

# Introduction to Database Searching using MASCOT

**MASCOT**

 **MATRIX  
SCIENCE**

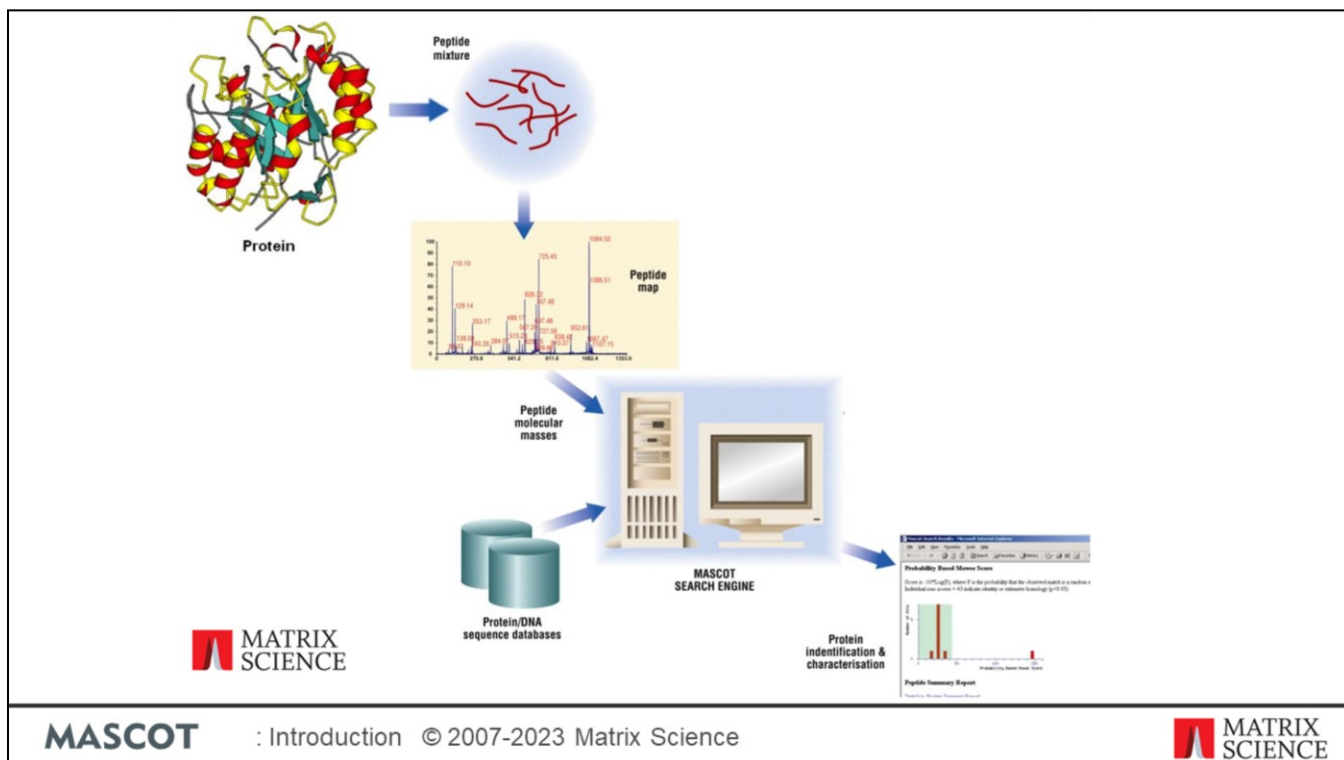
This presentation introduces the topics we will discuss in depth in subsequent talks. It assumes a working knowledge of protein chemistry and mass spectrometry, but no experience of database searching.

# Three ways to use mass spectrometry data for protein identification

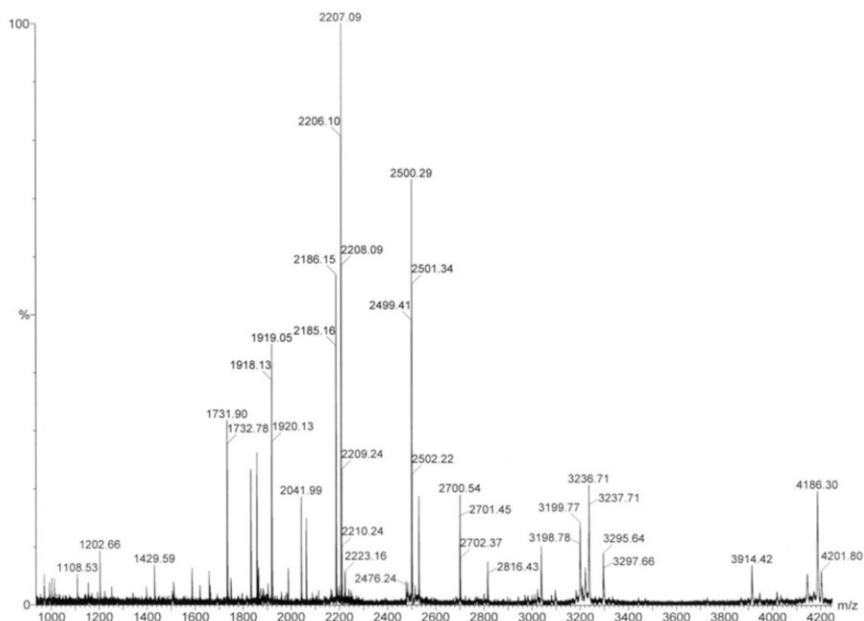
## 1. Peptide Mass Fingerprint

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.



Peptide mass fingerprinting can only be used with a pure protein or a very simple mixture, so the starting point will often be a spot off a 2D gel. The protein is digested with an enzyme of high specificity; usually trypsin, because this is reliable and inexpensive and produces peptides of a suitable size, but any specific enzyme can be used. The resulting mixture of peptides is analysed by mass spectrometry. This yields a set of molecular mass values, which are searched against a database of protein sequences using a search engine. For each entry in the protein database, the search engine simulates the known cleavage specificity of the enzyme, calculates the masses of the predicted peptides, and compares the set of calculated mass values with the set of experimental mass values. Some type of scoring is used to identify the entry in the database that gives the best match, and a report is generated. I will discuss the subject of scoring in detail later.



