

# Search Parameters

MASCOT



# Search Parameters

The screenshot shows the 'MASCOT Peptide Mass Fingerprint' search form. It includes a search title field, a database selection dropdown (currently set to 'SPMSD\_CapSino'), an enzyme selection dropdown (currently set to 'Trypsin'), and a taxonomy dropdown (currently set to 'All entries'). There are sections for fixed and variable modifications, each with a 'none selected' button and a list of modification options. The form also features a 'Peptide list' section with a 'Peptide list' dropdown and a 'Mass values' section with a 'Search' button and a 'Mass values' dropdown. At the bottom, there is a 'Start Search' button and a 'Reset Form' button.

The screenshot shows the 'MASCOT MS/MS Ions Search' form. It includes a search title field, a database selection dropdown (currently set to 'SPMSD\_CapSino'), an enzyme selection dropdown (currently set to 'Trypsin'), and a taxonomy dropdown (currently set to 'All entries'). There are sections for fixed and variable modifications, each with a 'none selected' button and a list of modification options. The form also features a 'Peptide list' section with a 'Peptide list' dropdown and a 'Mass values' section with a 'Search' button and a 'Mass values' dropdown. At the bottom, there is a 'Start Search' button and a 'Reset Form' button.

In this presentation, we will describe each of the Mascot search parameters.

If you submit a search from a web browser, you have a choice of three different search forms. All three forms submit to the same search engine, but they have been optimised for three different types of search. The form for a peptide mass fingerprint is shown on the left, and the form for a search of uninterpreted MS/MS data on the right. Most of the controls are common to both.

## Search Parameters

The screenshot shows the MASCOT Sequence Query web interface. At the top, there is a navigation bar with the Matrix Science logo and a search bar. Below the navigation bar, the page title is "MASCOT Sequence Query". The main form area contains several sections: "Your name" (set to "infared") and "Email" (set to "infared@matrixscience.com"); "Search file" (empty); "Database(s)" (a list box containing "IP28L\_C\_cysteine", "IP28L\_C\_cysteine", "IP28L\_C\_H12", "IP28L\_C\_cysteine", and "IP28L\_Hydrolyzed\_H12"); "Enzyme" (set to "Trypsin"); "Allow up to" (set to "3") missed cleavages; "Quantification" (set to "None"); "Taxonomy" (set to "All entries"); "Fixed modifications" (set to "none selected"); "Variable modifications" (set to "none selected"); "Peptide list" (set to "1") and "Peptide charge" (set to "1"); "MS/MS list" (set to "0.3") and "MS/MS list" (set to "0.3"); "Query" (a text area); "Instrument" (set to "Default"); and "Query" (a text area). At the bottom of the form, there are "Start Search..." and "Reset Form" buttons. The footer of the page contains the Matrix Science logo and the text "© 2007-2023 Matrix Science".

**MASCOT** : Search Parameters

© 2007-2023 Matrix Science



The third form is for a sequence query, such as a sequence tag search. The controls on this form are very similar to those on the MS/MS form. The main difference is that we have a text area to type in the queries, rather than a data file upload control.

# Help

PMF ✓ SQ ✓ MS/MS ✓

MASCOT Peptide Mass Fingerprint

Your name: [input] Email: [input]

Database(s): [input] Enzyme: Trypsin

Fixed modifications: [input] Allow up to: [input] missed cleavages

Variable modifications: [input] Display all modifications [input]

Protein names: [input] Peptide list: [input]

Mass values: [input] [input] [input] [input] [input] [input] [input] [input] [input] [input]

Peptide list: [input]

Start Search [button] Reset Form [button]

Mascot search parameters [P]

### Modifications

Select any known or suspected modifications.

Mascot supports two types of modification. Fixed modifications are applied universally, to every instance of the specified residue(s) or terminus. There is no computational overhead associated with a fixed modification, it is simply equivalent to using a different mass for the modified residue(s) or terminus. For example, selecting Carboxymethyl (C) means that all calculations will use 161 Da as the mass of cysteine.

Variable modifications are those which may or may not be present. Mascot tests all possible arrangements of variable modifications to find the best match. For example, if Oxidation (O) is selected, and a peptide contains 3 methionines, Mascot will test for a match with the experimental data for that peptide containing 0, 1, 2, or 3 oxidised methionine residues.

Variable modifications can be a very powerful means of finding a match, but there are also dangers to be aware of. Even a single variable modification will generate many possible additional peptides to be tested. More than one variable modification causes the number of arrangements to increase geometrically. This means that a search can take dramatically longer than the same search with fixed modifications. More importantly, testing all possible arrangements of modifications generates many more random matches, so that discrimination can be sharply reduced.

The best advice is to use variable modifications sparingly; never select a large number "just in case". Mascot allows up to 9 variable modifications to be specified but, in most

Click on any link for help

MASCOT : Search Parameters

© 2007-2023 Matrix Science



At the top of each slide, there is a key to show which search parameter applies to which type of search.

The labels on the search form are hyperlinks. Just click on them to get detailed help

## User details and title

PMF✓ SQ✓ MS/MS✓

|              |   |       |   |
|--------------|---|-------|---|
| Your name    | <input type="text" value="Expert User"/>              | Email | <input type="text" value="smartie@matrixscience..com"/> |
| Search title | <input type="text" value="Arabidopsis sample #3476"/> |       |   |

- Search form will ‘remember’ user name and email address in cookie
- If Mascot security is enabled, then this information taken from user database
- Email address used for sending results
- Search title is shown in the report, and can help locate a search in the search log

At the top of the form are a couple of fields for user information. The name and email are saved as a browser cookie when a search is submitted, so you don't need to complete them every time.

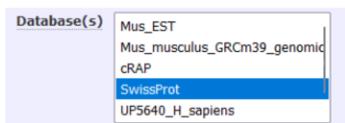
If you have an in-house server, and Mascot security is enabled, these fields will be populated automatically with the details of the user who is logged in.

When you use the Matrix Science public web site, you have to supply a name and email address. This is to allow the results of a search to be returned by email. Usually, search results are returned promptly to your browser window. However, if your connection to the web site is broken before the search is complete, they will be emailed to the supplied address. If you have an in-house server, you can enable this if you wish. It is turned off by default.

