

# Introduction to Database Searching using MASCOT

**MASCOT**



This presentation introduces the topics we will discuss in depth in subsequent talks. It assumes a working knowledge of protein chemistry and mass spectrometry, but no experience of database searching.

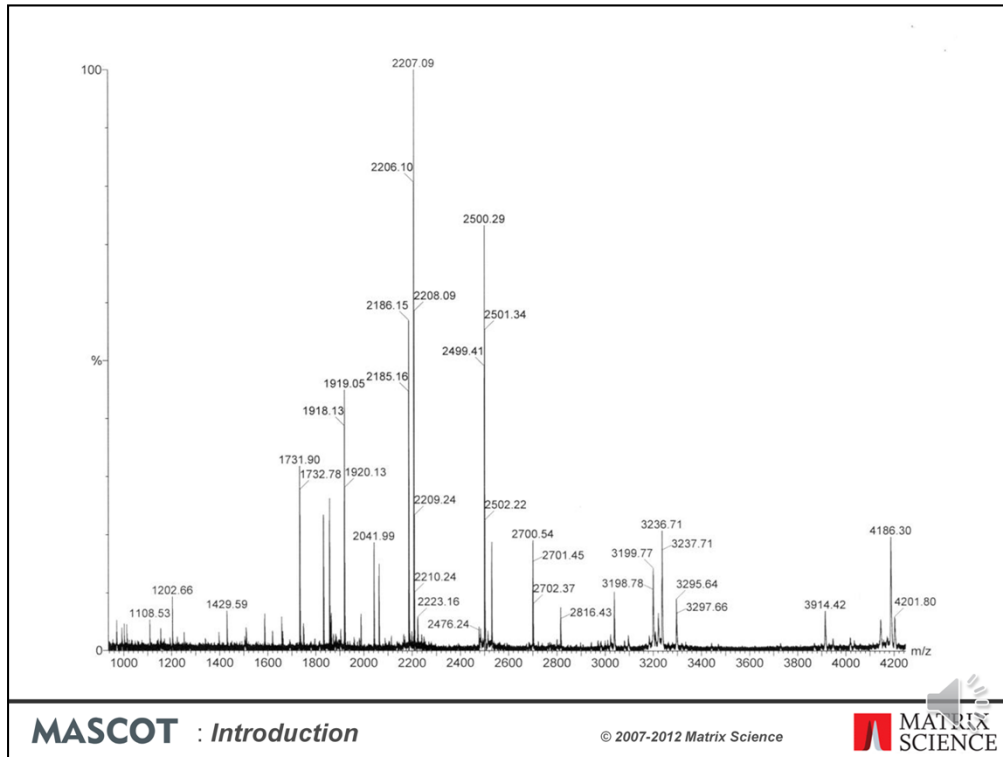
## Three ways to use mass spectrometry data for protein identification

### 1. Peptide Mass Fingerprint

A set of peptide molecular masses from an enzyme digest of a protein

There are three proven ways of using mass spectrometry data for protein identification. The first of these is known as a peptide mass fingerprint. This was the original method to be developed, and uses the molecular masses of the peptides resulting from digestion of a protein by a specific enzyme.





If the mass spectrum of your peptide digest mixture looks as good as this, and it is a single protein, and the protein sequence or something very similar is in the database, your chances of success are very high.

We don't submit the raw data to the search engine. First of all, the spectrum must be reduced to a peak list: a set of mass and intensity pairs, one for each peak. We call this procedure peak detection or peak picking.

In a peptide mass fingerprint, it is the mass values of the peaks that matter most. The peak area or intensity values are a function of peptide basicity, length, and several other physical and chemical parameters. There is no particular reason to assume that a big peak is interesting and a small peak is less interesting. The main use of intensity information is to distinguish signal from noise.

Mass accuracy is important, but so is coverage. Better to have a large number of mass values with moderate accuracy than one or two mass values with very high accuracy.

## PMF Servers on the Web

**ASCQ\_ME:** [https://www.genopole-lille.fr/logiciel/ascq\\_me/](https://www.genopole-lille.fr/logiciel/ascq_me/)

**Bupid:** <http://zlab.bu.edu/Amemee/>

**Mascot:** [http://www.matrixscience.com/search\\_form\\_select.html](http://www.matrixscience.com/search_form_select.html)

**MassSearch:** [http://www.cbrg.ethz.ch/services/MassSearch\\_new](http://www.cbrg.ethz.ch/services/MassSearch_new)

**MS-Fit (Protein Prospector):**

<http://prospector.ucsf.edu/prospector/mshome.htm>

**PepMAPPER:** <http://www.nwsr.manchester.ac.uk/mapper/>

**Profound (Prowl):** <http://prowl.rockefeller.edu/prowl-cgi/profound.exe>



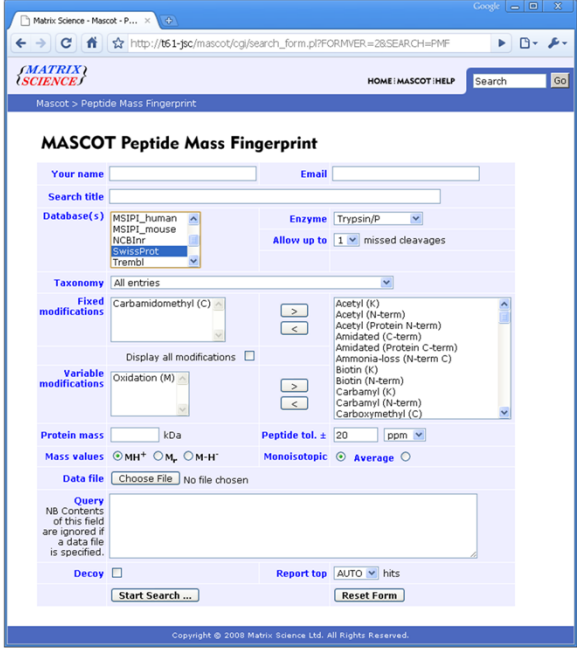
Mowse, PeptideSearch, Protocall, Aldente, XProteo

**MASCOT** : *Introduction*

© 2007-2012 Matrix Science



These presentations will focus on Mascot, but you should be aware that there are several other PMF search engines on the web. There are also software packages available for download to run locally or sold as commercial products. Some of the early search engines, such as Mowse and PeptideSearch, are no longer available.