**MASCOT**

: Introduction © 2007-2023 Matrix Science

**MATRIX  
SCIENCE**

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

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

MS-Fit (Protein Prospector): <https://prospector.ucsf.edu/prospector/mshome.htm>

PeptideMass (ExPASy): [https://web.expasy.org/peptide\\_mass/](https://web.expasy.org/peptide_mass/)

SpectrumMill: <https://proteomics.broadinstitute.org/millhtml/msfit.htm>



. Mowse, PeptideSearch, Profound (Prowl), Protocall, Aldente,  
Xproteo, Bupid, MassSearch

**MASCOT**

: Introduction © 2007-2023 Matrix Science



These presentations will focus on Mascot, but you should be aware that there are a few 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

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.

Concise Summary Report (/dal/ x + - □ ×

localhost/mascot/cgi/m/ 67% ☆ >> ☰

**Mascot Search Results**

User : Low Score  
Email : low@res.edu  
Search title : Fig1d001.mgf  
MS data file : Fig1d001.mgf  
Database : SwissProt 2021\_04 (565928 sequences; 204173280 residues)  
Timestamp : 24 May 2022 at 15:15:23 GMT  
Top Score : 83 for [BIP\\_YEAST](#), Endoplasmic reticulum chaperone BIP OS=Saccharomyces cerevisiae (strain ATCC 204508 / S288C) OS=

SwissProt [Details](#)

Protein hits above identity threshold : 1 0  
Highest scoring protein hit : 83 51

**Mascot Score Histogram**

Protein score is  $-10 \cdot \log(P)$ , where P is the probability that the observed match is a random event.  
Protein scores greater than 70 are significant ( $p < 0.05$ ).

**Concise Protein Summary Report**

Format As: Concise Protein Summary [Help](#)  
Significance threshold  $p < 0.05$  Max. number of hits: AUTO  
Preferred taxonomy: All entries

[Re-Search All](#) [Search Unmatched](#)

- [BIP\\_YEAST](#) Mass: 74422 Score: 83 Expect: 0.003 Matches: 14  
Endoplasmic reticulum chaperone BIP OS=Saccharomyces cerevisiae (strain ATCC 204508 / S288C) OS=SV01  
[CUGP\\_BUCOE](#) Mass: 18253 Score: 45 Expect: 20 Matches: 5  
Co-chaperonin GroES OS=Buchnera aphidicola subsp. Hypos persicae OS=98795 OS=SV01  
[SPER\\_CLOST](#) Mass: 31660 Score: 42 Expect: 40 Matches: 7  
Polyamine aminopropyltransferase OS=Clostridium tetani (strain Massachusetts / E88) OS=212717 OS=SV01  
[BIP\\_CANDIDA](#) Mass: 73285 Score: 40 Expect: 59 Matches: 9  
Endoplasmic reticulum chaperone BIP OS=Candida glabrata (strain ATCC 2001 / CBS 138 / JCH 3762 / HIRC 0622 / HIRL Y-65) OS=284

MASCOT : Introduction © 2007-2023 Matrix Science

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.

7

## 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.

If you want to learn more about the origins and development of peptide mass fingerprinting, I can 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 principle, 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.

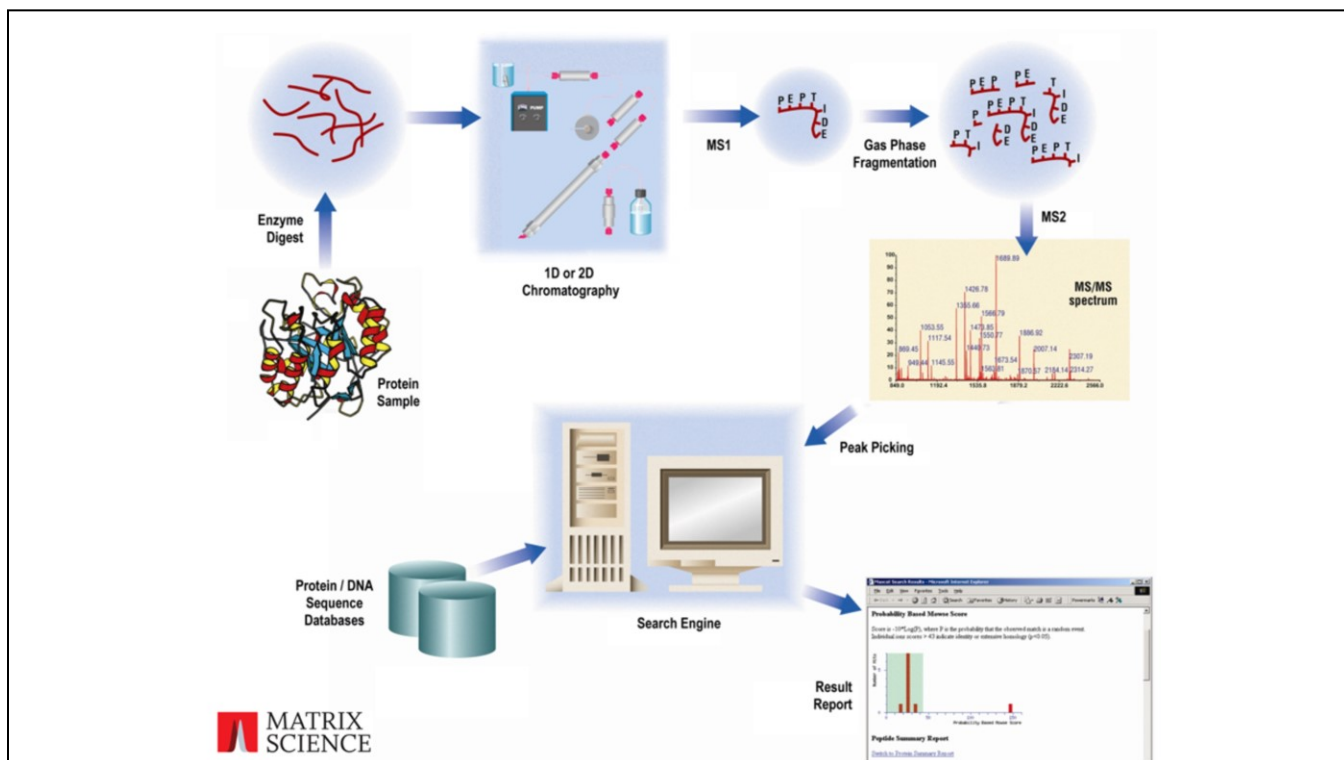
One of the limitations is that you need a database of proteins or nucleic acid sequences that are equivalent to proteins, e.g. mRNAs. In most cases, you will not get satisfactory results from an EST database, where most of the entries correspond to protein fragments, or genomic DNA, where there is a 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.