The search title is free text. You don't have to enter anything. However, it is a good idea to fill in all of these fields, because it makes it much easier to find your old search results in the search log.

## Database

PMF ✓ SQ ✓ MS/MS ✓



Database(s)

- Mus\_EST
- Mus\_musculus\_GRCm39\_genomic
- cRAP
- SwissProt
- UP5640\_H\_sapiens



Database(s) SwissProt (AA)

- Nucleic acid (NA)
- Mus\_EST
- Mus\_musculus\_GRCm39\_genomic
- Spectral library (SL)
- PRIDE\_Contaminants
- PRIDE\_Human
- Amino acid (AA)
- cRAP
- UP5640\_H\_sapiens
- UP589\_M\_musculus

### Choose the right database

- Swiss-Prot good for PMF and MS/MS of well characterised organisms
- NCBIprot or UniRef100 if you want to search all known protein sequences
- ESTs for MS/MS if genome not sequenced

**MASCOT**

: Search Parameters

© 2007-2023 Matrix Science

**MATRIX**  
SCIENCE

Choosing the right database is so important that there will be a complete presentation on this topic. On the left if the field for the PMF and Sequence Query forms and on the right for the MS/MS form

Very briefly, for a peptide mass fingerprint, search a comprehensive, non-redundant database, like SwissProt. If the data are any good, it won't matter if one or two mass values fail to find matches. The advantage of searching a small database is that the search is fast and the reports are concise.

For MS/MS of a well characterised organism, such as human or mouse or yeast, SwissProt is still a good choice. In other cases, search a comprehensive, non-identical database, where every single peptide is explicitly represented, such as NCBIprot or UniRef100.

If the genome of your organism of interest has not been sequenced, it won't be represented in the protein databases, but there may be lots of Expressed Sequence Tags (ESTs). Not advisable for PMF, because many sequences correspond to protein fragments.

You can select multiple databases for a search. This is particularly useful when you want to search a single organism database and include the sequences of common contaminants, such as BSA and trypsin. One restriction is that you cannot mix AA and DNA databases.

# Taxonomy

PMF✓ SQ✓ MS/MS✓

Taxonomy All entries ▾

- Speeds up the search
- Simplifies the result report
- The drop-down list is easily configurable
- Database Manager keeps taxonomy indexes up to date

If a database contains taxonomy information, we can use this to restrict the search to entries for a particular organism or family. This speeds up the search because, in effect, it makes the database smaller.

Limiting the taxonomy simplifies the result report, because you don't see all the homologous proteins from other species.

The drop down list in the search form is configurable. If you are working on a particular organism, you can easily add this to the list.

It is important that the taxonomy is as accurate as possible and kept up to date, this is maintained automatically by Database Manager. You can find out more about Database Manger in the "Sequence Database Administration" presentation.

# Taxonomy

PMF ✓ SQ ✓ MS/MS ✓

|  | Tax IDs | Count     |
|--|---------|-----------|
| Name = NCBIprot  | 0       | 12868     |
| Filename = NCBIprot_20210320.fasta   | 1       | 323887521 |
| Status = In use  | 2       | 126614020 |
| State Time = Wed May 11 09:36:56 # searches = 0                              | 3       | 27642421  |
| Mem mapped = YES Request to mem map = YES Request unmap = NO Mem locked = NO | 4       | 11073723  |
| Number of threads = -1 Current = YES Type = Amino acid                       | 5       | 10898910  |
| Name = NCBIvertebrate_other  | 6       | 15122345  |
| Filename = NCBIvertebrate_other_20220308.fasta                               | 7       | 533969    |
| Status = In use  | 8       | 263644    |
| State Time = Wed May 11 09:36:56 # searches = 0                              | 9       | 164323    |
| Mem mapped = YES Request to mem map = YES Request unmap = NO Mem locked = NO | 10      | 120357    |
| Number of threads = -1 Current = YES Type = Amino acid                       | 11      | 93580     |
| Name = noXtProt  | 12      | 74055     |
| Filename = noXtProt_2021-02-18.fasta   |         |           |
| Status = In use  |         |           |
| State Time = Wed May 11 09:36:56 # searches = 0                              |         |           |
| Mem mapped = YES Request to mem map = YES Request unmap = NO Mem locked = NO |         |           |
| Number of threads = -1 Current = YES Type = Amino acid                       |         |           |
| Name = NGF   |         |           |
| Filename = NGF_20160224.fasta  |         |           |
| Status = In use  |         |           |
| State Time = Wed May 11 09:36:56 # searches = 0                              |         |           |
| Mem mapped = YES Request to mem map = YES Request unmap = NO Mem locked = NO |         |           |
| Number of threads = -1 Current = YES Type = Amino acid                       |         |           |
| Name = NIST_Human_HCD  |         |           |
| Filename = NIST_Human_HCD_20160503.msp                                       |         |           |
| Status = In use  |         |           |
| State Time = Wed May 11 09:36:56 # searches = 0                              |         |           |
| Mem mapped = NO Request to mem map = NO Request unmap = NO Mem locked = NO   |         |           |

MASCOT

: Search Parameters

© 2007-2023 Matrix Science



From time to time, its a good idea to go to the database status page and check the stats file for each database. The stats file contains lots of useful information, like whether entries contain illegal characters or whether an entry is too long.

It also tells you how good your taxonomy is. Here are the numbers for the NCBI nr database on our web site in March 2022. There are 517 million entries, and 12868 have no taxonomy. In other words, 99.99% of the entries have a taxonomy assigned. If you look at your stats file and see that (say) 10% of the entries have no taxonomy, that's 10% of the entries that are going to be missed whenever you do a search with taxonomy specified.