## Search Parameters

- database
- taxonomy
- enzyme
- missed cleavages
- fixed modifications
- variable modifications
- protein MW
- estimated mass measurement error

**MASCOT** : Introduction

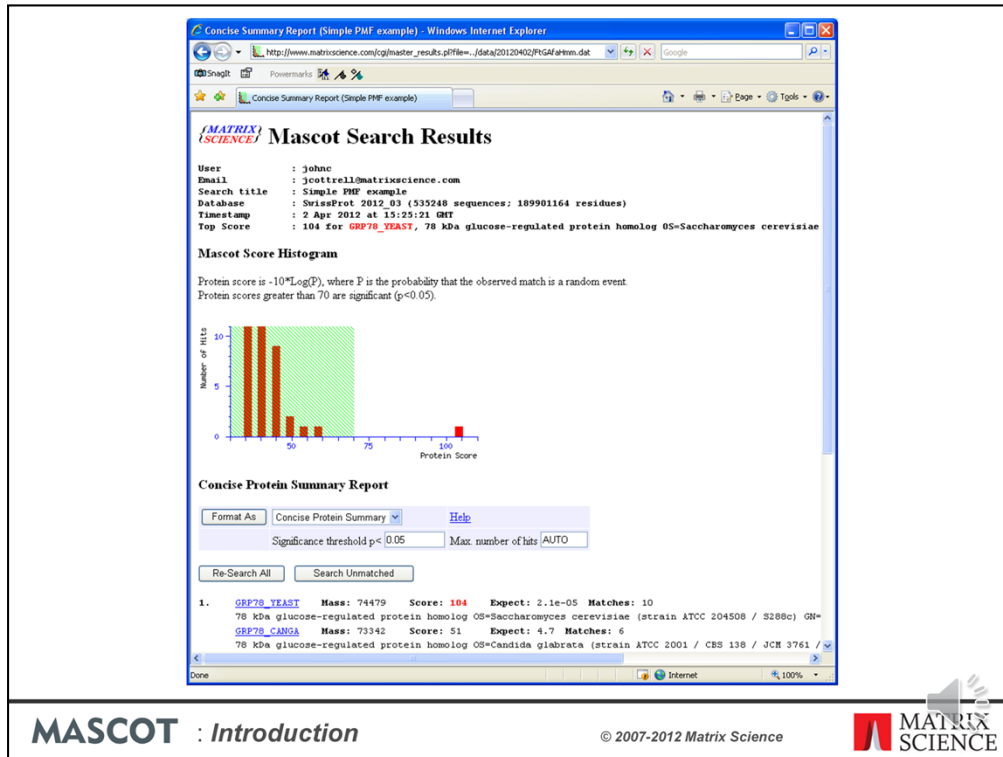
© 2007-2012 Matrix Science



This is the Mascot search form for a peptide mass fingerprint. Besides the MS data, a number of search parameters are required. Some search engines require fewer parameters, others require more. We'll be discussing most of these search parameters in detail in a later presentation.

In theory, you could design a search engine that didn't require search parameters, and tried to work everything out from the mass values, but this would be very inefficient. If you know the enzyme was trypsin, much easier to supply this information as part of the search.

To perform a search, you paste your peak list into the search form, or upload it as a file, enter values for the search parameters, and press the submit button.



A short while later, you receive the results.

A peptide mass fingerprint search will almost always produce a list of matching proteins, and something has to be at the top of that list. One of the main problems in the early days of the technique was how to tell whether the top match was “real”, or just the match at the top of the list ... that is, a false positive.

There have been various attempts to deal with this problem, which I will describe when we come to discuss scoring.

## Protein Identification: The Origins of Peptide Mass Fingerprinting

William J. Henzel and Colin Watanabe

Protein Chemistry Department and Bioinformatics Department, Genentech, Inc.,  
South San Francisco, California, USA

John T. Stults

Analytical Sciences Department, Biospect, Inc., South San Francisco, California, USA

Peptide mass fingerprinting (PMF) grew from a need for a faster, more efficient method to identify frequently observed proteins in electrophoresis gels. We describe the genesis of the idea in 1989, and show the first demonstration with fast atom bombardment mass spectrometry. Despite its promise, the method was seldom used until 1992, with the coming of significantly more sensitive commercial instrumentation based on MALDI-TOF-MS. We recount the evolution of the method and its dependence on a number of technical breakthroughs, both in mass spectrometry and in other areas. We show how it laid the foundation for high-throughput, high-sensitivity methods of protein analysis, now known as proteomics. We conclude with recommendations for further improvements, and speculation of the role of PMF in the future. (J Am Soc Mass Spectrom 2003, 14, 931-942) © 2003 American Society for Mass Spectrometry

➤Henzel, W. J., Watanabe, C., Stults, J. T., *JASMS* 2003, 14, 931-942.

**MASCOT** : Introduction

© 2007-2012 Matrix Science



If you want to learn more about the origins and development of peptide mass fingerprinting, I strongly recommend this review by the Genentech group. They discuss the history and the methodology in a very readable style.



## Peptide Mass Fingerprint



**Fast, simple analysis**

**High sensitivity**

**Need database of protein sequences**

- not ESTs or genomic DNA

**Sequence must be present in database**

- or close homolog

**Not good for mixtures**

- especially a minor component.

One of the strengths of PMF is that it is an easy experiment that can be performed using just about any mass spectrometer. The whole process is readily automated and MALDI instruments, in particular, can churn out high accuracy PMF data at a very high rate.

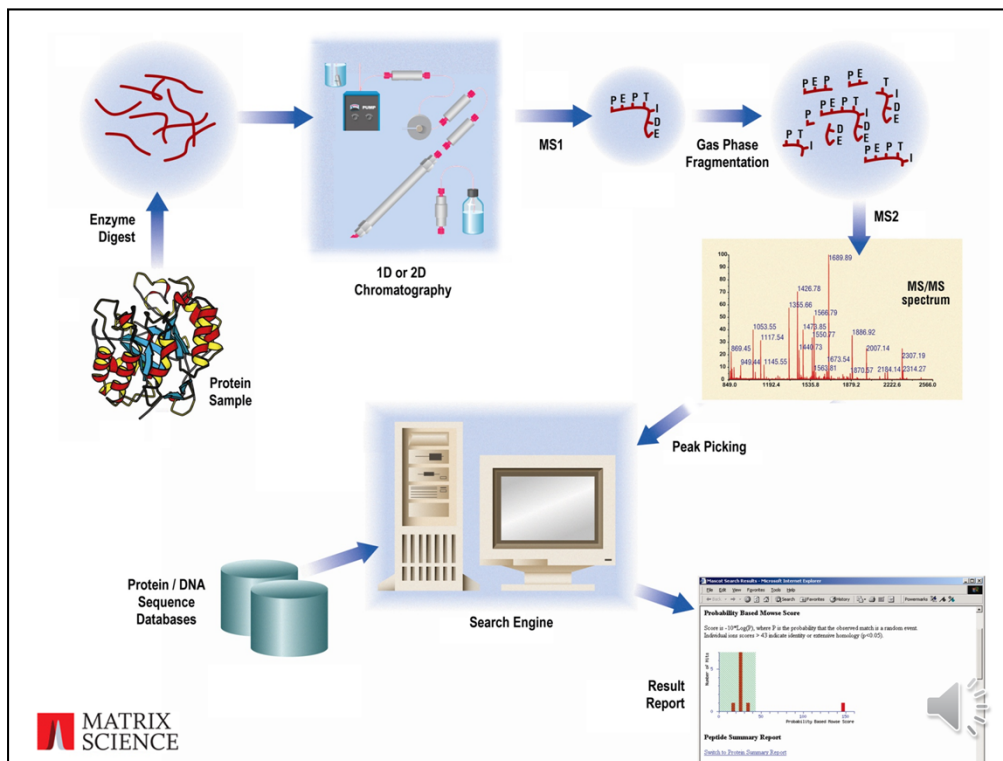
In principal, it is a sensitive technique because you don't need 100% coverage. It doesn't matter too much if a small part of the protein fails to digest or some of the peptides are insoluble or don't fly very well.