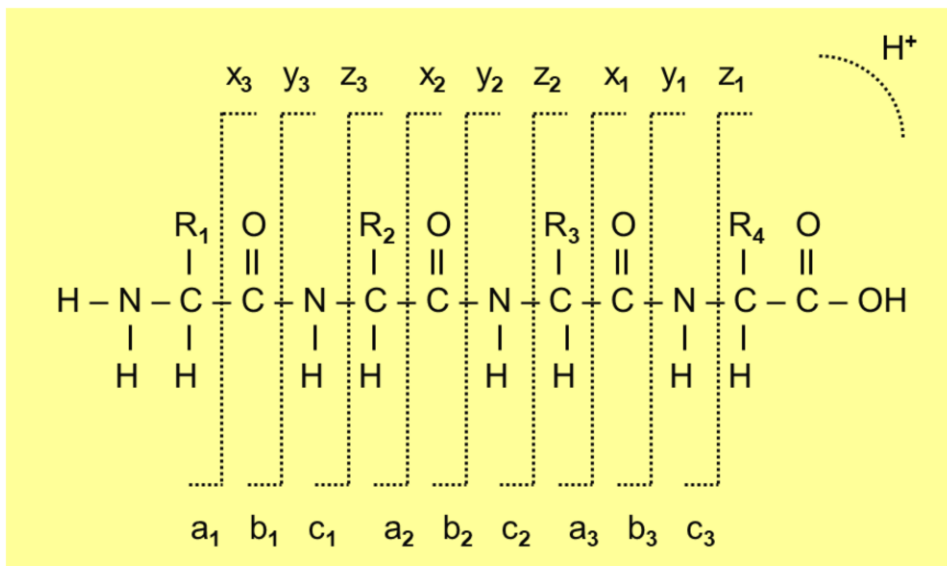
Again, we start with protein, which can now be a single protein or a complex mixture of proteins. We use an enzyme such as trypsin to digest the proteins to peptides. If it is a complex mixture, such as a whole cell lysate, we will probably use one or more stages of chromatography to regulate the flow 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.



➤Roepstorff, P. and Fohlman, J. (1984). *Proposal for a common nomenclature for sequence ions in mass spectra of peptides*. Biomed Mass Spectrom 11, 601.

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.

$$\begin{array}{c} \text{R1} \quad \text{O} \quad \text{R2} \\ | \quad || \quad | \\ \text{H}_2\text{N}-\text{C}-\text{C}-\text{N}^+=\text{C} \\ | \quad | \quad | \\ \text{H} \quad \text{H} \quad \text{H} \end{array}$$

**a<sub>2</sub>**

$$\begin{array}{c} \text{R3} \quad \text{O} \quad \text{R4} \\ | \quad || \quad | \\ \text{O}^+=\text{C}-\text{N}-\text{C}-\text{N}-\text{C}-\text{COOH} \\ | \quad | \quad | \quad | \\ \text{H} \quad \text{H} \quad \text{H} \quad \text{H} \end{array}$$

**x<sub>2</sub>**

$$\begin{array}{c} \text{R2} \quad \text{O} \quad \text{R3} \\ | \quad || \quad | \\ \text{H}_2\text{N}-\text{C}-\text{C}-\text{N}-\text{C}-\text{C}\equiv\text{O}^+ \\ | \quad | \quad | \quad | \\ \text{H} \quad \text{H} \quad \text{H} \quad \text{H} \end{array}$$

**Internal**

$$\begin{array}{c} \text{R3} \\ | \\ \text{H}_2\text{N}^+=\text{C} \\ | \\ \text{H} \end{array}$$

**Immonium**

$$\begin{array}{c} \text{R1} \quad \text{O} \quad \text{R2} \\ | \quad || \quad | \\ \text{H}_2\text{N}-\text{C}-\text{C}-\text{N}-\text{C}-\text{C}\equiv\text{O}^+ \\ | \quad | \quad | \quad | \\ \text{H} \quad \text{H} \quad \text{H} \quad \text{H} \end{array}$$

**b<sub>2</sub>**

$$\begin{array}{c} \text{R3} \quad \text{O} \quad \text{R4} \\ | \quad || \quad | \\ \text{H}_3\text{N}^+-\text{C}-\text{N}-\text{C}-\text{N}-\text{C}-\text{COOH} \\ | \quad | \quad | \quad | \\ \text{H} \quad \text{H} \quad \text{H} \quad \text{H} \end{array}$$

**y<sub>2</sub>**

$$\begin{array}{c} \text{R1} \quad \text{O} \quad \text{R2} \quad \text{O} \\ | \quad || \quad | \quad || \\ \text{H}_2\text{N}-\text{C}-\text{C}-\text{N}-\text{C}-\text{C}-\text{NH}_3^+ \\ | \quad | \quad | \quad | \\ \text{H} \quad \text{H} \quad \text{H} \quad \text{H} \end{array}$$

**c<sub>2</sub>**

$$\begin{array}{c} \text{R3} \quad \text{O} \quad \text{R4} \\ | \quad || \quad | \\ \text{C}^+-\text{C}-\text{N}-\text{C}-\text{COOH} \\ | \quad | \quad | \\ \text{H} \quad \text{H} \quad \text{H} \end{array}$$

**z<sub>2</sub>**

**Sequence Ions**

$$\begin{array}{c} \text{R1} \quad \text{O} \quad \text{CHR}^+ \\ | \quad || \quad | \\ \text{H}_2\text{N}-\text{C}-\text{C}-\text{N}-\text{C} \\ | \quad | \quad | \\ \text{H} \quad \text{H} \quad \text{H} \end{array}$$