# Taxonomy

In most cases, if the correct protein is not in the database, you'd like to see the closest match ... whatever the species

PMF ✓ SQ ✓ MS/MS ✓

```
Database statistics
Time files compressed : Mon May 02 12:50:35 2022
Time files compressed (int) : 1651510235
Time / date of fasta file : Mon May 02 12:49:56 2022
Time of fasta files (int) : 1651510196
Number of residues : 204698499
Number of sequences : 566996
Number with invalid residues: 0
Number of sequences too long: 0
Length of longest sequence : 35213
Maximum Accession Length : 11
Version of Mascot : 2.8.1
Version of this file : 6
Type of fasta file : AA
Parse rule for accession : >.[!~]*[!~]*
Seqs with invalid taxon tree: 0
Num sequences for taxonomy : All entries=566996
Num sequences for taxonomy : Archaea (Archaeobacteria)=19694
Num sequences for taxonomy : Eukaryota (eucaryotes)=194882
Num sequences for taxonomy : Alveolata (alveolates)=1194
Num sequences for taxonomy : Plasmodium falciparum (malaria parasite)=381
Num sequences for taxonomy : Other Alveolata=813
Num sequences for taxonomy : Metazoa (Animals)=108571
Num sequences for taxonomy : Caenorhabditis elegans=4353
Num sequences for taxonomy : Drosophila (fruit flies)=5991
Num sequences for taxonomy : Chordata (vertebrates and relatives)=86467
Num sequences for taxonomy : bony vertebrates=8581
Num sequences for taxonomy : lobe-finned fish and tetrapod clade=80318
Num sequences for taxonomy : Mammalia (mammals)=67503
Num sequences for taxonomy : Primates=27145
Num sequences for taxonomy : Homo sapiens (human)=20377
Num sequences for taxonomy : Other primates=6768
Num sequences for taxonomy : Rodentia (Rodents)=27009
Num sequences for taxonomy : Mus.=17167
Num sequences for taxonomy : Mus musculus (house mouse)=17114
Num sequences for taxonomy : Rattus=8162
Num sequences for taxonomy : Other rodentia=1680
Num sequences for taxonomy : Other mammalia=13349
Num sequences for taxonomy : Xenopus laevis (African clawed frog)=3471
```

MASCOT : Search Parameters

© 2007-2023 Matrix Science



A word of warning. Don't specify a very narrow taxonomy in a search. Think carefully about what you are trying to achieve when you do this.

If the correct protein from the correct species is not in the database, wouldn't you want to see a good match to a protein from a similar species?

This is especially important for poorly represented species. For example, look at these numbers for the Swiss-Prot 2022\_01: half a million entries; 27 thousand entries for rodents, but only 1680 are not either mouse or rat. So, even if you are studying hamster or porcupine, you don't want to choose 'Other rodentia'.

## Enzyme

PMF ✓ SQ ✓ MS/MS ✓

Enzyme    
Allow up to   missed cleavages

- **First choice should normally be the enzyme actually used, and 1 missed cleavage**
- **Large number of missed cleavages, try increasing to 2**
- **Use semi-trypsin rather than no enzyme**
- **No enzyme only in exceptional cases, and never for PMF**

MASCOT

: Search Parameters

© 2007-2023 Matrix Science

 MATRIX  
SCIENCE

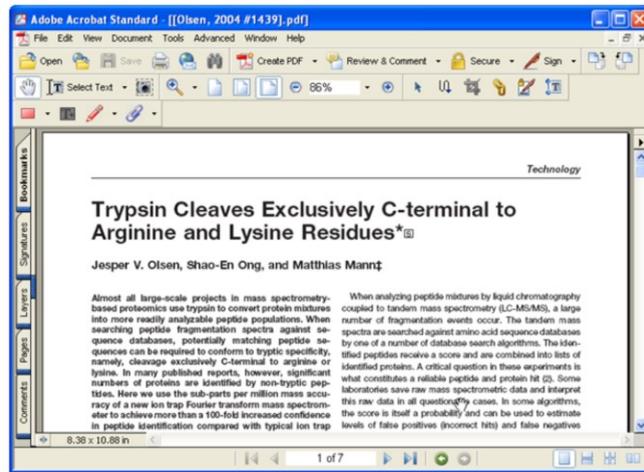
All the search forms have a drop down list for choosing an enzyme. If your peptides come from an enzyme digest, you need to know what the enzyme was and then choose it from the list.

Setting the number of allowed missed cleavage sites to zero simulates a limit digest. If you are confident that your digest is perfect, with no partial fragments present, this will give maximum discrimination and the highest score for a peptide mass fingerprint.

If experience shows that your digest mixtures usually include some partials, that is, peptides with missed cleavage sites, you should choose a setting of 1, or maybe 2 missed cleavage sites. Don't specify a higher number without good reason, because each additional level of missed cleavages increases the number of calculated peptide masses to be matched against the experimental data. In other words, the missed cleavage parameter should be set by looking at some successful search results to see how complete your digests really are.

# Enzyme

PMF ✓ SQ ✓ MS/MS ✓



Olsen, J. V., Ong, S.-E. and Mann, M., *Mol. and Cellular Proteomics*, 3, 608-14 (2004)

**MASCOT** : Search Parameters

© 2007-2023 Matrix Science



Although some people like to perform searches without enzyme specificity, and then gain confidence that a match is correct if the match is tryptic, this isn't a good idea. If there is evidence for a lot of non-specific cleavage, then a semi-specific enzyme allows one end of the peptide to be non-specific, but not both. Only abandon enzyme specificity if you have no other choice, such as when searching endogenous peptides.

You cannot perform a no-enzyme peptide mass fingerprint. It simply won't work, even if you have good mass accuracy.

There is some controversy over the level of non-specific peptides that can be expected in a tryptic digest. Our experience is that the levels of non-specific peptides are very low, less than 3%, unless there is something seriously wrong with the trypsin or the protocol.

Why do we advise so strongly against no-enzyme searches?

## Enzyme

PMF✓ SQ✓ MS/MS✓

Fixed modifications: Carbamidomethyl (C)  
Variable modifications: Oxidation (M)  
Peptide mass tolerance: ± 10 ppm (# 13C = 1)  
Fragment mass tolerance: ± 0.1 Da  
Max missed cleavages: 2  
Instrument type: ESI-TRAP  
Number of queries: 44,894  
Peptide FDR: 1%

| CLE          | candidate peptides | seconds | average identity score | matches above identity | Matches above homology |
|--------------|--------------------|---------|------------------------|------------------------|------------------------|
| Trypsin      | 4.4E6              | 42      | 26                     | 16,767                 | 17,437                 |
| Semi-trypsin | 6.9E7              | 150     | 38                     | 12,732                 | 15,242                 |
| none         | 3.9E8              | 670     | 44                     | 10,681                 | 14,074                 |

**MASCOT**

: Search Parameters

© 2007-2023 Matrix Science



Here are some numbers for an Orbitrap dataset when we search using strict trypsin, semi-specific trypsin, and no enzyme specificity

As you can see, the no enzyme search takes a lot longer and we get fewer reliable matches.