One of the limitations is that you need a database of proteins or nucleic acid sequences that are equivalent to proteins, e.g. mRNAs. In most cases, you will not get satisfactory results from an EST database, where most of the entries correspond to protein fragments, or genomic DNA, where there is a continuum of sequence, containing regions coding for multiple proteins as well as non-coding regions. This is because the statistics of the technique rely on the set mass values having originated from a defined protein sequence. If multiple sequences are combined into a single entry, or the sequence is divided between multiple entries, the numbers may not work.

If the protein sequence, or something very similar, is not in the database, the method will fail. If you are studying a well characterised organism, such as human or mouse or yeast, this is unlikely to be a problem. If you are studying a virus or plant with an unsequenced genome, it can be a major problem, and you depend on getting matches to homologous proteins from related organisms.

The most important limitation concerns mixtures. If the data quality is good, then it may be possible to identify a two component mixture, where both components are at a similar level, and on very rare occasions three. But if the data are poor, it can be difficult to get any match at all out of a mixture, and it is never possible to identify a minor component.

To identify proteins from mixtures reliably, it is necessary to work at the peptide level. That is, using MS/MS data.



The experimental workflow for database matching of MS/MS data is similar to that for PMF, but with an added stage of selectivity and fragmentation.

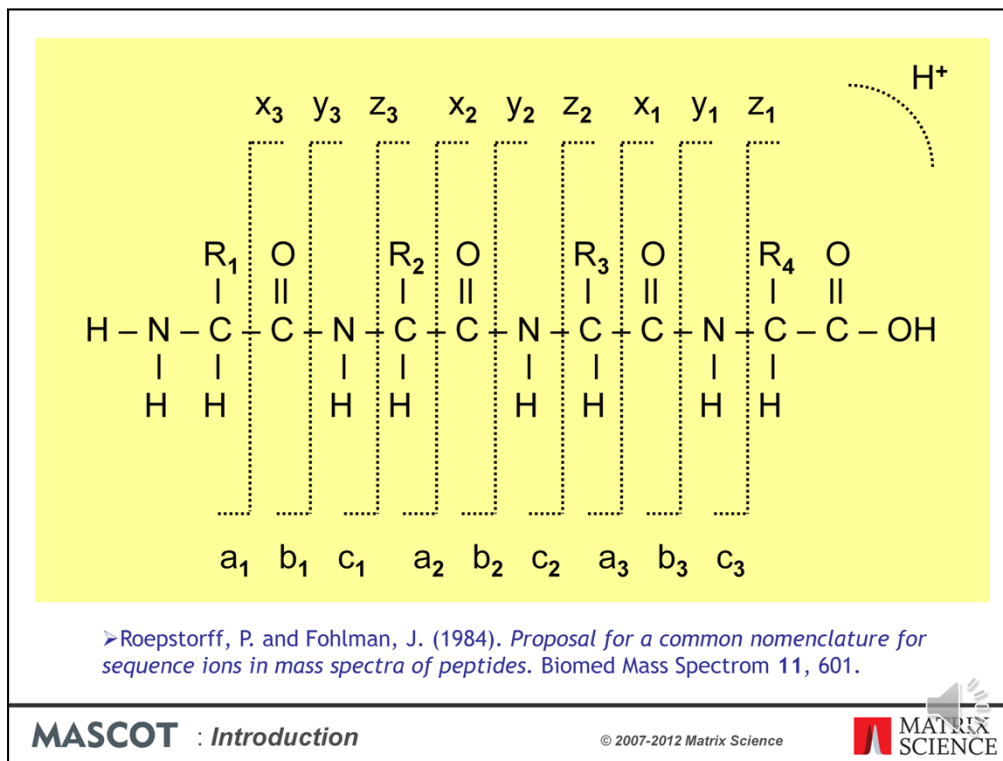
Again, we start with protein, which can now be a single protein or a complex mixture of proteins. We use an enzyme such as trypsin to digest the proteins to peptides. If it is a complex mixture, such as a whole cell lysate, we will probably use one or more stages of chromatography to regulate the flow of peptides into the mass spectrometer. We select peptides one at a time using the first stage of mass analysis. Each isolated peptide is then induced to fragment, possibly by collision, and the second stage of mass analysis used to collect an MS/MS spectrum.

Because we are collecting data from isolated peptides, it makes no difference whether the original sample was a mixture or not. We identify peptide sequences, and then try to assign them to one or more protein sequences. One consequence is that, unless a peptide is unique to one particular protein, there may be some ambiguity as to which protein it should be assigned to.

For each MS/MS spectrum, we use software to try and determine which peptide sequence in the database gives the best match. As in the case of a peptide mass fingerprint, each entry in the database is digested, *in silico*, and the masses of the expected peptides calculated. If a calculated peptide mass matches the experimental one, the mass values expected to result from the gas phase fragmentation of the peptide are calculated and the degree of matching to the peaks in the MS/MS spectrum scored.

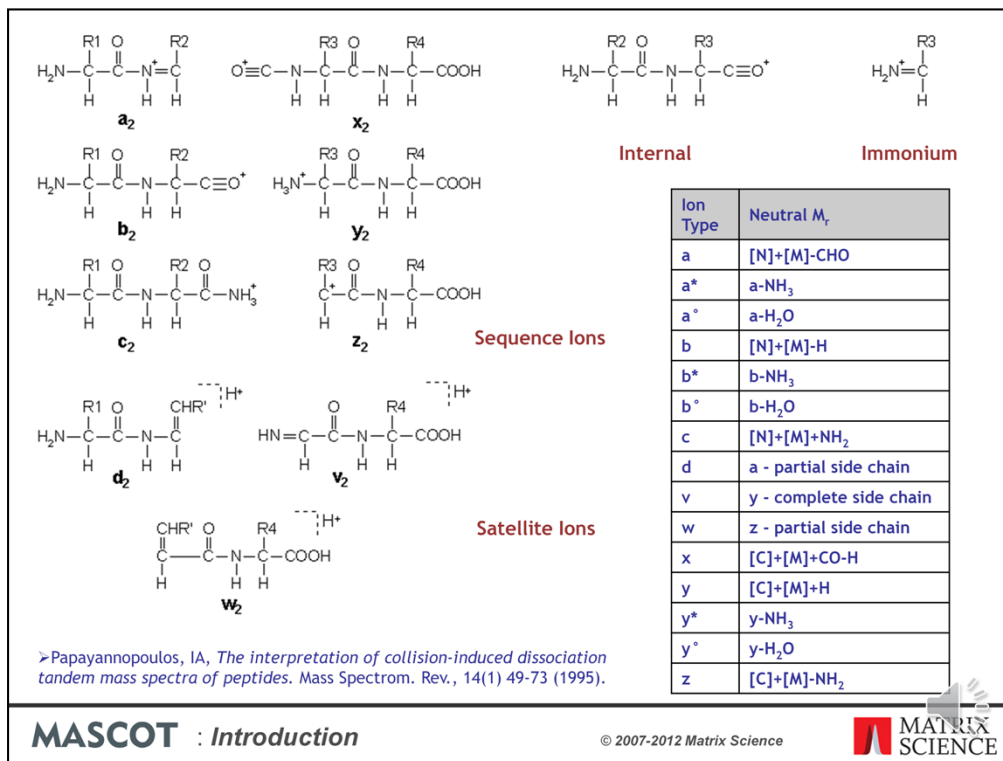
Unlike a peptide mass fingerprint, use of a specific enzyme is not essential. By looking at

all possible sub-sequences of each entry that fit the precursor mass, it is possible to match peptides when the enzyme specificity is unknown, such as endogenous peptides.