**d<sub>2</sub>**

$$\begin{array}{c} \text{O} \quad \text{R4} \\ || \quad | \\ \text{HN}=\text{C}-\text{C}-\text{N}-\text{C}-\text{COOH} \\ | \quad | \quad | \\ \text{H} \quad \text{H} \quad \text{H} \end{array}$$

**v<sub>2</sub>**

**Satellite Ions**

$$\begin{array}{c} \text{CHR}^+ \quad \text{O} \quad \text{R4} \\ || \quad || \quad | \\ \text{C}-\text{C}-\text{N}-\text{C}-\text{COOH} \\ | \quad | \quad | \\ \text{H} \quad \text{H} \quad \text{H} \end{array}$$

**w<sub>2</sub>**

Ion Type	Neutral M <sub>r</sub>
a	[N]+[M]-CHO
a*	a-NH <sub>3</sub>
a <sup>+</sup>	a-H <sub>2</sub> O
b	[N]+[M]-H
b*	b-NH <sub>3</sub>
b <sup>+</sup>	b-H <sub>2</sub> O
c	[N]+[M]+NH <sub>2</sub>
d	a - partial side chain
v	y - complete side chain
w	z - partial side chain
x	[C]+[M]+CO-H
y	[C]+[M]+H
y*	y-NH <sub>3</sub>
y <sup>+</sup>	y-H <sub>2</sub> O
z	[C]+[M]-NH <sub>2</sub>

➤ Papayannopoulos, IA, *The interpretation of collision-induced dissociation tandem mass spectra of peptides*. Mass Spectrom. Rev., 14(1) 49-73 (1995).

**MASCOT**

: Introduction © 2007-2023 Matrix Science



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

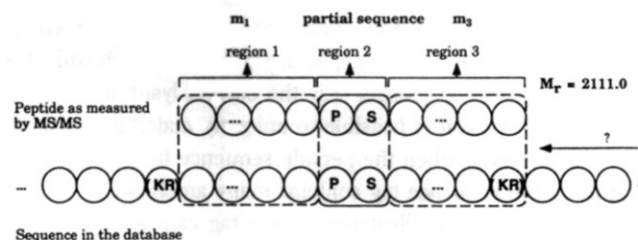
### 1. Peptide Mass Fingerprint

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.



**Figure 1.** Principle of matching peptide sequence tags to a proposed sequence. The upper chain of amino acids represents the peptide sequence as measured by MS/MS (from Table 1 in this example), and the lower chain represents amino acids in the sequence database that the tag is compared to. Note that the partial sequence divides the peptide into three regions. The added mass  $m_1$  of the residues in region 1, together with the N-terminus, is a match criterion as is the added mass in region three,  $m_3$ . In region 2, the sequence is known. Furthermore, it can be required that the peptide obey the cleavage condition of the proteolytic enzyme, marked by KR for trypsin. The left pointing arrow indicates that both search directions may have to be considered.

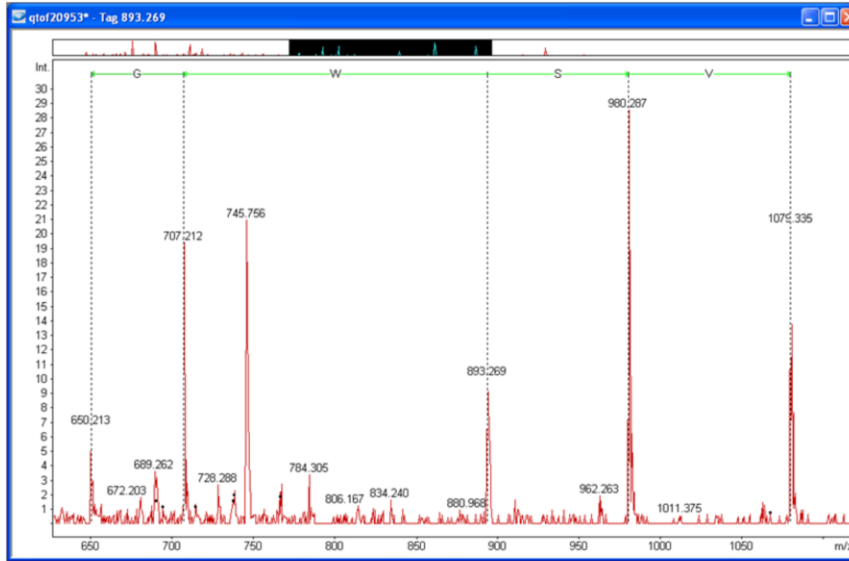
➤Mann, M. and Wilm, M., *Error-tolerant identification of peptides in sequence databases by peptide sequence tags*. Anal. Chem. 66 4390-9 (1994).

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.





1490.4 tag(650.2,GWSV,1079.3)

**MASCOT**

: Introduction © 2007-2023 Matrix Science

**MATRIX  
SCIENCE**

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

MS-Seq (Protein Prospector): <https://prospector.ucsf.edu/prospector/mshome.htm>



PeptideSearch, Multident

**MASCOT**

: Introduction © 2007-2023 Matrix Science



There are a few 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:  Email:

Search title:

Database(s):  Enzyme:

Allow up to:  missed cleavages Quantitation:

Taxonomy:

Fixed modifications:  Display all modifications: ☐

Variable modifications:

Peptide tol.  $\pm$   Da  $\#$   C  $\#$   MS/MS tol.  $\pm$   Da

Peptide charge:  Monoisotopic ☐ Average ☐

Query:

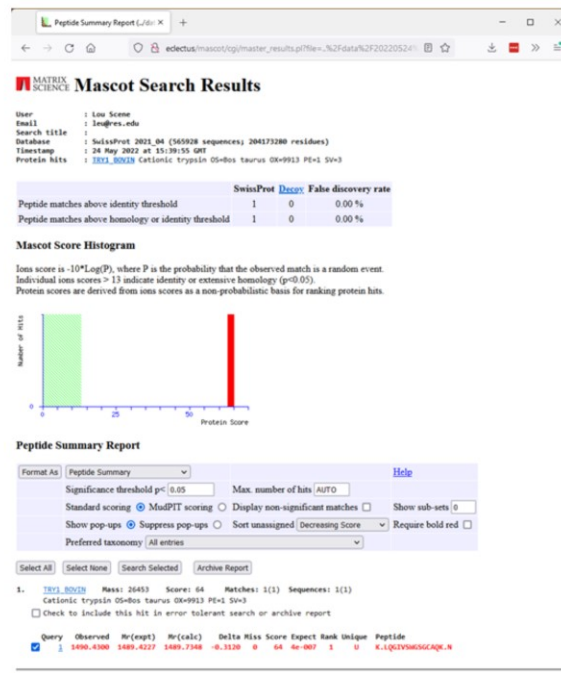
Instrument:  Decoy: ☒

**MASCOT**

: Introduction © 2007-2023 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: LQGVSWGSGCAQK 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.