The reason is simple, the search space for a no-enzyme search is much, much larger than for a tryptic search. This means that the thresholds are higher and we lose marginal matches. Unless you have a high level of non-specific peptides, you lose more than you gain.

So, doing a no-enzyme search in Mascot is not a good idea unless there is a very high level of non-specific peptides. Semi-trypsin will be a better choice if the peptides came from a tryptic digest but there is a high level of non-specific cleavage. Only use no enzyme if the peptides are not the products of an enzyme digest, e.g. MHC peptides or endogenous peptides.

# Enzyme

PMF ✓ SQ ✓ MS/MS ✓

| Title         | Sense  | Cleave at | Restrict                 | Independent | Semispecific |             |
|---------------|--------|-----------|--------------------------|-------------|--------------|-------------|
| Trypsin       | C-term | KR        | P                        | no          | no           | Edit Delete |
| TrypsinP      | C-term | KR        |                          | no          | no           | Edit Delete |
| Arg-C         | C-term | R         | P                        | no          | no           | Edit Delete |
| Asp-N         | N-term | BD        |                          | no          | no           | Edit Delete |
| Asp-N_ambic   | N-term | DE        |                          | no          | no           | Edit Delete |
| Chymotrypsin  | C-term | FLWV      | P                        | no          | no           | Edit Delete |
| CNBr          | C-term | M         |                          | no          | no           | Edit Delete |
| CNBr+Trypsin  | C-term | M         | P                        | no          | no           | Edit Delete |
| Formic_acid   | N-term | D         |                          | no          | no           | Edit Delete |
| Lys-C         | C-term | K         | P                        | no          | no           | Edit Delete |
| Lys-C/P       | C-term | K         |                          | no          | no           | Edit Delete |
| LysC+AspN     | N-term | BD        | P                        | no          | no           | Edit Delete |
| Lys-N         | N-term | K         |                          | no          | no           | Edit Delete |
| PepsinA       | C-term | FL        |                          | no          | no           | Edit Delete |
| semiTrypsin   | C-term | KR        | P                        | no          | yes          | Edit Delete |
| TrypChymo     | C-term | FKLRWV    | P                        | no          | no           | Edit Delete |
| TrypsinMSIP1  | N-term | J         | P                        | no          | no           | Edit Delete |
| TrypsinMSIP1  | C-term | J         |                          | no          | no           | Edit Delete |
| TrypsinMSIP1P | N-term | J         |                          | no          | no           | Edit Delete |
| TrypsinMSIP1P | C-term | JKR       |                          | no          | no           | Edit Delete |
| VB-DE         | C-term | BDEZ      | P                        | no          | no           | Edit Delete |
| VB-E          | C-term | EZ        | P                        | no          | no           | Edit Delete |
| NoCleave      | C-term | J         | ABCDEFGHIJKLMNPQRSTUWXYZ | no          | no           | Edit Delete |
| None          | C-term |           |                          |             |              |             |

[Add new enzyme](#) | [Main menu](#)

MASCOT : Search Parameters

© 2007-2023 Matrix Science



The list of enzymes is user configurable. Standard entries are described in the help. If you wish, you can modify the definitions or create new ones using the configuration editor.

Mascot supports two categories of mixed enzyme definitions. An independent mixed enzyme is used where multiple sample aliquots have been digested separately, and the digests combined for analysis. This means that the sample could contain (say) tryptic peptides and Asp-N peptides, but no peptides that are tryptic at one end and Asp-N at the other. The second category simulates a single sample aliquot being digested simultaneously or serially by more than one cleavage agent. For example CNBr followed by trypsin.

Remember that enzyme type None simulates cleavage at every peptide bond. For top down searches, where you don't want any cleavage, choose NoCleave.

# Enzyme

PMF ✓ SQ ✓ MS/MS ✓

The screenshot shows the Matrix Science Mascot web interface. The page title is "MASCOT Sequence Query". The form includes the following fields and options:

- Search & MSID:** A text input field containing "GLU-FIB".
- Database(s):** A dropdown menu with "SWISSPROT" selected.
- Enzyme:** A dropdown menu with "Trypsin" selected.
- Taxonomy:** A dropdown menu with "All entries" selected.
- Fixed modifications:** A dropdown menu with "None selected" selected.
- Variable modifications:** A dropdown menu with "None selected" selected.
- Peptide list:** A text input field containing "1-3" and "1-4".
- Query:** A text input field containing "GLU-FIB".
- Peptide list:** A dropdown menu with "1-3" and "1-4" selected.
- Taxonomy:** A dropdown menu with "All entries" selected.
- Fixed modifications:** A dropdown menu with "None selected" selected.
- Variable modifications:** A dropdown menu with "None selected" selected.
- Peptide list:** A dropdown menu with "1-3" and "1-4" selected.
- Query:** A dropdown menu with "GLU-FIB" selected.

**MASCOT** : Search Parameters

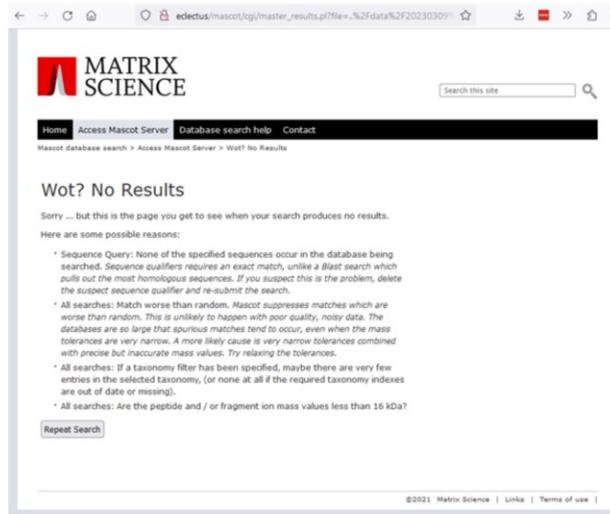
© 2007-2023 Matrix Science



Remember that enzyme specificity also applies to Sequence Queries.