Database matching of MS/MS data is only possible because peptide molecular ions fragment at preferred locations along the backbone. In many instruments, the major peaks in an MS/MS spectrum are b ions, where the charge is retained on the N-terminus, and y ions, where the charge is retained on the C-terminus.

However, this depends on the ionisation technique, the mass analyser, and the peptide structure. Electron capture dissociation, for example, produces predominantly c and z ions.



If peptides fragmented cleanly and uniformly along the backbone, we wouldn't need database search. We would see a ladder of peaks for each ion series, where the distance from one peak to the next was the mass of an amino acid residue, allowing the sequence to be read off the spectrum. In real life, fragmentation is rarely perfect, and the spectrum will usually show significant peaks from side chain cleavages and internal fragments, where the backbone has been cleaved twice. More importantly, the backbone may fail to cleave at certain locations, so that the MS/MS spectrum has no evidence for some of the residues.

This slide shows the most common fragment ion structures and the table is a "ready reckoner" that can be used to calculate the masses. N is mass of the N-terminal group, (hydrogen for free amine). C is the mass of the C-terminal group, (hydroxyl for free acid). M is the sum of the residue masses

This review by Ioannis Papayannopoulos is a good introduction to the fragmentation chemistry of peptide ions in the gas phase

----

To determine the neutral mass of, say a 'b' ion with just two glycines, add the mass of the n terminal group, which is normally just a hydrogen, so '1', the mass of two glycines is 114 and subtract a hydrogen which leaves a mass of 114. To get the singly charged ion, we need to add a proton, which gives a mass/charge of approximately 115.

## Three ways to use mass spectrometry data for protein identification

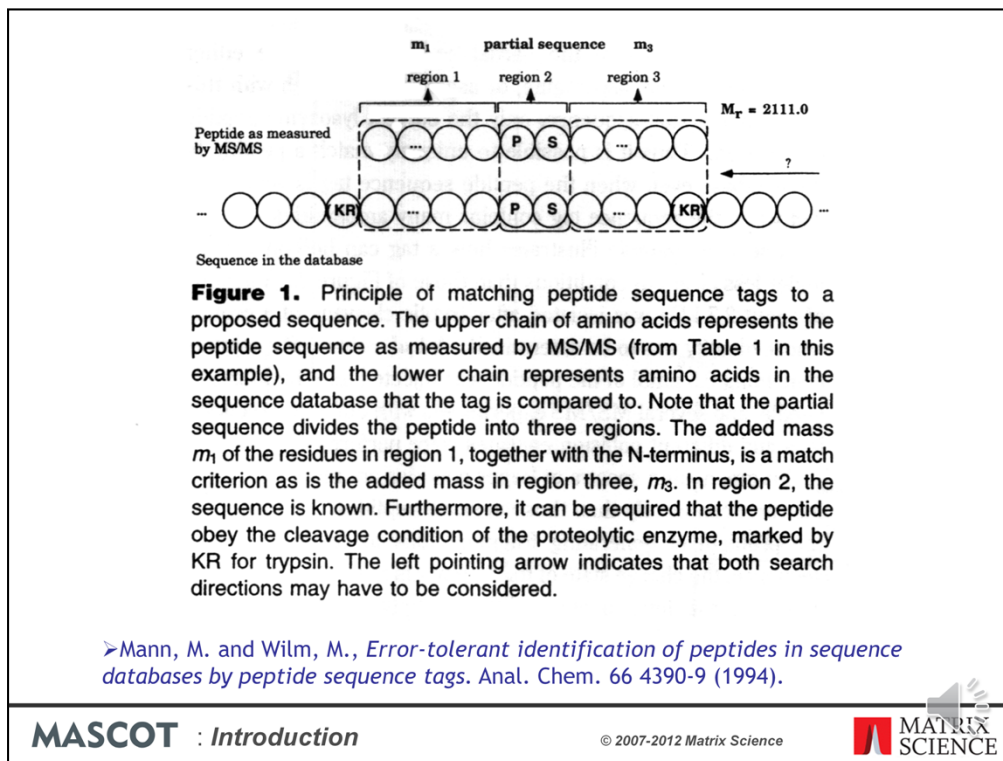
### 1. Peptide Mass Fingerprint

A set of peptide molecular masses from an enzyme digest of a protein

### 2. Sequence Query

Mass values combined with amino acid sequence or composition data

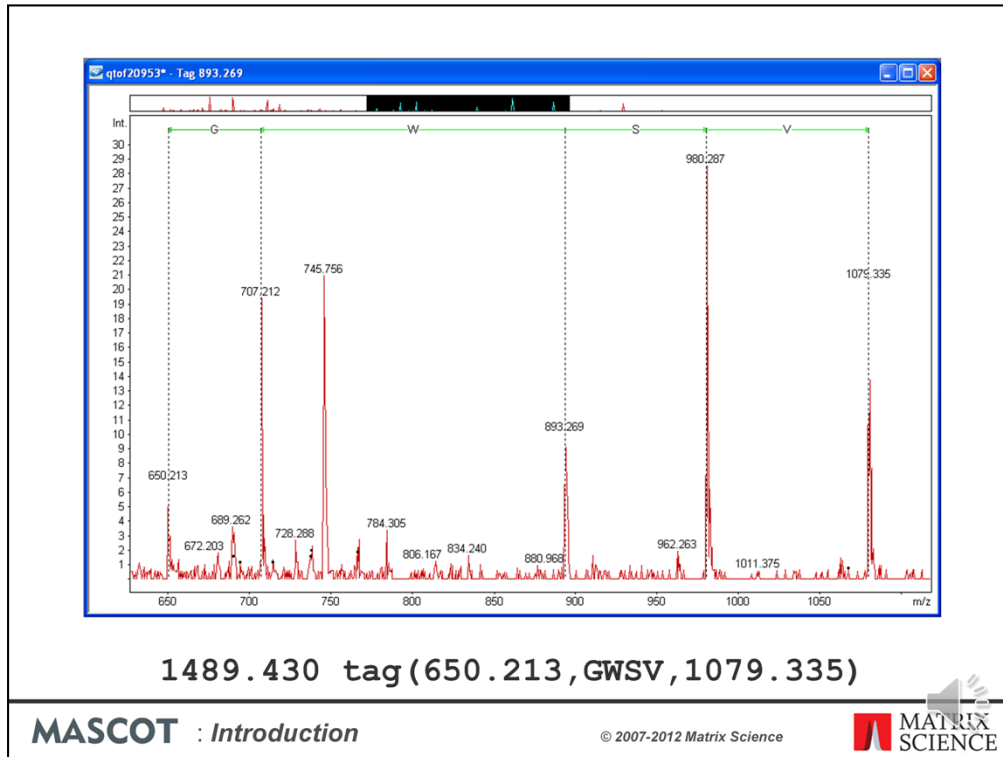
Which brings us to the second method of using mass spectrometry data for protein identification: a sequence query in which mass information is combined with amino acid sequence or composition data. The most widely used approach in this category is the sequence tag, developed by Matthias Mann and Matthias Wilm at EMBL.