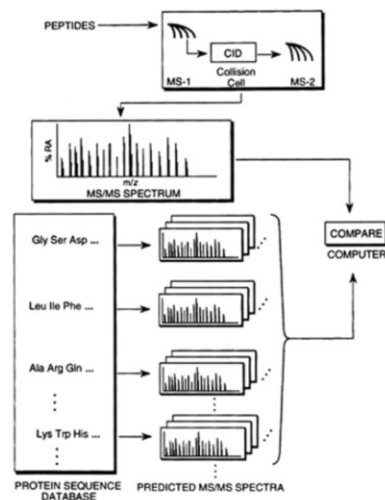


Figure 1. Flow chart that depicts the algorithm for searching protein databases with tandem mass spectrometry data.

## SEQUEST

➤Eng, J. K., McCormack, A. L. and Yates, J. R., 3rd., *An approach to correlate tandem mass spectral data of peptides with amino acid sequences in a protein database.* J. Am. Soc. Mass Spectrom. 5 976-89 (1994)

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 / MS-GFDB	<a href="https://proteomics.ucsd.edu/ProteoSAFe/">https://proteomics.ucsd.edu/ProteoSAFe/</a>
Mascot	<a href="https://www.matrixscience.com/search_form_select.html">https://www.matrixscience.com/search_form_select.html</a>
MS-Tag (Protein Prospector)	<a href="https://prospector.ucsf.edu/prospector/mshome.htm">https://prospector.ucsf.edu/prospector/mshome.htm</a>
PepFrag (Prowl)	<a href="https://prowl.rockefeller.edu/prowl/pepfrag.html">https://prowl.rockefeller.edu/prowl/pepfrag.html</a>
RAId_DbS	<a href="https://www.ncbi.nlm.nih.gov/CBBResearch/qmbp/RAId_DbS/index.html">https://www.ncbi.nlm.nih.gov/CBBResearch/qmbp/RAId_DbS/index.html</a>

Not on-line	Andromeda, Byonic, Comet, greylag, Morpheus, Myrimatch, MSFragger, OMSSA, Paragon, Peaks DB, PepSplice, pFind, Phenyx, ProbiD, ProLuCID, ProteinLynx GS, Sequest, SIMS, SpectrumMill, Tide, X!Tandem (The GPM), etc. etc.
-------------	---

**MASCOT**

: Introduction © 2007-2023 Matrix Science



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:  Email:

Search title:

Database(s):

Taxonomy:

Enzyme:  Allow up to:  missed cleavages

Quantitation:

Crosslinking:

Fixed modifications:

Display all modifications: ☐

Variable modifications:

Peptide tol.:  ppm   $\delta^{13}C$

MS/MS tol.:  ppm

Peptide charge:

Monoisotopic: ☒ Average: ☐

Data file:  No file selected.

Data format:

Instrument:

Precursor:

Error tolerant: ☒

Decoy: ☒

Target PSM FDR:

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.

[illegible]

MS/MS Example (Mascot Search)

Protein family 1 (out of 1)

10 per page 1 Expand all Collapse all

Accession contains Find

▼1 CH60\_HUMAN 1109 60 kDa heat shock protein, mitochondrial OS=Homo sapiens GR=HSPD1 PE=1 SV=2

1.1 cCH60\_HUMAN 1109 61016 30 (30) 18 (18) 60 kDa heat shock protein, mitochondrial OS=Homo sapiens GR=HSPD1 PE=1 SV=2

▼30 peptide matches (25 non-duplicate, 5 duplicate)

Auto-fit to window

Query Digests	Observed	Mr (exp1)	Mr (calc)	Delta M	Score	Expect	Rank	U	Peptide
d11	417.1822	832.3498	832.3828	-0.0329	0	43	0.0096	1	U K.AMPQDNR.E
d12	422.7433	843.4720	843.5066	-0.0346	0	46	0.014	1	U K.VGRVTVVK.S
d13	430.7328	859.4510	859.4837	-0.0327	0	36	0.048	1	U K.IPAHTAE.K + Oxidation (O)
d15	451.2499	900.4853	900.5280	-0.0428	0	52	0.008	1	U K.LSDQVWLE.V
d16	456.7806	911.5447	911.5804	-0.0337	0	59	0.0011	1	U K.VQLQVWVE.A
d21	480.7447	959.4748	959.5036	-0.0288	0	45	0.0099	1	U R.VDMLAKTS.A
d24	595.7855	1189.5565	1189.6012	-0.0447	0	57	0.0019	1	U K.EIGNIISDNR.E
d25	603.7720	1205.5294	1205.5942	-0.0648	0	60	2.7e-005	1	U K.EIGNIISDNR.E + Oxidation (O)
d26	608.3099	1214.6052	1214.6507	-0.0455	0	73	1.5e-005	1	U K.NAGVDSLVSE.I
d27	617.2807	1232.5569	1232.5885	-0.0316	0	81	8.5e-006	1	U K.VQSTQVWNR.E
d31	672.4575	1343.4605	1343.7085	-0.0440	0	64	0.0017	1	U R.PVLEKQWQVE.V
d35	714.8938	1427.7730	1427.8058	-0.0327	0	73	1.4e-005	1	U R.DNGLAVDAVLE.E
d36	722.8849	1443.7552	1443.8007	-0.0455	0	75	2.7e-006	1	U R.DNGLAVDAVLE.E + Oxidation (O)
d39	752.8643	1503.7141	1503.7490	-0.0349	0	90	9.8e-008	1	U K.TUNDELEIKNR.F
d40	760.8461	1519.6777	1519.7439	-0.0662	0	89	1.7e-007	1	U K.TUNDELEIKNR.F + Oxidation (O)
d45	640.3281	1917.9625	1918.0636	-0.1010	0	102	3.9e-009	1	U K.ISSISQVFALEIARNR.E
d46	940.0327	1918.0509	1918.0636	-0.0127	0	87	1.8e-007	1	U K.ISSISQVFALEIARNR.E
d48	1019.5106	2037.0067	2037.0153	-0.0087	0	52	0.00068	1	U R.TQELIQGLVTTSEVE.E
d51	1057.0537	2112.0929	2112.1323	-0.0394	0	116	3.9e-011	1	U R.ALMQVOLLADAVVWNR.G
d52	1045.0399	2128.0453	2128.1272	-0.0619	0	72	8e-007	1	U R.ALMQVOLLADAVVWNR.G + Oxidation (O)
d54	1073.0677	2144.0919	2144.1221	-0.0412	0	93	4.1e-009	1	U R.ALMQVOLLADAVVWNR.G + 2 Oxidation (O)
d59	1183.1570	2364.2994	2364.3264	-0.0270	0	65	1.3e-005	1	U R.RFLVLIARDVGRALSTVLVLR.L
d60	789.1094	2364.3043	2364.3264	-0.0201	0	95	5.9e-009	1	U R.RFLVLIARDVGRALSTVLVLR.L