Quite often, we receive a support email along the lines of "Mascot is broken. I did a search for this peptide and I know its in the database but Mascot failed to find it".

For example, here's a search for glu-fib, a very common sequencing standard. The mass is correct and the sequence is correct. But, when we do a search of SwissProt ...



No results!  
Why?

# Enzyme

PMF ✓ SQ ✓ MS/MS ✓

**MASCOT Search Results**

**Protein View: FIBB\_HUMAN**

Fibrinogen beta chain OS=Homo sapiens OX=9606 GN=FGB PE=1 SV=2

Database: SwissProt  
Score: 72  
Expect: 0.037  
Molecular mass (M<sub>r</sub>): 55892  
Calculated pI: 8.54  
Taxonomy: [Homo sapiens](#)

Sequence similarity is available as an [NCBI BLAST search of FIBB\\_HUMAN against nr](#).

**Search parameters**

Enzyme: semiTrypsin: cuts C-term side of KR unless next residue is R.  
Cleavage is semi-specific. (Peptide can be non-specific at one terminus only.)

Protein sequence coverage: 2%

Matched peptides shown in **bold red**.

```
1  HHHVHRRFP  KLTWDELL  LLVCTVYS  QQNNRRSP  PSAGRRFLD
51  HFKKAKIA  PAFVLSGG  TRARAPAA  TQVTERSP  DAQCLRQD
101  DLGVLRFQ  QLQALLQE  RFRNVEEL  HNPVATSQ  RSRFQDEL
151  LKCLKRFQ  QVDEHVVH  EYSELEHQ  LYDETVSH  IPIGLAVL
201  ILEGLRFQ  KLEDPVAG  RYVRFCTV  CHFPVYGR  CEIIRPGE
251  TKEELIQV  SEVREYVQ  DSHDGGHT  VYRQGVSY  DQSRVQVH
301  QQVAVATW  DQVVCLEF  EYLGWDEIS  QLEHGFEL  LIEDEWGD
351  RYVAVYDPT  VQREARVQ  SVYRGRAG  HALDGAQL  HSEIRMTI
401  HNPFFSTH  DNDKLTSP  RQCRKEDG  SHVYRCSA  HPRRTWGG
451  QTVWGRW  TDDVYDWR  RQVYDWRD  DSHRFFVQ  Q
```

Unformatted sequence string: [SLL residues](#) (for pasting into other applications).

Sort by  residue number  increasing mass  decreasing mass

| Query | Start - End | Observed  | Mr (expt) | Mr (calc) | Delta H | Score | Peptide                 |
|-------|-------------|-----------|-----------|-----------|---------|-------|-------------------------|
|       | 31 - 44     | 1549.7000 | 1549.6927 | 1549.4855 | 0.0072  | 0     | <b>S-QQNNRRSIFPSAAG</b> |

**MASCOT** : Search Parameters

© 2007-2023 Matrix Science



Because glu-fib in SwissProt is not a tryptic peptide. The N-terminus is created by a post-translational cleavage after serine. If you now go back to the search form and select semi-trypsin or enzyme type none, you'll get the match.

## Modifications

PMF ✓ SQ ✓ MS/MS ✓

|                               |  |        |  |
|-------------------------------|--|--------|--|
| <b>Fixed modifications</b>    | --- none selected ---                              | ><br>< | Acetyl (K)<br>Acetyl (N-term)<br>Acetyl (Protein N-term)<br>Amidated (C-term)<br>Amidated (Protein C-term)<br>Ammonia-loss (N-term C)<br>Biotin (K)<br>Biotin (N-term)<br>Carbamidomethyl (C)<br>Carbamyl (K)<br>Carbamyl (N-term) |
|                               | Display all modifications <input type="checkbox"/> |        |  |
| <b>Variable modifications</b> | --- none selected ---                              | ><br>< |  |

- Get details of current modifications, download updates, and define new entries at <https://www.unimod.org>
- User definable with an in-house Mascot installation

MASCOT

: Search Parameters

© 2007-2023 Matrix Science

MATRIX  
SCIENCE

This screen shot shows how modifications are displayed in the search form in Mascot 2.3 and later. If you are using an earlier version, there are just two list boxes, one for fixed modifications and one for variable. In the current arrangement, you move modifications from the single list on the right to and from the lists on the left. This makes it easier to see at a glance what has been selected for the search. If the checkbox labelled 'Display all modifications' is clear, as shown here, you get a relatively short list of the most common modifications. If you check the box, a much longer list is available. You can keep your list of modifications up-to-date by downloading the latest information from Unimod. If you have a modification which you don't want to share with others, you can add it to the local configuration file. We'll describe how to go about doing this in detail in the Mascot Server Administration talk.

## Modifications

PMF✓ SQ✓ MS/MS✓

### Modifications

- Fixed / static modifications cost nothing
- Variable / differential modifications are very expensive
- Use minimum variable modifications, especially for PMF

Maybe oxidation of M

Maybe alkylation of C

Modifications in database searching are handled in two ways. First, there are the fixed or static or quantitative modifications. An example would be the efficient alkylation of cysteine. Since all cysteines are modified, this is effectively just a change in the mass of cysteine. It carries no penalty in terms of search speed or specificity.

In contrast, most post-translational modifications do not apply to all instances of a residue. For example, phosphorylation might affect just one serine in a peptide containing many serines. These variable or differential or non-quantitative modifications are expensive in the sense that they increase the time taken for a search and reduce its specificity. This is because the software has to permute out all the possible arrangements of modified and unmodified residues that fit to the peptide molecular mass. As more and more modifications are considered, the number of combinations and permutations increases geometrically. A so-called combinatorial explosion.

Hence, it is very important to be as sparing as possible with variable modifications. Especially in a peptide mass fingerprint, where the increase in the number of calculated peptides quickly makes it impossible to find a statistically significant match.

# Quantitation

PMF✗ SQ✓ MS/MS✓

Quantitation  ▾

•More later ...

Quantitation is the subject of a separate presentation.

# Crosslinking

PMF✗ SQ✗ MS/MS✓

Crosslinking None ▾

•More later ...

