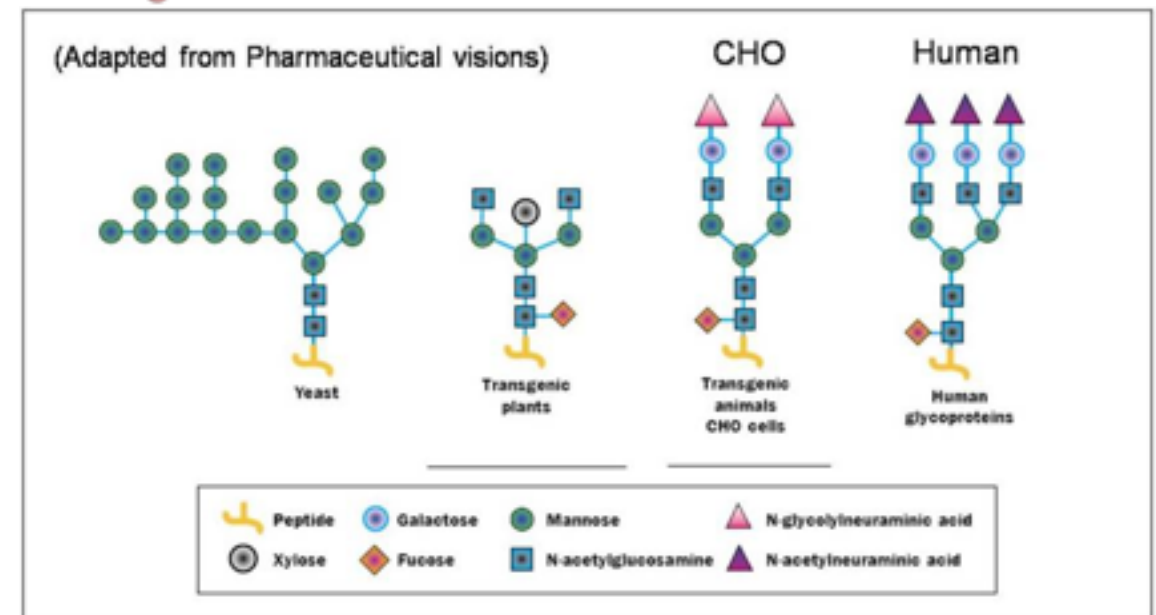
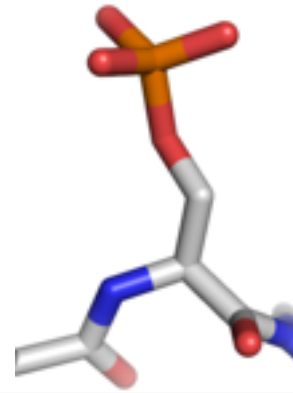
# q-p plot



Higher fraction incorrect identifications

# Some common PTMs

- Phosphorylations
  - Phosphate attached to serine or threonine
- Glycosylations
  - Glycans attached to a nitrogen (N-linked) of asparagine or arginine side-chains
- Acetylations
  - Acetyl group attached to lysine or N-terminus