In a sequence tag search, a few residues of amino acid sequence are interpreted from the MS/MS spectrum.

Even when the quality of the spectrum is poor, it is often possible to pick out four clean peaks, and read off three residues of sequence. In a sequence homology search, a triplet would be worth almost nothing, since any given triplet can be expected to occur by chance many times in even a small database.

What Mann and Wilm realised was that this very short stretch of amino acid sequence might provide sufficient specificity to provide an unambiguous identification if it was combined with the fragment ion mass values which enclose it, the peptide mass, and the enzyme specificity.



Picking out a good tag is not trivial, and requires both luck and experience. In this spectrum, we can see a promising four residue tag. The syntax used by Mascot for a sequence tag is shown below the spectrum. We'll discuss this format in greater detail in the Sequence Query presentation



## Sequence Query Servers on the Web

### Mascot

- [http://www.matrixscience.com/search\\_form\\_select.html](http://www.matrixscience.com/search_form_select.html)

### MS-Seq (Protein Prospector)

- <http://prospector.ucsf.edu/prospector/mshome.htm>

### TagIdent

- <http://web.expasy.org/tagident/>



**PeptideSearch**

**MASCOT** : Introduction

© 2007-2012 Matrix Science



There are a number of software packages for sequence query searches. As with PMF, I have limited my list to servers that are publicly available on the web. Not such a wide choice as for PMF.

**MASCOT Sequence Query**

Your name: Lou Scene Email: lou@res.edu

Search title:

Database(s): MSIP1\_mouse, NCBItr, UniRef100 Enzyme: Trypsin/P

Allow up to: 1 missed cleavages Quantitation: None

Taxonomy: All entries

Fixed modifications: Carbamidomethyl (C)

Variable modifications: none selected

Peptide tol.: 0.6 Da MS/MS tol.: 0.6 Da

Peptide charge: Mr Monoisotopic ☒ Average ☐

Query: 1489.430 tag(650.213, GWSV, 1079.335)

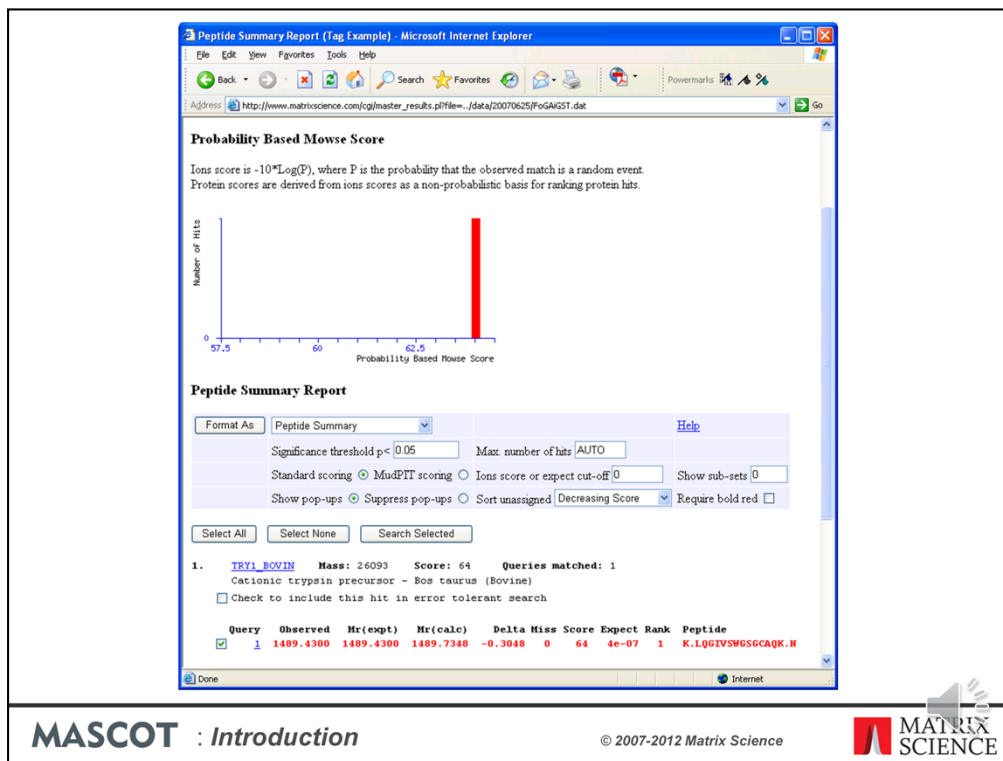
Instrument: ESI-TRAP Decoy: ☐ Report top: AUTO hits

Start Search ... Reset Form

Copyright © 2008 Matrix Science Ltd. All Rights Reserved.

**MASCOT : Introduction** © 2007-2012 Matrix Science **MATRIX SCIENCE**

I entered the tag shown earlier into the Mascot Sequence Query search form. As with a PMF, several search parameters are required, such as the database to be searched and an estimate of the mass accuracy.



This is the result report from the search. There is just one peptide in the database that matches: LQGIVSWGSGCAQK from bovine trypsinogen.

The score is good, but even if it wasn't, you are on very safe ground accepting any match to trypsin, keratin, or BSA ;-)

## Sequence Tag

### Rapid search times

- Essentially a filter

### Error tolerant

- Match peptide with unknown modification or SNP

### Requires interpretation of spectrum

- Usually manual, hence not high throughput

### Tag has to be called correctly

- Although ambiguity is OK

2060.78 tag(977.4,[Q|K][Q|K][Q|K]EE,1619.7).

A sequence tag search can be rapid, because it is simply a filter on the database.

However, the standard sequence tag is essentially obsolete. It is easier and more reliable to skip the interpretation step and pass the peak list to the search engine. The reason the sequence tag is still important is because it can be used in an “error tolerant” mode. This consists of relaxing the specificity, by removing the peptide molecular mass constraint. The tag is effectively allowed to float within the candidate sequence, so that a match is possible even if there is a difference in the calculated mass to one side or the other of the tag. This is one of the few ways of getting a match to a peptide when there is an unsuspected modification or a variation in the primary amino acid sequence.

Tags can be called by software. But, in most cases, they are called manually, which requires time and skill.

If the tag is not correct, then no match will be found. In Mascot, ambiguity is OK, as long as it is recognised and the query is formulated correctly. Obviously, I=L and, in most cases, Q=K and F=MetOx. Software or a table of mass values can help identify the more common ambiguities. Even so, it is very difficult to identify all possible ambiguities, especially when we allow for missing peaks.