Query 45

Score > 40 indicates identity

Score > 30 indicates homology

Score	Expect	Rank	U	Peptide
0.0110	0	11	4.3	U K.ISSISQVFALEIARNR.E
-0.1197	1	3	25	U K.ISSISQVFALEIARNR.E
-0.0548	0	2	35	U K.ISSISQVFALEIARNR.E
0.0881	0	2	35	U K.ISSISQVFALEIARNR.E
-0.0449	1	2	36	U K.ISSISQVFALEIARNR.E
-0.1547	0	1	43	U K.ISSISQVFALEIARNR.E
0.0007	1	1	44	U K.ISSISQVFALEIARNR.E
-0.0026	1	1	45	U K.ISSISQVFALEIARNR.E
-0.1111	0	0	53	U K.ISSISQVFALEIARNR.E

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 if:**

- 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.

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.

Another approach we can take with MS/MS peak lists is to search spectral libraries.

## Spectral Library Search

Hertz, H. S, Hites, R. A., and Biemann, K (1971)

*Identification of mass spectra by computer-searching a file of known spectra*  
Anal. Chem. 43, 6, 681-691

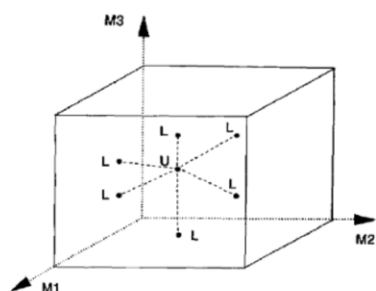


Figure 2. Point representation of library search results (L) for a hypothetical three-peak unknown (U) spectrum (masses M1, M2, and M3)

Stein, S. and Scott, D. R. (1994).  
*Optimization and testing of mass spectral library search algorithms for compound identification*, J. Am. Soc. Mass Spectrom., 5, 859-66

MASCOT

: Introduction © 2007-2023 Matrix Science

MATRIX  
SCIENCE

Spectral library searching of small molecules has been around for a long time as this Biemann paper from 1971 attests, but only more recently used for proteomics and peptide identification.

The earliest one was NIST MS Search, introduced in 1994, which also provided mass spectral libraries to search.

## MS/MS Spectral Library Search Servers on the Web

**SpectraST Spectrum Library Search** <https://www.peptideatlas.org/spectrast/>

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

**Protein Prospector Batch-Tag** <https://prospector.ucsf.edu/prospector/mshome.htm>

### Not on-line

BiblioSpec, COSS, Epsilon-Q, Mistle, MS PepSearch, NIST MS Search, PEAKS Studio, Pepitome, PeptideAtlas, Pmatch, Progenesis, Scaffold, Skyline, SpectraST, etc

**MASCOT**

: Introduction © 2007-2023 Matrix Science



Initially, spectral library searching was not as popular as other methods because the peptide libraries were incomplete. Around 2006 things started to change with mass spec data becoming available on a large scale. Multiple spectral library search engine papers were published that year.

Most spectral library search engines are only available to install and run locally. Of the two originally on the web, SpectraST and X!Hunter, only SpectraST is still available.

Mascot integrates NIST MS PepSearch, which is available for use on our public website.

There are number of other SL search engines that are available for download and local searching.

## MASCOT MS/MS Ions Search

Your name: richard  Email: richardj@matrixscience.com

Search title:

Database(s): NIST\_Human\_HCD (SL)

Taxonomy: All entries

Enzyme: Trypsin  Allow up to: 1  missed cleavages

Quantitation: None

Crosslinking: None

Fixed modifications: --- none selected ---

Mus\_musculus\_GRCm39\_geno  
 Test\_DNA  
 ZEST\_human  
 Spectral library (SL)  
 MassIVE\_HumanHCD  
 NIST\_Human\_IonTrap  
 NIST\_Rat\_QToF  
 PiM\_PiZ\_SL  
 PRIDE\_Contaminants  
 PRIDE\_Human

6C-CysPAT (C)  
 6C-CysPAT (N-term)  
 Acetyl (K)

**MASCOT**

: Introduction © 2007-2023 Matrix Science

**MATRIX  
SCIENCE**

The input and output of a spectral library search uses the same input form and report format as MS/MS Ions Search. The Mascot search engine and NIST MS PepSearch are integrated. Combined searches can be run by selecting a spectral library database and a FASTA database in the search form.

Spectral libraries are in section separate from sequence databases.

Most search parameters – modifications, enzyme, missed cleavages, taxonomy, and instrument – simply don't apply to a library search, although they are used for the database search portion.

The precursor and fragment ion tolerances apply to both.