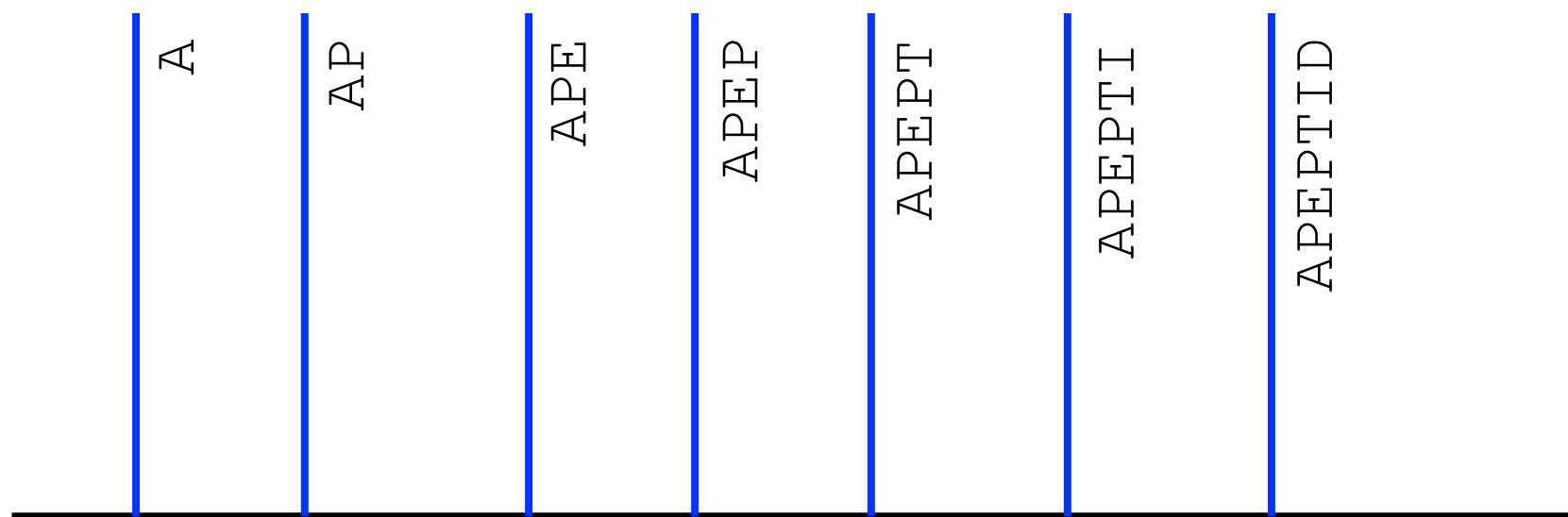




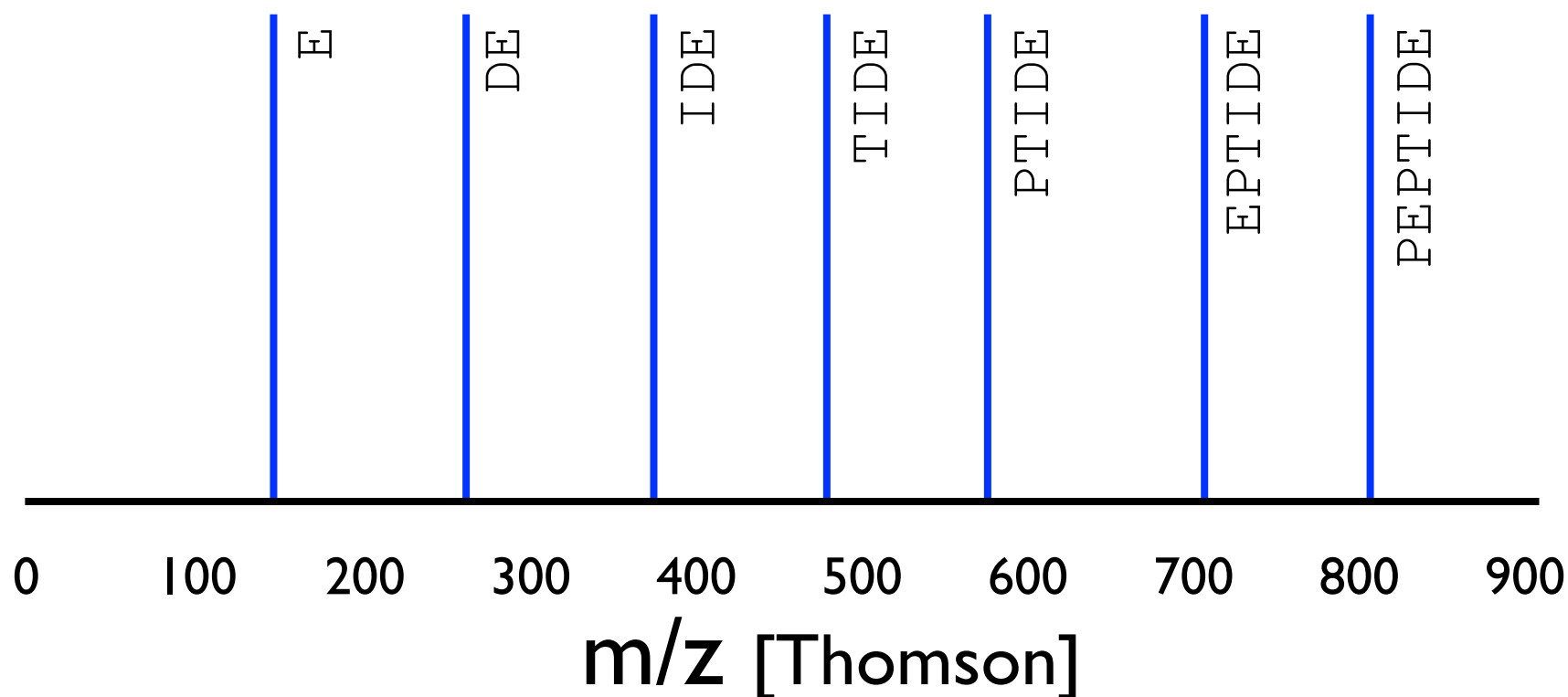
# Theoretical Spectrum of a peptide

A|P|E|P|T|I|D|E

b:



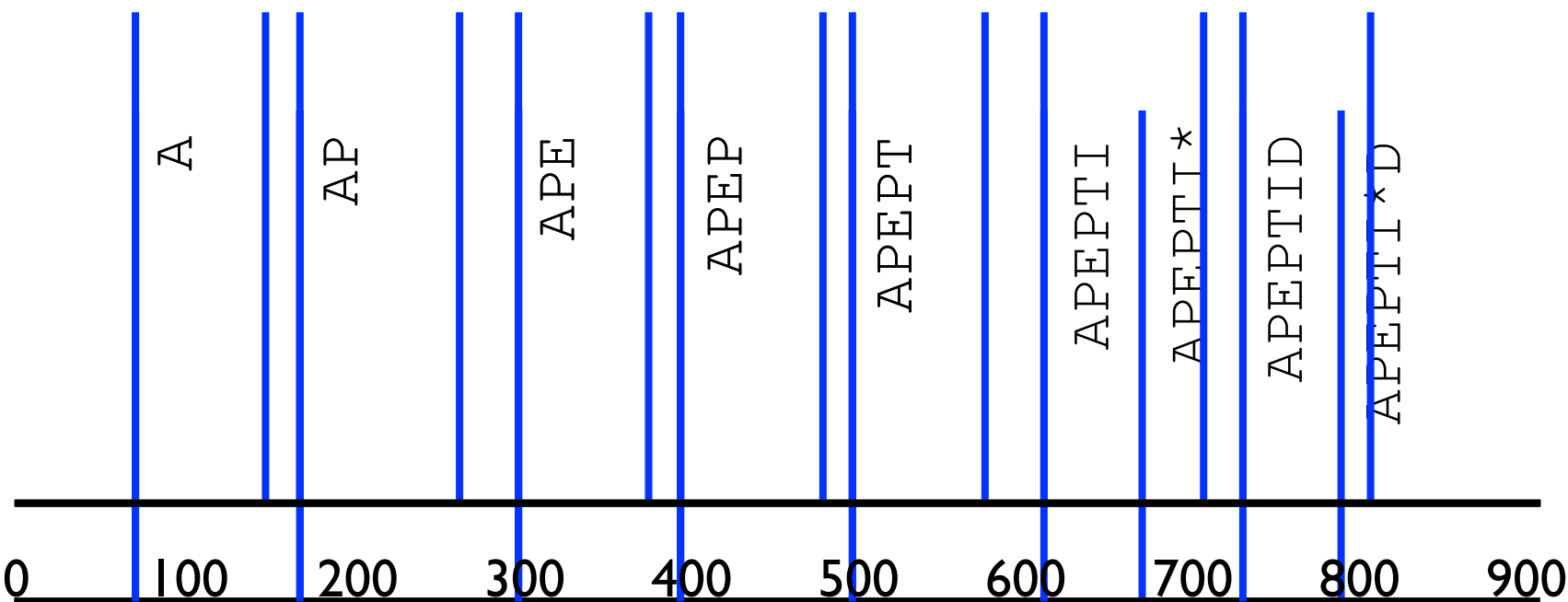
y:



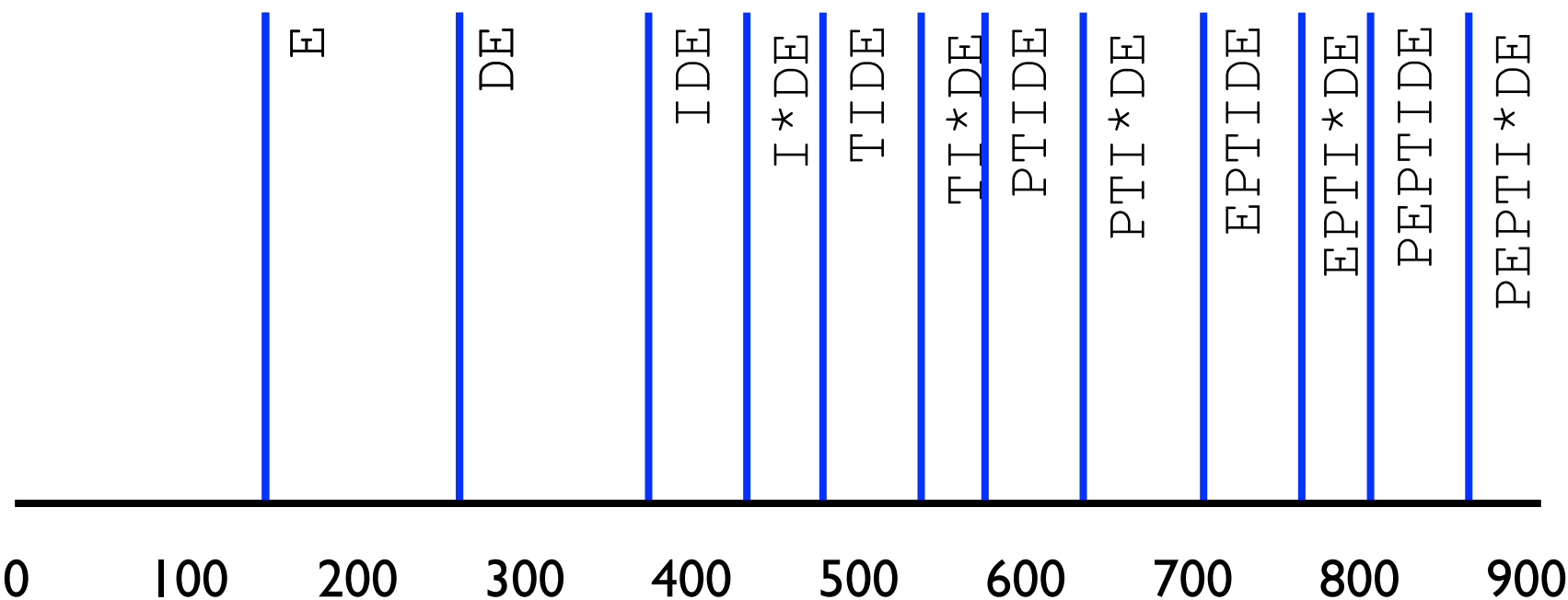
# Theoretical Spectrum of a PTM Peptide

A|P|E|P|T|I\*|D|E

Unmodified:  
b:

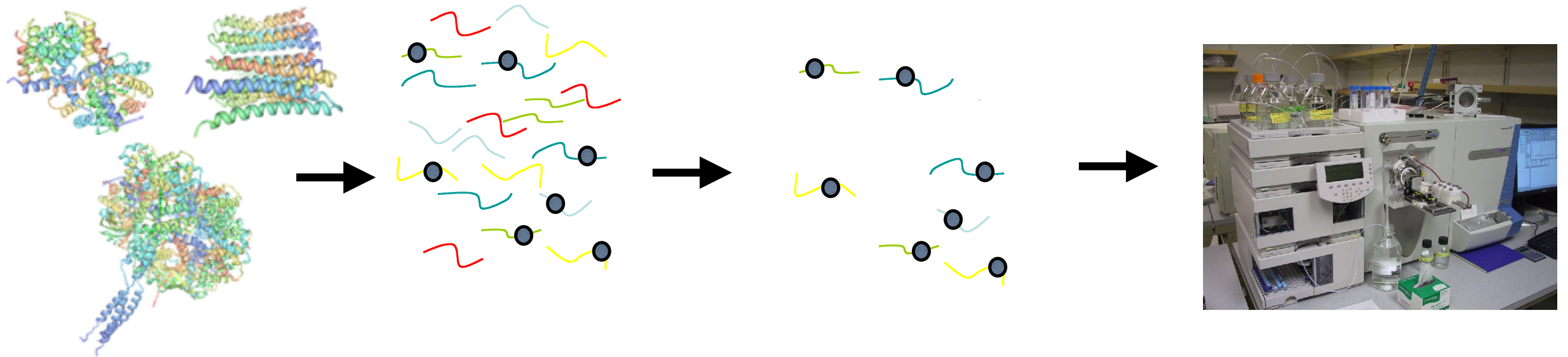


Modified:  
y:



m/z [Thomson]

# Identification Post-Translational Modifications



- Large-scale identification of PTMs normally involve digestion, PTM enrichment and identification by MS/MS

# File format for spectral data



Spectral data

Peptide Spectrum Matches

- XML-based: .mzXML, .mzData
- XML-based: pepXML, mzIdentML
- tab delimited: .ms2
- tab delimited: .sqt