## Three ways to use mass spectrometry data for protein identification

### 1. Peptide Mass Fingerprint

A set of peptide molecular masses from an enzyme digest of a protein

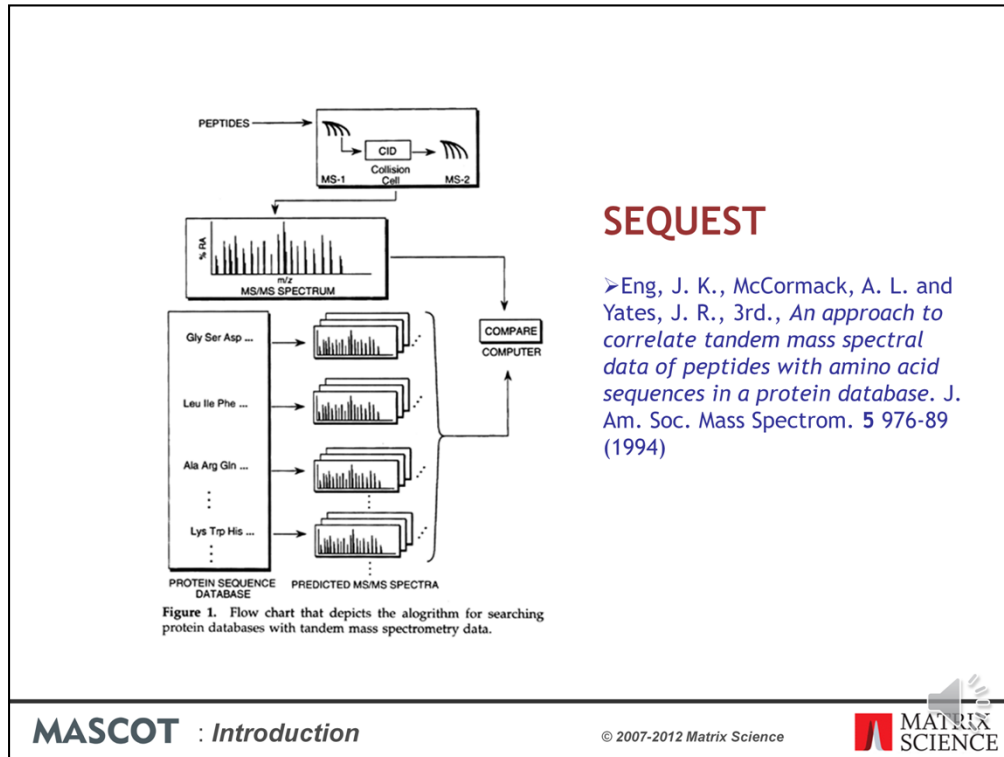
### 2. Sequence Query

Mass values combined with amino acid sequence or composition data

### 3. MS/MS Ions Search

Uninterpreted MS/MS data from a single peptide or from a complete LC-MS/MS run

Which brings us to the third category: Searching the uninterpreted MS/MS data from a single peptide or from a complete LC-MS/MS run. That is, using software to match the peak list, without any manual sequence calling.



This approach was pioneered by John Yates and Jimmy Eng at the University of Washington, Seattle. They used a cross correlation algorithm to compare an experimental MS/MS spectrum against spectra predicted from peptide sequences from a database. Their ideas were implemented as the Sequest program.

## MS/MS Ions Search Servers on the Web

Inspect	<a href="http://proteomics.ucsd.edu/LiveSearch/">http://proteomics.ucsd.edu/LiveSearch/</a>
Mascot	<a href="http://www.matrixscience.com/search_form_select.html">http://www.matrixscience.com/search_form_select.html</a>
MS-Tag (Protein Prospector)	<a href="http://prospector.ucsf.edu/prospector/mshome.htm">http://prospector.ucsf.edu/prospector/mshome.htm</a>
Omssa	<a href="http://pubchem.ncbi.nlm.nih.gov/omssa/index.htm">http://pubchem.ncbi.nlm.nih.gov/omssa/index.htm</a>
PepFrag (Prowl)	<a href="http://prowl.rockefeller.edu/prowl/pepfrag.html">http://prowl.rockefeller.edu/prowl/pepfrag.html</a>
PepProbe	<a href="http://bart.scripps.edu/public/search/pep_probe/search.jsp">http://bart.scripps.edu/public/search/pep_probe/search.jsp</a>
RAId_DbS	<a href="http://www.ncbi.nlm.nih.gov/CBBResearch/qmbp/RAId_DbS/index.html">http://www.ncbi.nlm.nih.gov/CBBResearch/qmbp/RAId_DbS/index.html</a>
Sonar (Knexus)	<a href="http://hs2.proteome.ca/prowl/knexus.html">http://hs2.proteome.ca/prowl/knexus.html</a>
X!Tandem (The GPM)	<a href="http://thegpm.org/TANDEM/index.html">http://thegpm.org/TANDEM/index.html</a>
Not on-line	Byonic, Crux, greytag, MassMatrix, Myrimatch, Paragon, Peaks, PepSplice, pFind, Phenyx, ProbiD, ProLuCID, ProteinLynx GS, Sequest, SIMS, SpectrumMill

There is a wide choice of search engines on the web for performing searches of uninterpreted MS/MS data. I've also listed some of the packages that are not on the web, which includes Sequest.

As with a peptide mass fingerprint, the starting point is a peak list. There are several different formats for MS/MS peak lists, and this may constrain your choice of search engine

**MASCOT MS/MS Ions Search**

Your name: Expert User      Email: smarbe@matrixscience.com

Search title:

Database(s): MSIPI\_human, MSIPI\_mouse, TrEMBL, SwissProt

Enzyme: Trypsin/P

Allow up to: 1 missed cleavages

Quantitation: None

Taxonomy: All entries

Fixed modifications: Carbamidomethyl (C)

Variable modifications: Oxidation (M)

Peptide tol.  $\pm$  20 ppm      #  $^{13}\text{C}$  0      MS/MS tol.  $\pm$  0.1 Da

Peptide charge: 2+ and 3+      Monoisotopic ☒ Average ☐

Data file: Choose File      No file chosen

Data format: Mascot generic

Instrument: ESI-QUAD-TOF

Decoy: ☐

Precursor:  m/z

Error tolerant: ☐

Report top: AUTO hits

Copyright © 2008 Matrix Science Ltd. All Rights Reserved.

This is the Mascot search form for an MS/MS search. Fairly similar the previous two and, as before, you must also specify the database, mass accuracy, modifications to be considered, etc.



Peptide Summary Report (MS/MS Example) - Microsoft Internet Explorer

Address: [http://www.matrixscience.com/cgi/master\\_results.pl?file=...&data#981123.dat](http://www.matrixscience.com/cgi/master_results.pl?file=...&data#981123.dat)

1. **A32800** Mass: 61016 Score: 1195 Queries matched: 31  
 chaperonin GroEL precursor - human  
☐ Check to include this hit in error tolerant search