With Mascot Server 2.7 we introduced Crosslinking. It is also a subject of a separate talk.

## Protein mass

PMF✓ SQ✗ MS/MS✗

Protein mass  kDa

- Applied as sliding window because there is no guarantee that the database entry represents the processed protein
- Slows down the search
- Never useful for MS/MS search. Only useful for Peptide Mass Fingerprint when
  - Analyte is small fragment of very large entry
  - Low complexity entry.

MASCOT

: Search Parameters

© 2007-2023 Matrix Science

MATRIX  
SCIENCE

The protein mass is the mass of the intact protein in kDa applied as a sliding window. That is, the mass of the contiguous stretch of sequence which contains all of the matched peptide mass values. This will generally be less than the mass of the entire sequence entry. Consequently, if you specify a value for the protein mass, this acts only as a ceiling. Not only will you see smaller proteins on the hit list, you will also see larger ones, but all of the reported matches will be within a stretch of sequence less than or equal to the specified mass.

If this field is left blank, there is no restriction on protein mass.

Specifying a protein mass will slow down the search a little.

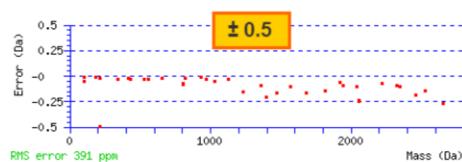
It's hard to find examples where this parameter is useful. We include it mainly because many people requested it. It could give a better score if the analyte was small fragment of very large entry, or a low complexity protein. But, you can't know this in advance, so our general recommendation is to leave the protein mass open.

## Peptide tolerance

PMF ✓ SQ ✓ MS/MS ✓

Peptide tol.  $\pm$

Specifying too tight a mass tolerance is a common reason for failing to get a match



MASCOT : Search Parameters

© 2007-2023 Matrix Science

MATRIX SCIENCE

This is the error window on experimental peptide mass values, not the error window for MS/MS fragment ion mass values, which is set using the MS/MS tol.  $\pm$  parameter.

Units can be selected from: percentage, milli-mass units, parts per million, or Daltons.

Specifying too tight a tolerance is a very common reason for failing to get a match.

Making an estimate of the mass accuracy doesn't have to be a guessing game. Protein View includes a graph of the mass errors for intact peptides. Just search a strong standard and look at the error graph. You'll normally see some kind of trend. Add on a safety margin and this is your error estimate. If you see something that looks like this, a mass tolerance of  $\pm 0.5$  Da is about right. It gives some safety margin. Remember that there will always be the odd outlier, like the data point at the lower left. It is the general trend and distribution of the majority of the data points that is important.

For a peptide mass fingerprint, the score depends on the peptide tolerance. In an MS/MS search, this parameter has no effect on the ions score. However, it does affect the search time. The larger the tolerance, the longer the search will take.

## Peptide tolerance <sup>13</sup>C

PMF ✗ SQ ✗ MS/MS ✓

# <sup>13</sup>C 0 ▾

### Sometimes, peak detection chooses the <sup>13</sup>C peak

The normal test for a precursor match is:

$TOL > \text{absolute}(\text{exp} - \text{calc})$

If this field is set to 1, the test will also succeed for

$TOL > \text{absolute}(\text{exp} - \text{calc} - 1)$

If this field is set to 2, the test will succeed for the above two conditions, plus:

$TOL > \text{absolute}(\text{exp} - \text{calc} - 2)$

MASCOT

: Search Parameters

© 2007-2023 Matrix Science

MATRIX  
SCIENCE

Sometimes, peak detection chooses the <sup>13</sup>C peak rather than the <sup>12</sup>C. In extreme cases, it may pick the <sup>13</sup>C2 peak. The normal test for a precursor match is:

$TOL > \text{absolute}(\text{exp} - \text{calc})$

Assuming the mass values and tolerance are in Da, if this field is set to 1, the test will also succeed for

$TOL > \text{absolute}(\text{exp} - \text{calc} - 1)$

If this field is set to 2, the test will succeed for the above two conditions, plus:

$TOL > \text{absolute}(\text{exp} - \text{calc} - 2)$

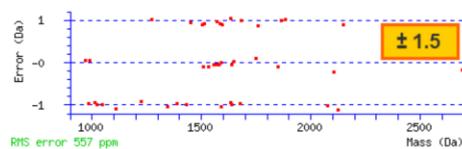
This means that you can use a tight mass tolerance and still get a match to a <sup>13</sup>C peak. If you are using a very high accuracy instrument, note that the precise shifts are the carbon isotope spacings of 1.00335 and 2.00670, rather than 1 and 2.

## MS/MS tolerance

PMF ✗ SQ ✓ MS/MS ✓

MS/MS tol. ±  Da

Specifying too tight or too loose a mass tolerance will reduce the ions score



MASCOT : Search Parameters

© 2007-2023 Matrix Science

MATRIX SCIENCE

This is the error window on MS/MS fragment mass values.

Units can be milli-mass units, Daltons, or ppm (Mascot 2.5 and later).

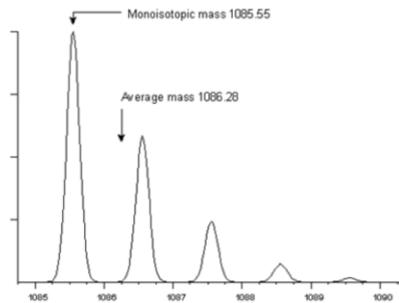
Specifying too tight or too loose a mass tolerance will reduce the ions score. Peptide View includes a graph of the mass errors for fragment ions.

Here, the mass tolerance is much too high. A more appropriate tolerance might be +/- 0.3. Having a tolerance which is much too high can sometimes lead to artefacts and false positives

## Mass type

PMF✓ SQ✓ MS/MS✓

Monoisotopic  Average



If you get this setting wrong, the mass errors will be very large and show a strong trend

MASCOT : Search Parameters

© 2007-2023 Matrix Science

MATRIX SCIENCE

Mass type specifies whether the experimental mass values are average or monoisotopic. Monoisotopic mass is the mass of the peptide where all atoms are the most abundant natural isotopes of their elements, e.g. Carbon 12, Nitrogen 14, Hydrogen 1, etc. In most cases, this is the first peak of the natural isotope distribution. Average mass is the chemical mass, which is the centre of gravity of the isotope distribution.