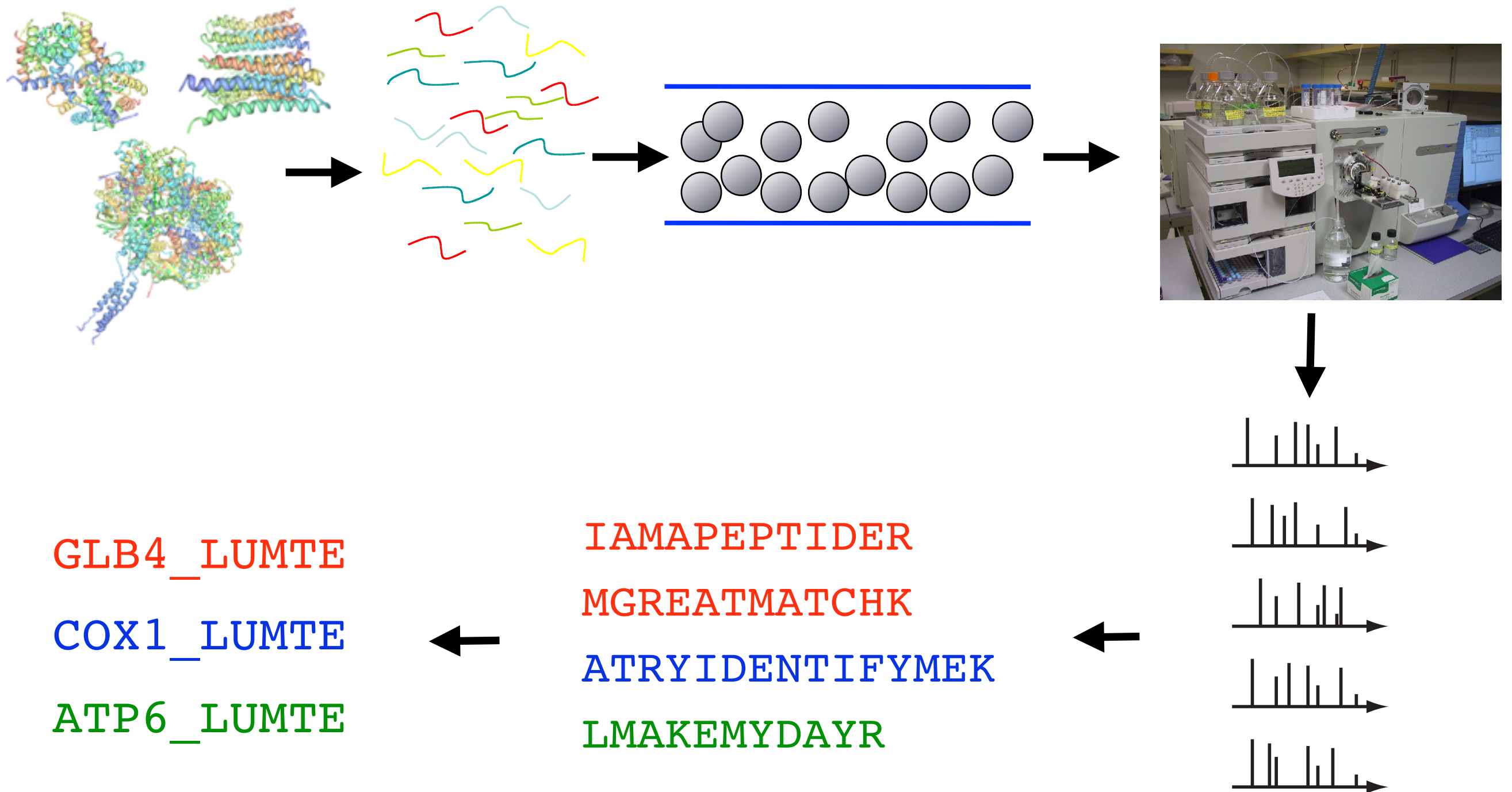
S	45894	45894	2	1	maccoss007	2038.59	9199.5	147.1	153628	
M	1	27	2040.244	0.0000	1.5881	245.6	11 34		V.YKCAADKQDATVVELTNL.T	U
L	YCR102C									
M	2	68	2038.265	0.0116	1.5698	208.4	11 36		S.TQSGIVAEQALLHSLNENL.S	U
L	YGR080W									
M	3	34	2039.247	0.1582	1.3369	239.3	11 36		I.NEKTSPALVIPTPDAENEI.S	U
L	YLR035C									
M	4	322	2040.365	0.1699	1.3183	160.0	9 36		I.LKESKSVQPGKAIPDIIES.P	U
L	YJL126W									