Query	Observed	M <sub>r</sub> (exp't)	M <sub>r</sub> (calc)	Delta	Miss	Score	Expect	Rank	Peptide
11	417.1822	832.3490	832.3827	-0.0329	0	45	0.1	1	K.APGGDMR.K
12	422.7433	843.4720	843.5065	-0.0345	0	46	0.11	1	K.VGEIVTK.D
13	430.7328	859.4510	859.4837	-0.0327	0	12	2.5e+02	3	K.IPANTIAK.H + Oxidation (H)
15	451.2499	900.4853	900.5280	-0.0427	0	52	0.025	1	K.LSDGVAVLK.V
16	456.7806	911.5467	911.5803	-0.0337	0	59	0.0041	1	K.VGLQVAVK.A
21	480.7447	959.4748	959.5036	-0.0288	0	45	0.11	1	R.VTDLNATR.A
24	595.7855	1189.5565	1189.6012	-0.0447	0	57	0.0068	1	K.EIGHTISDAHK.K
25	603.7720	1205.5294	1205.5961	-0.0668	0	(50)	0.027	1	K.EIGHTISDAHK.K + Oxidation (H)
26	608.3099	1214.6052	1214.6506	-0.0454	0	73	0.00015	1	K.NAGVEGLIVEK.I
27	617.2857	1232.5569	1232.5884	-0.0315	0	81	2.7e-05	1	K.VGTSDFVEHK.K
31	672.8375	1343.6605	1343.7085	-0.0480	0	64	0.001	1	R.VTIEQSWGSPK.V
34	714.8884	1427.7623	1427.8057	-0.0434	0	(65)	0.00086	1	R.GVHLAVDAVIAELK.K
35	714.8938	1427.7730	1427.8057	-0.0327	0	(73)	0.00013	1	R.GVHLAVDAVIAELK.K
36	722.8849	1443.7752	1443.8006	-0.0454	0	73	0.00014	1	R.GVHLAVDAVIAELK.K + Oxidation (H)
37	722.8934	1443.7722	1443.8006	-0.0284	0	(70)	0.00025	1	R.GVHLAVDAVIAELK.K + Oxidation (H)
39	752.8643	1503.7141	1503.7490	-0.0349	0	90	2.7e-06	1	K.TLHDELEIEGDK.F
40	760.8461	1519.6777	1519.7439	-0.0662	0	(84)	8.9e-06	1	K.TLHDELEIEGDK.F + Oxidation (H)
45	640.3281	1917.9625	1918.0636	-0.1010	0	102	1.3e-07	1	K.ISSIQSVPALEIAHNR.K
46	960.0327	1918.0509	1918.0636	-0.0127	0	(87)	3.2e-06	1	K.ISSIQSVPALEIAHNR.K
48	1019.5106	2037.0067	2037.0153	-0.0086	0	52	0.01	1	R.IQEITIEQLDVTTEYK.E
51	1057.0537	2112.0929	2112.1322	-0.0393	0	116	4.6e-09	1	R.ALHLQGVDLADAVAVTHGPK.G
52	1065.0399	2128.0653	2128.1271	-0.0618	0	(69)	0.00022	1	R.ALHLQGVDLADAVAVTHGPK.G + Oxidation (H)
53	1065.0623	2128.1100	2128.1271	-0.0172	0	(26)	3.9	1	R.ALHLQGVDLADAVAVTHGPK.G + Oxidation (H)
54	1073.0477	2144.0809	2144.1220	-0.0411	0	(70)	0.00018	1	R.ALHLQGVDLADAVAVTHGPK.G + 2 Oxidation (H)
58	789.1062	2364.2968	2364.3263	-0.0296	0	(56)	0.0038	1	R.KPLVTIAEDVDGEALSTLVLR.L
59	1183.1570	2364.2994	2364.3263	-0.0269	0	(65)	0.00038	1	R.KPLVTIAEDVDGEALSTLVLR.L
60	789.1094	2364.3063	2364.3263	-0.0200	0	95	4.5e-07	1	R.KPLVTIAEDVDGEALSTLVLR.L
61	828.1238	2481.3495	2481.3941	-0.0446	0	(26)	2.8	1	R.TALLDAAGVASLTIAEVVTEIPK.E
62	828.1322	2481.3748	2481.3941	-0.0193	0	48	0.02	1	R.TALLDAAGVASLTIAEVVTEIPK.E
64	854.0588	2559.1545	2559.2412	-0.0867	0	75	3.4e-05	1	K.LVQVAMTHREAGDGTITATVLAR.S

MASCOT : Introduction

© 2007-2012 Matrix Science

MATRIX SCIENCE

The results from this type of search tend to be more complicated to report. This is because the results usually represent a number of MS/MS spectra, rather than a single spectrum.

We match peptide sequences to individual MS/MS spectra, then try to assign these peptide sequences to proteins. Usually, there is ambiguity, and we aim to report a minimalk list of proteins. That is, the shortest list of proteins that can account for all the observed peptide matches. So, the report lists a series of proteins and, for each protein, the peptide matches that have been assigned. But, there is an additional dimension to the data.

Peptide Summary Report (MS/MS Example) - Microsoft Internet Explorer

Address: [http://www.matrixscience.com/cgi/master\\_results.cgi?file=...&data%961123.dat](http://www.matrixscience.com/cgi/master_results.cgi?file=...&data%961123.dat)

1. **A32800** Mass: 61016 Score: 1195 Queries matched: 31  
 chaperonin GroEL precursor - human  
☐ Check to include this hit in error tolerant search

Query	Observed	Mr(calc)	Mr(calc)	Delta	Miss	Score	Expect	Rank	Peptide
11	417.1822	832.3490	832.3827	-0.0329	0	45	0.1	1	K.APGEGDNR.K
12	422.7433	843.4720	843.5065	-0.0345	0	46	0.11	1	K.VGEIVTVK.D
13	430.7328	859.4510	859.4837	-0.0327	0	12	2.5e+02	3	K.IFANTIAK.H + Oxidation (H)
15	451.2499	900.4853	900.5280	-0.0427	0	52	0.025	1	K.LSDGVAVLK.V
16	456.7806	911.5467	911.5803	-0.0337	0	59	0.0041	1	K.VGLQVAVK.A
21	480.7447	959.4748	959.5036	-0.0288	0	45	0.11	1	K.VTDLNATK.A
24	595.7855	1189.5565	1189.6012	-0.0447	0	57	0.0068	1	K.EIGHIISDANK.K
25	603.7720	1205.5294	1205.5961	-0.0668	0	(50)	0.027	1	K.EIGHIISDANK.K + Oxidation (H)
26	608.3099	1214.6052	1214.6506	-0.0454	0	73	0.00015	1	K.NAGVEGSLIVEK.I
27	617.2857	1232.5569	1232.5884	-0.0315	0	81	2.7e-05	1	K.VGTSDEVEVEK.K
31	672.8375	1343.6605	1343.7085	-0.0480	0	64	0.001	1	R.VTIEQSWGSPK.V
34	714.8884	1427.7623	1427.8057	-0.0434	0	(65)	0.00086	1	R.GVHLAVDAVIAELK.K
35	714.8938	1427.7730	1427.8057	-0.0327	0	(73)	0.00013	1	R.GVHLAVDAVIAELK.K
36	722.8849	1443.7752	1443.8006	-0.0254	0	73	0.00014	1	R.GVHLAVDAVIAELK.K + Oxidation (H)
37	722.8934	1443.7722	1443.8006	-0.0284	0	(70)	0.00025	1	R.GVHLAVDAVIAELK.K + Oxidation (H)
39	752.8643	1503.7141	1503.7490	-0.0349	0	90	2.7e-06	1	K.TLHDELEIEGK.F
40	760.8461	1519.6777	1519.7439	-0.0662	0	(84)	8.9e-06	1	K.TLHDELEIEGK.F + Oxidation (H)
41	640.3281	1917.9625	1918.0636	-0.1010	0	102	1.3e-07	1	K.ISSIQTIVPALEIANHR.K
42	860.8322	1818.0500	1818.0636	-0.0137	0	(82)	3.2e-06	1	K.ISSIQTIVPALEIANHR.K
43							0.01	1	R.IQEIEQLDVTSEYK.E
44							4.6e-09	1	R.ALHLQGVLLADAVAVTHGPK.G
45							0.00022	1	R.ALHLQGVLLADAVAVTHGPK.G + Oxidation (H)
46							3.9	1	R.ALHLQGVLLADAVAVTHGPK.G + Oxidation (H)
47							0.00018	1	R.ALHLQGVLLADAVAVTHGPK.G + 2 Oxidation (H)
48							0.0038	1	R.KPLVTIAEDVDGEALSTLVLR.L
49							0.00038	1	R.KPLVTIAEDVDGEALSTLVLR.L
50							4.5e-07	1	R.KPLVTIAEDVDGEALSTLVLR.L
51							2.8	1	R.TALLDAAVASLITAEVVTVEIPK.E
52							0.02	1	R.TALLDAAVASLITAEVVTVEIPK.E
53							3.4e-05	1	K.LVQVAMTHREAGDGTITATVLAR.S