**Protein Family Summary**

Format: Significance threshold: 300 Max. number of families: AUTO [help]  
 Display non-sig. matches: ☐ Min. number of sig. unique sequences: 1  
 Dendrograms cut at: 0

**Sensitivity**

Proteins (976) [Report Builder] [Unassigned (3884)] [permalink]

**Protein families 1-10 (out of 929)**

10 per page 1 2 3 4 5 6 ... 93 [Next] [Expand all] [Collapse all]

Accession: contains: Find Clear

**▼ 1 KPYK1\_YEAST 24012** Pyruvate kinase 1 OS=Saccharomyces cerevisiae (strain ATCC 204508 / S288c) OX=5592...

**1.1 KPYK1\_YEAST** Score: 24012 Mass: 54510 Matches: 68 (68) Sequences: 43 (43) Pyruvate kinase 1 OS=Saccharomyces cerevisiae (strain ATCC 204508 / S288c) OX=559292 GN=CDK19 PE=1 SV=2

**▼ 68 peptide matches (59 non-duplicate, 9 duplicate)**

☒ Auto-fit to window

Query	Dupes	Observed	Mr (expt)	Mr (calo)	ppm	M	Score	Source	Expect	Rank	U	Peptide
d244		371.7316	741.4487	741.4498	-1.38	0	887	SL	6.7e-08	1	U	K.AGLNIVRL
d451		387.7085	773.4025	773.4032	-0.86	0	771	SL	9.7e-07	1	U	K.SVINDNLK
d935		841.3965	2521.1660	2521.1675	-0.59	0	561	SL	0.00012	1	U	R.AEYSDVGNALIDGACVMSLGRTAK.G + Carboxidomethyl
d1077		430.2297	858.4448	858.4447	0.047	0	946	SL	1.7e-08	1	U	R.RVGGQK.D
d1176		435.7791	869.5436	869.5448	-1.32	0	322	SL	0.03	2	U	R.KAGLNIVRL
d1451		453.2757	904.5369	904.5381	-1.37	0	561	SL	0.00012	1	U	K.AKEFGILK.K
d1452		302.5196	904.5370	904.5380	-1.16	0	589	SL	6.4e-05	1	U	K.AKEFGILK.K
d1720		466.2325	930.4505	930.4521	-1.81	0	772	SL	9.5e-07	1	U	K.DMYDYK.N
d1828		471.2864	940.5583	940.5594	-1.09	0	680	SL	7.9e-06	1	U	R.PLAIALQYK.G
d1978		474.2314	946.4482	946.4469	1.33	0	427	SL	0.0027	1	U	K.DMYDYK.N + Oxidation
d1981		474.2517	946.4809	946.4906	-1.72	0	700	SL	5e-06	1	U	K.VDSQWNLK.G
d2261		493.7971	985.5796	985.5808	-1.18	0	782	SL	7.6e-07	1	U	R.TSLIOTIQK.T
d2421		501.7816	1001.5486	1001.5506	-1.96	0	568	SL	0.0001	1	U	R.TANEVQLTK.E
d2724		345.2177	1032.6313	1032.6332	-1.89	0	499	SL	0.00051	1	U	K.AKEFGILK.G
d2831		522.7806	1043.5467	1043.5500	-3.08	0	424	SL	0.0029	1	U	F.QVSLVQDK.T
d2979		529.7866	1057.5587	1057.5589	-0.22	0	554	SL	0.00014	1	U	K.SNLGKQPVIC.A + Carboxidomethyl
d3066	2	535.2899	1068.5652	1068.5643	0.83	0	890	SL	6.3e-08	1	U	R.QVFPVQK.E
d3721		571.8472	1141.6799	1141.6819	-1.79	0	700	SL	5e-06	1	U	R.RSLIOTIQK.T
d3722		581.8476	1141.6808	1141.6817	-0.95	0	796	SL	9.8e-06	1	U	R.SLPIYVQVQK.E

**MASCOT** : Introduction © 2007-2023 Matrix Science

Results are presented in the same format as an MS/MS search with a list of inferred proteins and the peptide matches.

Spectral libraries typically contain only peptide-level information. Mascot maps spectral library matches to a reference sequence database, which enables robust protein inference. Each library entry is assigned to all the accessions from the reference database that contain the peptide sequence, ignoring enzyme specificity. This means that, in favourable cases, protein inference will be just as good as if the matches had been found in a search of the reference sequence database.

When an entry has no reference accession, we use the protein metadata from the spectral library, if present. When the reference file is well chosen, as in the example, few or no spectral library accessions will be visible in the summary report.