A nice file format converter - Proteowizard

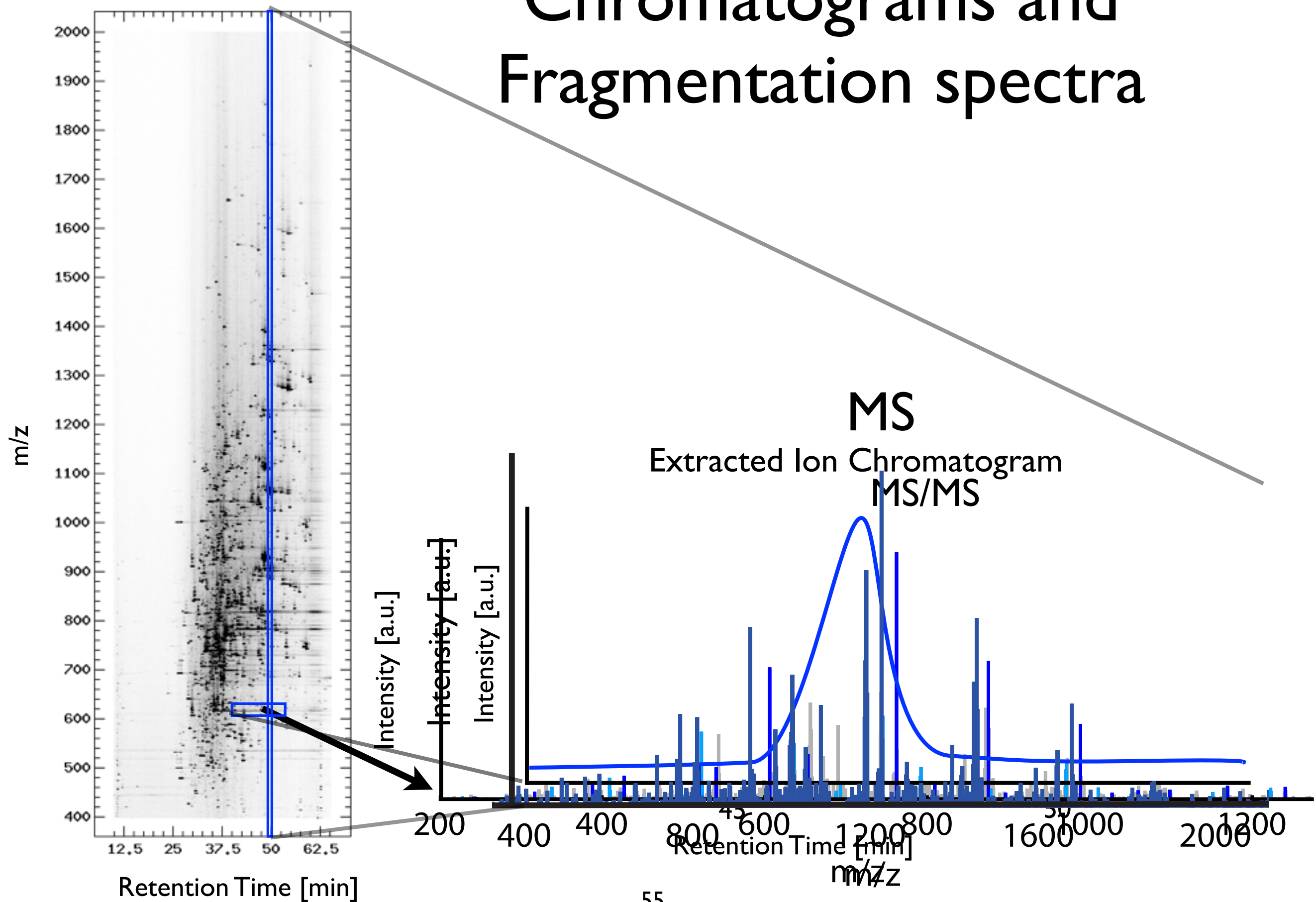


# Inferring proteins

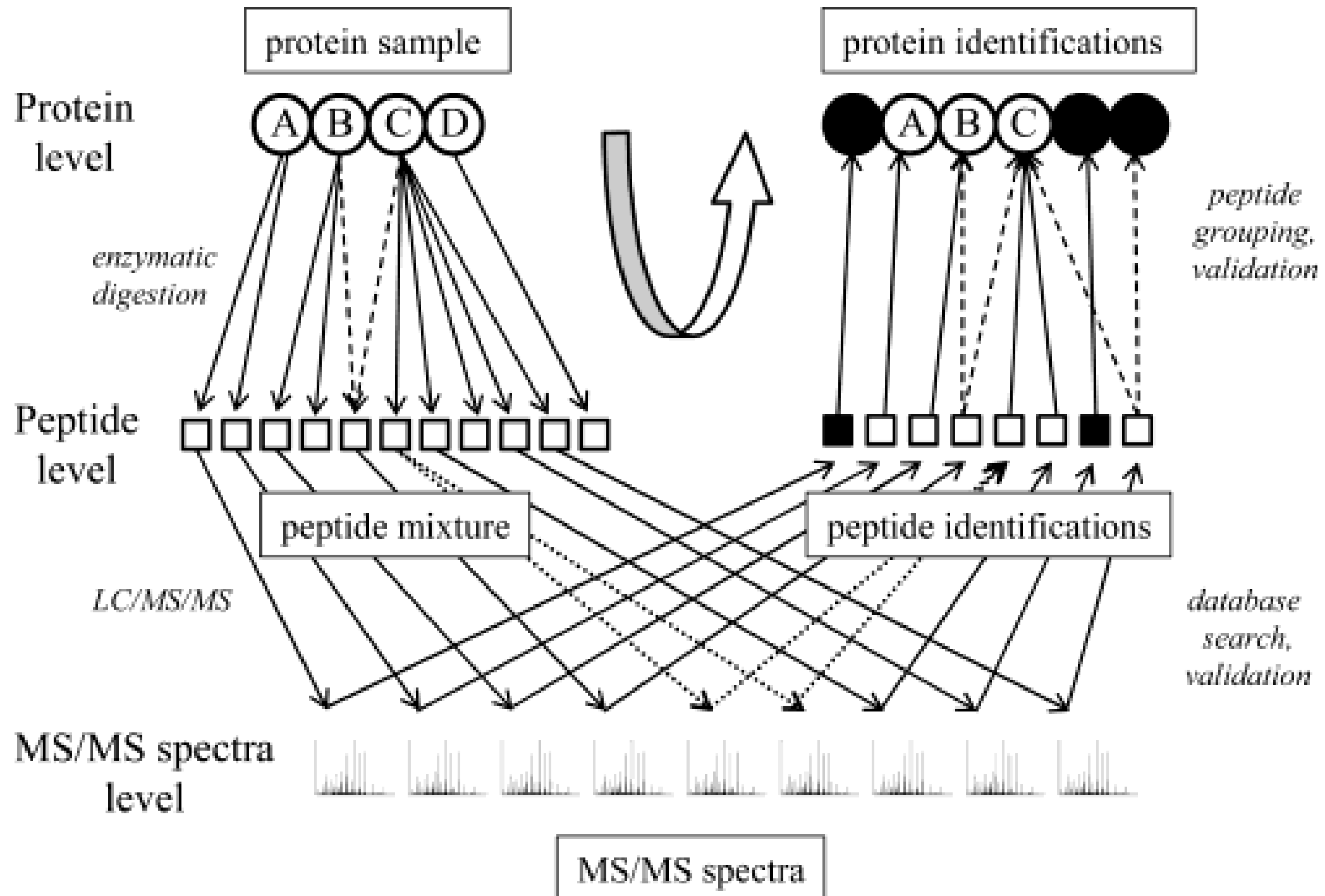
# Shotgun proteomics



# Chromatograms and Fragmentation spectra

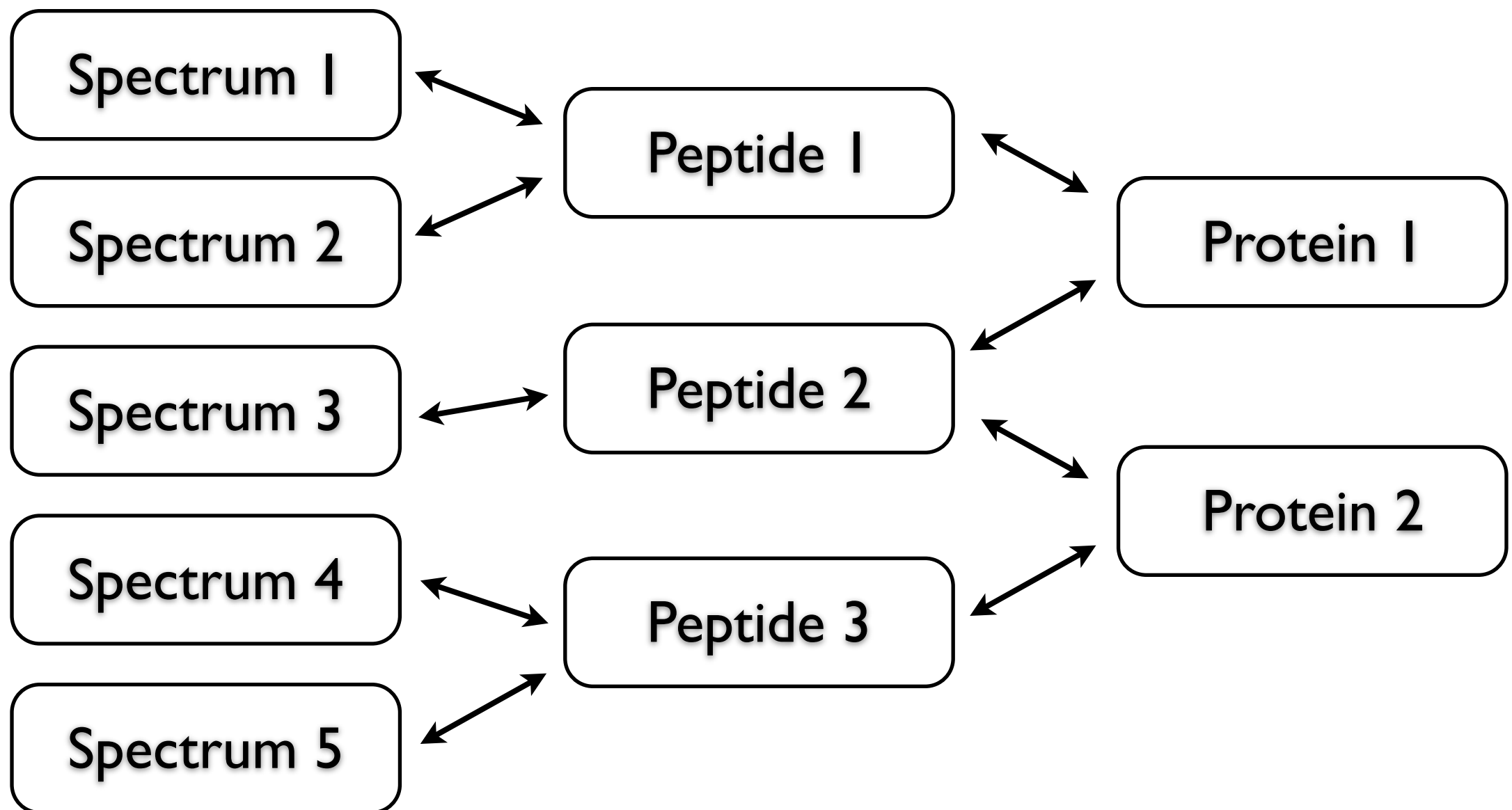


# protein identification

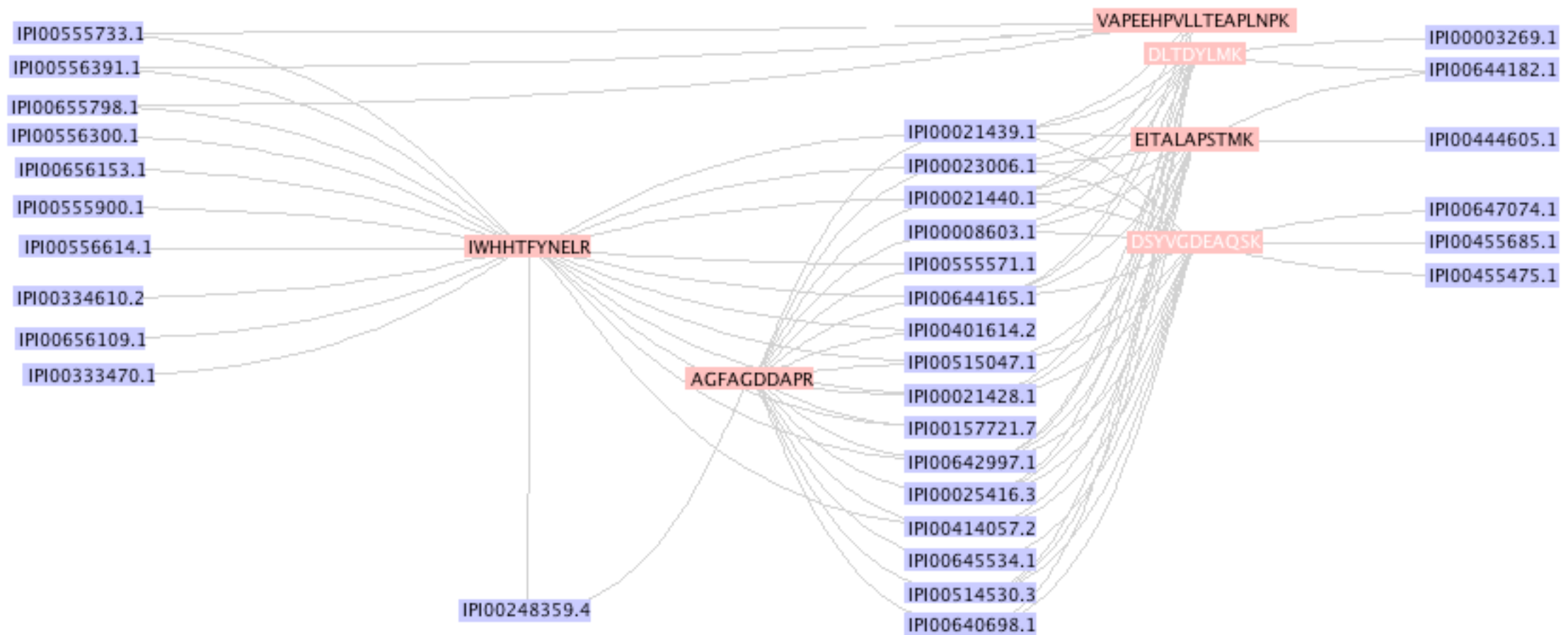




# PSM/Peptide/Protein level



# Shared peptides



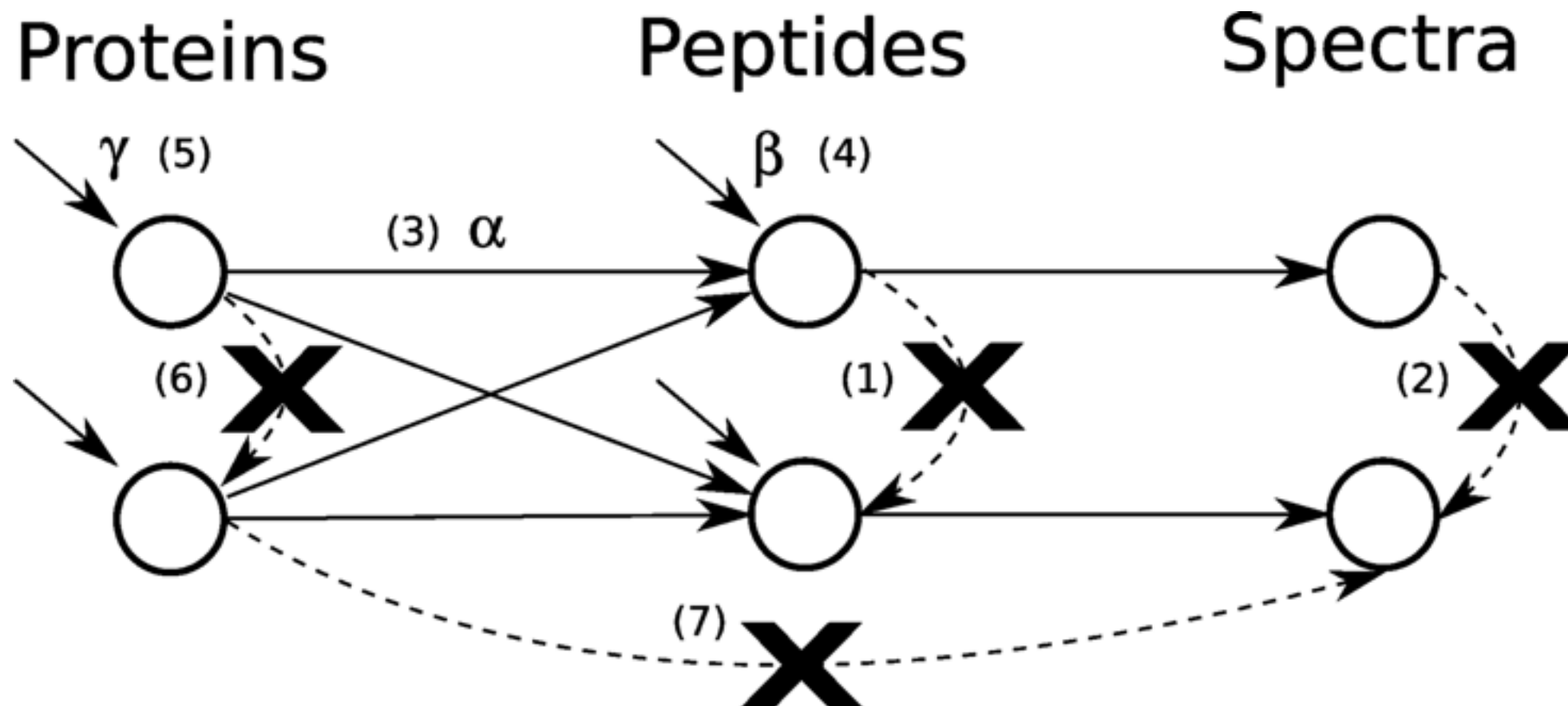
[MassSieve]

# Quality assessment of identified proteins

Two strategies:

1. Design a probabilistic models from which we may infer protein level probabilities
2. Target-decoy competition on protein level

# Bayesian Approach



- 1 Conditional Independence of Peptides Given Proteins
- 2 Conditional Independence of Spectra Given Peptides