In Mascot, you cannot mix the two, and have (say) average precursors and monoisotopic fragments.

Most modern instruments produce monoisotopic mass values. You will only have an average mass if the entire isotope distribution has been centroided into a single peak, which usually implies very low resolution. If you get this setting wrong, the mass errors will be very large and show a strong trend, because the difference between an average and a monoisotopic mass for peptides and proteins is approximately 0.06%.

## Charge

PMF✓ SQ✓ MS/MS✓

Mass values  MH<sup>+</sup>  M<sub>r</sub>  M-H<sup>-</sup>

Peptide charge

- 1+ means MH<sup>+</sup>, 1- means M-H<sup>-</sup>, etc.
- For MS/MS, this setting is a default, which is rarely used.

MASCOT

: Search Parameters

© 2007-2023 Matrix Science

MATRIX  
SCIENCE

These fields are used to specify the peptide charge state. The radio buttons are from the peptide mass fingerprint form. The drop down list is used on the sequence query and MS/MS forms.

The notation "1+", "2+", etc. is used to save space and because some HTML form fields do not support the use of superscripts and subscripts. "1+" always means MH<sup>+</sup>, "1-" always means M-H<sup>-</sup>, etc.

For MALDI-PSD, the precursor peptides will generally be MH<sup>+</sup>, so the charge state should be set to "1+“

For an MS/MS search, the value specified here is a default. Most peak lists always specify a charge state, so default is never used.

## Data (PMF)

PMF ✓ SQ ✓ MS/MS ✗



- Mass [ intensity] [additional text]
- Applied Biosystems Data Explorer (.pkm)
- Bruker Analysis AutoExecute Data Report
- Bruker XML
- mzData (1.05)
- mzML

**MASCOT**

: Search Parameters

© 2007-2023 Matrix Science

 **MATRIX  
SCIENCE**

The contents of the query window on the peptide mass fingerprint form are only used when no data file has been specified.

The data format for a peptide mass fingerprint is auto detected. It can be a simple list of mass values, one per line. If a second value is present, it is assumed to be intensity. Any further values on the same line are ignored.

Mascot also supports other peak list formats, as listed.

mzML is the standard interchange format sponsored by the HUPO Proteomics Standards Initiative working group. The earlier standard was mzData.

## Data (MS/MS)

PMF ✗ SQ ✗ MS/MS ✓

Data file  No file selected.  
Data format Mascot generic ▾ Precursor  m/z

- Mascot Generic Format (.MGF)
- Finnigan (.ASC)
- Sequest (.DTA)
- PerSeptive (.PKS)
- Micromass (.PKL)
- Sciex API III
- Bruker (.XML)
- mzData (.XML)
- mzML (.mzML)

**MASCOT**

: Search Parameters

© 2007-2023 Matrix Science

 **MATRIX  
SCIENCE**

Data for MS/MS ion searches must be supplied as an ASCII file in one of these supported formats. The format cannot be auto-detected, and must be specified using the drop down list.

Certain data file formats, SCIEX API III, PerSeptive (.PKS), and Bruker (.XML), do not include m/z information for the precursor peptide. For these formats only, the Precursor field is used to specify the m/z value of the parent peptide.

A data file may include embedded search parameters. Most embedded parameters can only appear once, at the head of the data file. In a Mascot generic format file, a few parameters can appear within an MS/MS dataset. See the Data File Format help page for further details

If there is a conflict between the values of the embedded parameters and values entered into search form fields, the embedded parameters always take precedence. The search form fields are essentially defaults for values missing from the data file.

## Data URL

PMF ✗ SQ ✗ MS/MS ✓

Data file  No file chosen

**Data input**  Data URL (http or ftp)

**Data format** Mascot generic   m/z

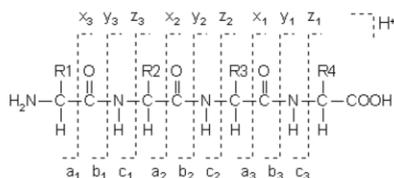
In Mascot 2.5 and later, if security is enabled, it is possible to specify a URL to the peak list. This means that the peak list file doesn't have to be downloaded to the client PC then uploaded to the Mascot server, which is useful for very large peak lists or when the client network connection is slow.

# Instrument

PMF ✗ SQ ✓ MS/MS ✓

Instrument ESI-QUAD-TOF ▾

- Click on the help link to see which ions series are used



MASCOT

: Search Parameters

© 2007-2023 Matrix Science

MATRIX  
SCIENCE

For an MS/MS Ions Search, choose the description which best matches the type of instrument used to acquire the data. This setting determines which fragment ion series will be used for scoring, according to the following table.