Massot database search: Spec1: X    Next 16 example Massot Spec: X

Accession:    contains:    Find    Clear

Auto-fit to window

Query Dupes	Observed	Mr (expt)	Mr (calc)	ppm	M	Score	Source	Expect	Rank	U	Peptide
d5427	452.5922	1354.7549	1354.7567	-1.36	0	527	SL	0.00027	1	U	P.KTNKCTEIVLALR.K
d5702	465.2392	1392.4958	1392.4971	-0.95	0	639	SL	2e-05	1	U	K.NGVNRYVAFIR.T + Oxidation
d5394	500.5497	1496.4272	1496.4296	-1.72	0	412	SL	3.8e-05	1	U	R.NWFSRQVETDK.R
d5395	375.6646	1496.4292	1496.4299	-0.48	0	685	SL	7.1e-06	1	U	R.NWFSRQVETDK.R
d5405	750.9274	1499.8402	1499.8419	-1.16	0	666	SL	1.1e-05	1	U	R.LTSLAVVAGSLAR.T
d5406	500.9544	1499.8415	1499.8421	-0.42	0	462	SL	0.0012	1	U	R.LTSLAVVAGSLAR.T
d5412	753.4172	1500.8199	1500.8235	-2.42	0	665	SL	1.1e-05	1	U	K.YRNCFFILVTR.C + Carbamidomethyl
d5415	501.2814	1500.8224	1500.8234	-0.68	0	667	SL	1.1e-05	1	U	K.YRNCFFILVTR.C + Carbamidomethyl
d5501	759.8380	1517.4415	1517.4433	-1.24	0	666	SL	1.1e-05	1	U	R.RPVNVRQVDAK.I
d5417	575.6394	1723.8963	1723.8992	-1.71	0	579	SL	8.1e-05	1	U	K.GNGLPQGVSLALAEK.D
d5418	842.9558	1723.8969	1723.8992	-1.28	0	470	SL	0.001	1	U	K.GNGLPQGVSLALAEK.D
d5705	583.6620	1747.9641	1747.9718	-4.40	0	496	SL	0.00055	1	U	R.GSLGIEIRAFVLAQK.K
d5707	874.9940	1747.9734	1747.9720	0.83	0	739	SL	2e-06	1	U	R.GSLGIEIRAFVLAQK.K
d5747	880.9423	1759.8701	1759.8739	-2.17	0	649	SL	1.6e-05	1	U	K.IENQGVNNVDELK.V
d5757	937.0048	1871.9950	1871.9991	-2.21	0	681	SL	7.7e-06	1	U	K.SKEELPQGVSLALAEK.G
d57918	425.0043	1871.9970	1871.9993	-1.27	0	657	SL	1.3e-05	1	U	K.SKEELPQGVSLALAEK.G
d57933	426.3417	1876.0633	1876.0667	-1.84	0	657	SL	0.0013	1	U	R.GSLGIEIRAFVLAQK.L
d5136	456.6795	1947.0167	1947.0210	-2.18	0	469	SL	0.001	1	U	K.GNGLPQGVSLALAEK.E
d5195	667.7041	2000.0905	2000.0942	-1.86	0	809	SL	4.1e-07	1	U	R.KSKEELPQGVSLALAEK.G
d5196	1001.0531	2000.0916	2000.0941	-1.25	0	762	SL	1.2e-06	1	U	R.KSKEELPQGVSLALAEK.G
d5197	501.8305	2000.0929	2000.0943	-0.70	0	750	SL	1.6e-06	1	U	R.KSKEELPQGVSLALAEK.G
d5206	670.8016	2007.0100	2007.0158	-2.89	0	696	SL	5.2e-06	1	U	R.PYTFSTVVAAGVAFYQK.A
d5207	1004.5141	2007.0176	2007.0159	0.85	0	593	SL	6.2e-05	1	U	K.PYTFSTVVAAGVAFYQK.A
d5233	1011.4730	2020.9314	2020.9377	-3.12	0	717	SL	3.4e-06	1	U	F.VYKQSPVWMTQVDAK.I
d5550	421.0749	2480.2705	2480.2759	-2.18	0	458	SL	0.0013	1	U	K.GNGLPQGVSLALAEKQKDA.F
d5618	870.1836	2607.2890	2607.2849	1.56	0	676	SL	8.7e-06	1	U	R.NCTPKPTSTTVVAAGVAFYQK.A + Carbamidomethyl
d5646	941.4308	2821.2705	2821.2752	-1.69	0	758	SL	1.3e-06	1	U	R.TQYTFNVDVTFPNNEMFTYQK.Y
d5654	946.7616	2837.2630	2837.2703	-2.57	0	699	SL	5.1e-06	1	U	R.TQYTFNVDVTFPNNEMFTYQK.Y + Oxidation

6927: Score 9457 (rc=4433.54) [C:\Downloads\slc-03130p\_optac\_wstudy\_08011.RAM]  
Score > 300 indicates identity

Score	Ident	SL	U	Peptide
12.2	0	192	SL	0.6 2
-15.4	0	188	SL	0.46 3
-1.27	0	162	SL	1.2 4
-6.56	0	160	SL	1.3 5
-7.20	0	145	SL	1.8 6
13.6	0	139	SL	2 7
20.4	0	129	SL	2.6 8
13.5	0	112	SL	3.8 9
-13.3	0	111	SL	3.9 10

1 G3P3\_YEAST  
2 G3P2\_YEAST

17399 Glyceroldehyde-3-phosphate dehydrogenase 3 OS=Saccharomyces cerevisiae [...]  
14945 Glyceroldehyde-3-phosphate dehydrogenase 2 OS=Saccharomyces cerevisiae [...]

As in MS/MS reports there is an additional dimension to the data. For each spectrum, there may be multiple possible peptide matches which are displayed in the pop-up window.

## Spectral Library Search

**Easily automated for high throughput**

**Very Fast**

**Can only match spectra in the database**

**Predetermined enzyme specificity or fixed and variable modifications**

**Can match peptides in the library with:**

Semi-specific cleavage or different enzyme(s)

Uncommon variable modifications

**Spectral Library search is peptide identification**

Proteins by inference.

**MASCOT**

: Introduction © 2007-2023 Matrix Science

 **MATRIX  
SCIENCE**

Just like MS/MS searches, the process is easily automated for high throughput. This is one of the fastest ways to search uninterpreted MS/MS data. One of the reasons why it is so fast is that the search only considers spectra in the library, which is a smaller search space than traditional database searching.

Most search parameters – modifications, enzyme, missed cleavages, taxonomy, and instrument – simply don't apply to a library search. These are predetermined when the library is created. However, if the library contains peptides that are outside the database search space, such as different enzyme or uncommon variable modifications, then they will automatically be included in the search.

As in MS/MS searching, the spectral library search is identifying peptides not proteins. Unlike an MS/MS search against proteins from a FASTA file, the spectral library contains limited protein information beyond an accession number and the peptide sequence. Instead, a reference library is selected when the library is set up and protein data is taken from there.

	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

To complete this overview, I'd like to compare the fundamental characteristics of database searching using MS data versus MS/MS data. The MS/MS analysis covers both database searching and spectral library searching.

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.

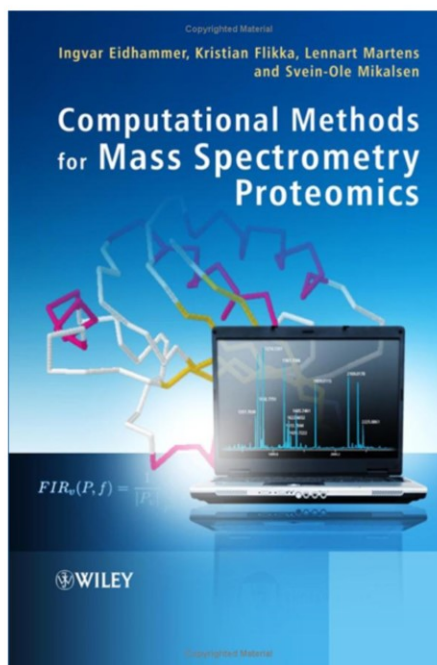
In the peptide mass fingerprint, the specificity comes from the predictable cleavage behaviour of the proteolytic enzyme. Thus, we want an enzyme with a low cutting frequency, such as trypsin. In the MS/MS ions search, the specificity comes from the mostly

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 © 2007-2023 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.