- 3 Emission of a Peptide Associated with a Present Protein
- 4 Creation of a Peptide from Noise
- 5 Prior Belief a Protein Is Present in the Sample
- 6 Independence of Prior Belief between Proteins
- 7 Dependence of a Spectrum Only on the Best-Matching Peptide

Find MAP estimate protein set by evaluating  
 $\Pr(\text{Proteins}|\text{Spectra})$

[Serang *et al* JPR 2010]

# Papers for L9

## Assigning Significance to Peptides Identified by Tandem Mass Spectrometry Using Decoy Databases

Lukas Käll,<sup>†</sup> John D. Storey,<sup>†,‡</sup> Michael J. MacCoss,<sup>†</sup> and William Stafford Noble<sup>\*,†,§</sup>

*Departments of Genome Sciences, Biostatistics, and Computer Science and Engineering, University of Washington, Seattle, Washington 98195*

Received September 18, 2007

Automated methods for assigning peptides to observed tandem mass spectra typically return a list of peptide–spectrum matches, ranked according to an arbitrary score. In this article, we describe methods for converting these arbitrary scores into more useful statistical significance measures. These methods employ a decoy sequence database as a model of the null hypothesis, and use false discovery rate (FDR) analysis to correct for multiple testing. We first describe a simple FDR inference method, and then describe how estimating and taking into account the percentage of incorrectly identified spectra in the entire data set can lead to increased statistical power.