# Instrument

PMF  SQ  MS/MS

Mascot Configuration: Instruments

| Ion series        | Default  | ESI MALDI | ESI FTMS    | ETD MALDI | MALDI MALDI | MALDI MALDI | CID+ETD | ETD    | ETD    |
|-------------------|----------|-----------|-----------|-----------|-----------|-----------|-------------|-----------|-------------|-------------|---------|--------|--------|
|                   | QUAD TOF | TRAP TOF  | QUAD TOF  | TRAP TOF  | QUAD TOF  | TRAP TOF  | ASECTOR TOF | ECD       | TRAP        | QUAD        | QEI     | ESI    | ESI    |
| 1+                | X        | X         | X         | X         | X         | X         | X           | X         | X           | X           | X       | X      | X      |
| 2+ (precursor=2+) | X        | X         | X         | X         | X         | X         | X           | X         | X           | X           | X       | X      | X      |
| 2+ (precursor=3+) |          |           |           |           |           |           |             |           |             |             |         |        |        |
| monomers          |          | X         |           |           |           | X         | X           |           |             | X           | X       |        |        |
| a                 | X        | X         |           |           |           | X         | X           |           |             | X           | X       |        |        |
| a*                | X        | X         |           |           |           | X         | X           |           |             | X           | X       |        |        |
| ab                |          | X         |           |           |           | X         |             |           |             | X           |         |        |        |
| b                 | X        | X         | X         | X         | X         | X         | X           |           |             | X           | X       |        |        |
| b*                | X        | X         | X         | X         | X         | X         | X           |           |             | X           | X       |        |        |
| bd                | X        | X         | X         | X         | X         | X         | X           |           |             | X           | X       |        |        |
| c                 |          |           |           |           |           |           |             | X         | X           |             |         | X      | X      |
| x                 |          |           |           |           |           |           |             |           |             |             |         |        |        |
| y                 | X        | X         | X         | X         | X         | X         | X           | X         | X           | X           | X       | X      | X      |
| y*                | X        | X         | X         | X         | X         | X         | X           |           |             | X           | X       | X      | X      |
| y0                | X        | X         | X         | X         | X         | X         | X           |           |             | X           | X       | X      | X      |
| z                 |          |           |           |           |           |           |             | X         |             |             |         |        |        |
| z0                |          |           |           |           |           |           |             | X         | X           |             |         | X      | X      |
| z1                |          |           |           |           |           |           |             | X         | X           |             |         | X      | X      |
| d                 |          |           |           |           |           |           | X           |           |             |             |         |        |        |
| v                 |          |           |           |           |           |           | X           |           |             |             |         |        |        |
| w                 |          |           |           |           |           |           | X           |           |             |             |         | X      | X      |
| z=2               |          |           |           |           |           |           | X           | X         |             |             |         | X      | X      |
| Min internal mass | 700      | 700       | 700       | 700       | 700       | 700       | 700         | 700       | 700         | 700         | 700     | 700    | 700    |
|                   | Delete   | Delete    | Delete    | Delete    | Delete    | Delete    | Delete      | Delete    | Delete      | Delete      | Delete  | Delete | Delete |
|                   | Edit     | Edit      | Edit      | Edit      | Edit      | Edit      | Edit        | Edit      | Edit        | Edit        | Edit    | Edit   | Edit   |

New Instrument (Main menu)

MASCOT : Search Parameters

© 2007-2023 Matrix Science



"Default" corresponds to the configuration used in Mascot version 1.7 and earlier.

Many of the instruments are very similar.

You can modify instrument settings or create new ones using the configuration editor. In this screenshot, the right hand column is an experiment to see how the addition of w ions affects ETD matching.



# Decoy

PMF ✓ SQ ✓ MS/MS ✓

Decoy

**Sensitivity and FDR (reversed protein sequences)**

|                        | Target | Decoy    | FDR           |
|------------------------|--------|----------|---------------|
| Protein family members | 302    | 11       | 3.64%         |
| PSMs                   | above  | homology | 1821 18 0.99% |

Decoy results are available in [the decoy report](#).

For large scale experiments, provide the results of any additional statistical analyses that indicate or establish a measure of identification certainty, or allow a determination of the false-positive rate, e.g., the results of randomized database searches or other computational approaches.

MASCOT : Search Parameters

© 2007-2023 Matrix Science



The decoy checkbox enables you to validate the false discovery rate according to the approach recommended in the Molecular & Cellular Proteomics Guidelines for Publication: “For large scale experiments, provide the results of any additional statistical analyses that indicate or establish a measure of identification certainty, or allow a determination of the false-positive rate, e.g., the results of randomized database searches or other computational approaches”.

# Setting defaults

Access Mascot | Protein identifi: x

edectus/mascot/search\_form\_select\_windows.html

- Mascot overview
- Search parameter reference
- Data file format
- Results report overview

**Peptide Mass Fingerprint**  
The experimental data are a list of peptide mass values from the digestion of a protein by a specific enzyme such as trypsin.  
Perform search | Example of results report | More information

**Sequence Query**  
One or more peptide mass values associated with information such as partial or ambiguous sequence strings, amino acid composition information, MS/MS fragment ion masses, etc. A super-set of a sequence tag query.  
Perform search | Example of results report | More information

**MS/MS Ions Search**  
Identification based on raw MS/MS data from one or more peptides.  
Perform search | Example of results report | More information

**Xcalibur RAW file**  
This form can be used to convert a Thermo Xcalibur RAW file into a peak list file, which is then loaded into the MS/MS search form. Requires a local copy of the ExtractMm utility.

**Search form defaults**  
Save your preferred default settings as a browser cookie.

PMF ✓ SQ ✓ MS/MS ✓

Matrix Science - Mascot - Set : x

edectus/mascot/cgi/form\_defaults.pl

### Set Mascot search form defaults

Database: PRIDE\_human, PRIDE\_contaminants, Mus\_musculus\_GRCv38\_genomic, Mus\_EST, GRAP

Taxonomy: All entries

Enzyme: Trypsin

Allow up to: 1 missed cleavages

Fixed modifications: Acetyl (N-term), Acetyl (Protein N-term), Amidated (C-term), Amidated (Protein C-term)

Variable modifications: Acetyl (N), Acetyl (Protein N-term), Amidated (C-term), Amidated (Protein C-term)

Show all mods:

Quantification: None

Peptide tol. s: 1.2 Da, 1 kDa

MS/MS tol. s: 0.6 Da

Peptide charge: 1+

Monoisotopic:  Average

Data format: Mascot generic (MS/MS only)

Instrument: Default (MS/MS only)

Decoy:

Error tolerant:

Report top: AUTO hits

Save defaults as cookie

**MASCOT** : Search Parameters

© 2007-2023 Matrix Science



You can choose your own defaults for the search forms. Look for the link at the bottom of the search form selection page.

When you save the defaults, they are saved as a browser cookie. If you go to a different PC, or switch to a different browser, you'll need to repeat this step.

## Final Tip

### DANGER!

- Iteratively adjusting search parameters to get a better score can give misleading results
- Beware of
  - Narrowing the taxonomy
  - Reducing mass tolerances
  - Removing modifications
  - Selecting spectra or mass values

**Set search parameters using standard samples**

A final word of advice: It is easy to distort the search results without realising.

Basically, it is risky to adjust the search parameters interactively to get a better score for an unknown.

For example, you search the complete database and don't get a significant match. However, a very interesting looking protein is near the top of the list, surrounded by some others that are clearly wrong. You change the taxonomy filter so as to exclude the "wrong" proteins. Sorry, but this is cheating.

Search parameters should be set using standards. Broadening the search if you get a negative result is usually OK, but not narrowing the search.