Top scoring peptide matches to query 45  
 Score greater than 32 indicates homology  
 Score greater than 45 indicates identity  
 Status bar shows all hits for this peptide

Score Delta Hit Protein Peptide  
 101.5 -0.10 1 A32800 K.ISSIQTIVPALEIANHR.K  
 13.6 -0.10 K.TLHDELEIEGK.F  
 13.5 -0.17 R.TALPDAILAIRVPLTR.E  
 11.1 0.01 K.TEVTIYCIAPVGEAPPEAR.A  
 10.2 -0.06 K.OLVLAQAGGTHLEIVNDR.L  
 6.5 -0.16 R.LAPTOALGEVITQPTILR.E  
 5.5 -0.08 K.HUTDELIAKVPVTILR.V  
 5.1 -0.01 R.RITGCEGAPQNTAROR.O

MASCOT : Introduction

© 2007-2012 Matrix Science

MATRIX SCIENCE

For each spectrum, there may be multiple possible peptide matches. This particular Mascot report uses a pop-up window to show the alternative peptide matches to each spectrum. In this case, the top match has a high score and the other matches are low scoring, random matches, so no ambiguity. In other cases, the top two or three matches may all be interesting, such as a phosphopeptide where there are several potential phosphorylation sites and moving the phosphate from one site to another only changes the score slightly.

## MS/MS Ions Search

**Easily automated for high throughput**

**Can get matches from marginal data**

**Can be slow**

- No enzyme
- Many variable modifications
- Large database
- Large dataset

**MS/MS is peptide identification**

Proteins by inference.


To summarise, searching of uninterpreted MS/MS data is readily automated for high throughput work. Most “proteomics pipelines” use this approach.

It offers the possibility of getting useful matches from spectra of marginal quality, where it would not be possible to call a reliable sequence tag. Imagine a weak or noisy spectrum that gets a match with a poor score. In isolation, this might be insufficient evidence for the presence of a protein. But, if there are other, similar quality spectra with matches to the same peptide or to other peptides from the same protein, taken together and with the right safeguards, they can provide a degree of confidence that the protein has been identified.

On the down side, such searches can be slow. Particularly if performed without enzyme specificity or with several variable modifications.

Finally, always remember that it is peptides that are being identified, not proteins. From the peptides that have been identified, we try to infer which proteins were present in the sample.

	PMF	MS/MS
Information content	20 to 200 mass values	20 to 200 mass values
Boundary condition	Single protein sequence	Single peptide sequence
Cleavage specificity	Enzyme	Gas-phase dissociation
Major unknown	Protein length	Fragmentation channels
Unique strength	Shotgun protein identification	Residue level characterisation

**MASCOT** : Introduction
 © 2007-2012 Matrix Science


To complete this overview, I'd like to compare the fundamental characteristics of database searching using MS data versus MS/MS data.

The mass spectrum of a tryptic digest of a protein of average size might contain 50 peptide masses, not dissimilar from the MS/MS spectrum of an average sized tryptic peptide. Thus, the "information content" of the individual spectra is similar. The reason MS/MS searches are perceived to be more powerful is mainly that the data set often contains many spectra, multiplying the information content. However, at the single spectrum level, there is little to choose.

In a peptide mass fingerprint, the boundary condition is that the peptides all originate from a single protein. In an MS/MS search, the boundary condition is that the fragments all originate from a single peptide. The weakness of the peptide mass fingerprint is that this boundary condition is often violated, and the spectrum actually represents the digest products of a protein mixture. The MS/MS boundary condition can also be violated, when we analyse co-eluting, isobaric peptides. If this happens, and we have a mixture, the MS/MS search is just as likely to fail as the PMF. We tend not to notice this, because there are many reasons why spectra fail to get matches, such as unsuspected modifications or incorrect precursor mass or charge. We don't normally investigate these failures, so don't see that some of these are due to acquiring a mixed MS/MS spectrum.

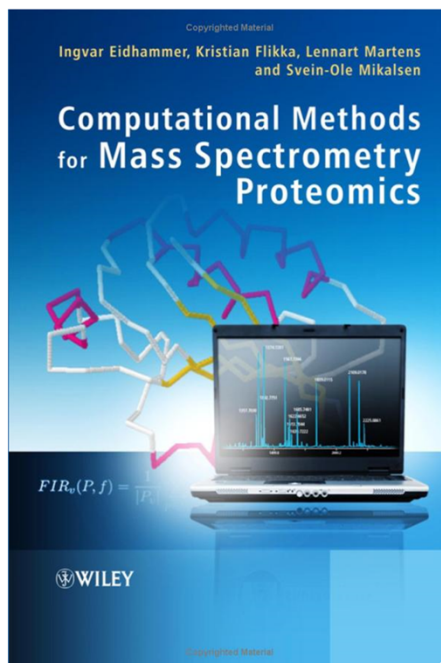
In the peptide mass fingerprint, the specificity comes from the predictable cleavage behaviour of the proteolytic enzyme. Thus, we want an enzyme with a low cutting frequency, such as trypsin. In the MS/MS ions search, the specificity comes from the mostly predictable gas-phase fragmentation behaviour of peptide molecular ions.

In a PMF, we tend to be unsure of the protein mass. Even if the sample is a spot off a 2D gel, there is no guarantee that the database sequence corresponds to the fully processed protein. For MS/MS, we tend to be unsure which types of fragment are present. There is often a dependency on the peptide sequence or a modification. We might get b ions or y ions or a combination of b and y ions. There may be neutral losses and multiple charge states.

Arguably, the major strength of PMF is that it really is shotgun protein identification. The higher the coverage,

the more confident one can be that the protein in the database is the one in the sample. The unique strength of searching MS/MS data is that one gets residue level information. A good match can reveal the presence and location of post translational modifications, which is simply not possible with a PMF.

## Further Reading



**MASCOT** : *Introduction*

© 2007-2012 Matrix Science

**MATRIX**  
SCIENCE

Finally, if you are looking for a recommendation for a text book, this one is fairly recent and covers the whole field clearly and systematically. It isn't just a loose collection of research papers.