**Keywords:** *q*-value • decoy database • false discovery rate • statistical significance • peptide identification

ion

The problem in the analysis of tandem mass spectra is to identify the peptide that gave rise to an observed fragmentation. The most commonly used tools for solving this problem, such as SEQUEST,<sup>1</sup> Mascot,<sup>2</sup> or X!Tandem,<sup>3</sup> search a peptide sequence database for the peptide whose theoretical fragmentation best matches the observed spectrum. The output of these tools is a collection of peptide–spectrum matches (PSMs), each with an associated score (Table 1). The subsequent question is, “Which of these PSMs are correct?”

These algorithms are very powerful, the problem is that

**Table 1.** Terminology

PSM	A peptide–spectrum match, with an associated score
target PSM	A PSM created by searching the observed database
decoy database	A shuffled or reversed version of the observed database
decoy PSM	A PSM created by searching a decoy database
accepted PSM	A PSM whose score is above some threshold
correct PSM	A PSM whose peptide corresponds to the peptide that generated the observed spectrum

## ARTICLE

doi:10.1038

## A draft map of the human proteome

Min-Sik Kim<sup>1,2</sup>, Sneha M. Pinto<sup>3</sup>, Derese Getnet<sup>1,4</sup>, Raja Sekhar Nirujogi<sup>3</sup>, Srikanth S. Manda<sup>3</sup>, Raghothama Chaerkasuri<sup>3</sup>, Anil K. Madugundu<sup>3</sup>, Dhanashree S. Kelkar<sup>3</sup>, Ruth Isserlin<sup>5</sup>, Shobhit Jain<sup>5</sup>, Joji K. Thomas<sup>3</sup>, Babylakshmi Muthusamy<sup>3</sup>, Pamela Leal-Rojas<sup>1,6</sup>, Praveen Kumar<sup>3</sup>, Nandini A. Sahasrabudhe<sup>3</sup>, Lavanya Balakrishnan<sup>3</sup>, Jayshree Advani<sup>3</sup>, Bijesh S. Santosh Renuse<sup>3</sup>, Lakshmi Dhevi N. Selvan<sup>3</sup>, Arun H. Patil<sup>3</sup>, Vishalakshi Nanjappa<sup>3</sup>, Aneesha Radhakrishnan<sup>3</sup>, Samarjeet Tejaswini Subbannayya<sup>3</sup>, Rajesh Raju<sup>3</sup>, Manish Kumar<sup>3</sup>, Sreelakshmi K. Sreenivasamurthy<sup>3</sup>, Arivusudar Marimuthu<sup>3</sup>, Gajanan J. Sathe<sup>3</sup>, Sandip Chavan<sup>3</sup>, Keshava K. Datta<sup>3</sup>, Yashwanth Subbannayya<sup>3</sup>, Apeksha Sahu<sup>3</sup>, Soujanya D. Yelamanchi<sup>3</sup>, Savita Jayaram<sup>3</sup>, Pavithra Rajagopalan<sup>3</sup>, Jyoti Sharma<sup>3</sup>, Krishna R. Murthy<sup>3</sup>, Nazia Syed<sup>3</sup>, Renu Goel<sup>3</sup>, Aafaque A. Khan<sup>3</sup>, Sartaj Ahmad<sup>3</sup>, Gourav Dey<sup>3</sup>, Keshav Mudgal<sup>7</sup>, Aditi Chatterjee<sup>3</sup>, Tai-Chung Huang<sup>1</sup>, Jun Zhong<sup>1</sup>, Xinyan Wu<sup>1,2</sup>, Patrick J. Donaldson<sup>1</sup>, Muhammad S. Zahari<sup>2</sup>, Kanchan K. Mukherjee<sup>8</sup>, Subramanian Shankar<sup>9</sup>, Anita Mahadevan<sup>10,11</sup>, Henryk J. Christensen<sup>1</sup>, Christopher J. Mitchell<sup>1</sup>, Susarla Krishna Shankar<sup>10,11</sup>, Parthasarathy Satishchandra<sup>13</sup>, John T. Schroeder<sup>14</sup>, Ravi Sirdekar<sup>15,16</sup>, Anirban Maltra<sup>15,16</sup>, Steven D. Leach<sup>1,17</sup>, Charles G. Drake<sup>16,18</sup>, Marc K. Halushka<sup>15</sup>, T. S. Keshava Prasad<sup>3</sup>, Ralph H. H. Heisterkamp<sup>19</sup>, Candace L. Kerr<sup>19</sup>, Gary D. Bader<sup>5</sup>, Christine A. Iacobuzio-Donahue<sup>15,16,17</sup>, Harsha Gowda<sup>3</sup> & Akhilesh Pandey<sup>1,2,3,4</sup>

The availability of human genome sequence has transformed biomedical research over the past decade. However, a draft map for the human proteome with direct measurements of proteins and peptides does not exist yet. Here we present a draft map of the human proteome using high-resolution Fourier-transform mass spectrometry. In-depth proteomic profiling of 30 histologically normal human samples, including 17 adult tissues, 7 fetal tissues and 6 purified primary hematopoietic cells, resulted in identification of proteins encoded by 17,294 genes accounting for approximately 84% of the total annotated protein-coding genes in humans. A unique and comprehensive strategy for proteogenomics enabled us to discover a number of novel protein-coding regions, which includes translated pseudogenes, non-coding RNAs and upstream open reading frames. This large human proteome catalogue (available as an interactive web resource at <http://www.humanproteomemap.org>) will complement available human genome and transcriptome data to accelerate biomedical research in health and disease.

Analysis of the complete human genome sequence has thus far led to the identification of approximately 20,687 protein-coding genes<sup>1</sup>, although the annotation still continues to be refined. Mass spectrometry has revolutionized proteomics studies in a manner analogous to the impact of next-generation sequencing on genomics and transcriptomics<sup>2–4</sup>. Several groups, including ours, have used mass spectrometry to catalogue complete proteomes of unicellular organisms<sup>5–7</sup> and to explore proteomes

sets—PeptideAtlas<sup>11</sup>, GPMDB<sup>12</sup> and neXtProt<sup>13</sup> (which integrates data from the Human Protein Atlas<sup>14</sup>).

A general limitation of current proteomics methods is dependence on predefined protein sequence databases for identification of proteins. To overcome this, we also used a comprehensive proteomic analysis strategy to identify novel peptides/proteins that are not part of annotated protein databases. This approach