# *Mascot Search Results FAQ*

*{MATRIX}{SCIENCE}*

# Why are some peptide matches shown in red or bold face?

2

## Why are some peptide matches shown in red or bold face?

- Red indicates the top scoring peptide match for this spectrum
  - Not necessarily a significant match
- Bold indicates the first time any match to this spectrum has appeared in the report
  - So, if a protein hit *doesn't* have any bold red matches, all the assigned spectra have better scoring matches elsewhere or the same matches have already appeared in the report, assigned to higher scoring protein(s)

| Format As | Peptide Summary | | | Help | |
|---|---|---|---|---|---|
| | Significance threshold p< 0.05 | Max. number of hits 50 | | | |
| | Standard scoring ⊙ MudPIT scoring ○ | Ions score cut-off 0 | | Show sub-sets ☐ | |
| | Show pop-ups ⊙ Suppress pop-ups ○ | Sort unassigned Decreasing Score | | Require bold red ☐ | |

**ASMS 2005**

*MATRIX SCIENCE*

Interpretation of the results from an LC-MS/MS search can be complex, because it is not always clear which peptide "belongs" to which protein. The use of red and bold typefaces is intended to highlight the most logical assignment of peptides to proteins. The first time a match to a spectrum appears in the report, it is shown in bold face. Whenever the top scoring peptide match for a spectrum appears, it is shown in red. This means that peptide matches which are both bold and red are the most likely assignments of the best matches. Conversely, if a protein hit *doesn't* have any bold red matches, all the assigned spectra have better scoring matches elsewhere or the same matches have already appeared in the report, assigned to higher scoring protein(s). This means that the protein hit is likely to be spurious, and would collapse into a higher scoring hit except for the presence of one or more weak, noisy matches. Such hits can be filtered out of the report by ticking the 'require bold red' checkbox.
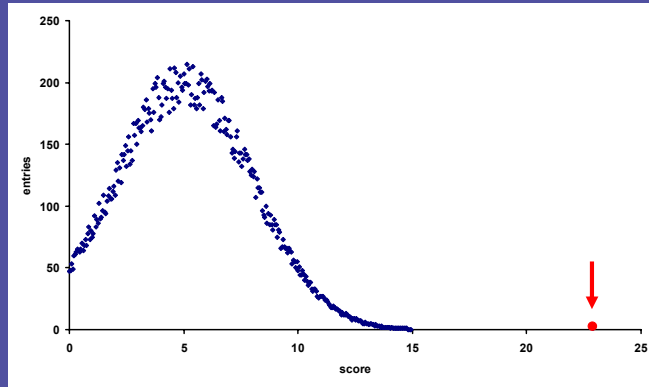
3

If a score is a true probability, then assigning a significance threshold is very simple ... its just a function of the number of trials - the number of times we test for a match.

If we are comfortable with a 1 in a 20 chance of getting a false positive match, and we are doing a MS/MS search of a database that contains 5000 peptides that fit to the precursor molecular weight , then we are looking for a probability of less than $1 / (20 \times 5000)$ which is a Mascot score of 50

If we could only tolerate a false positive rate of 1 in 200 then the threshold would be 60, 1 in 2000 70, etc.

*What is the difference between the identity threshold and the homology threshold?*

- The homology threshold is an empirical measure of whether the match is an outlier

ASMS 2005

Unfortunately, mass spectra are often far from ideal, with poor signal to noise or gaps in the fragmentation. In such cases, it may not be possible to reach this statistical threshold score, even though the best match in the database is a clear outlier from the distribution of random scores. To assist in identifying these outliers, we also report a second, lower threshold, the 'homology' threshold. This simply says the match is an outlier.

In practice, from measuring the actual false positive rate by searching large data sets against reversed or randomised databases, we find that the identity threshold is usually conservative, and the homology threshold can provide a useful number additional true positive matches without exceeding the specified false positive rate.

# *What is an expectation value?*

- The number of times you could expect to get this score or better by chance

$$E = P_{threshold} * (10 ** ((S_{threshold} - score) / 10))$$

If $P_{threshold}$ = 0.05 and $S_{threshold}$ = 50
  - score = 40 corresponds to E = 0.5
  - score = 50 corresponds to E = 0.05
  - score = 60 corresponds to E = 0.005

**ASMS 2005**

*MATRIX SCIENCE*

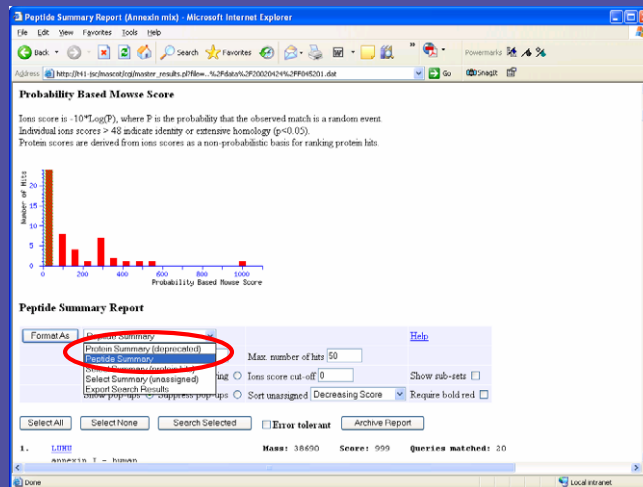The expectation value does not contain new information. It can be derived directly from the score and the threshold. The advantage is that it tells you everything you need to know in a single number.

It is the number of times you could expect to get this score or better by chance.

A completely random match has an expectation value of 1 or more

The better the match, the smaller the expectation value.

# Why does it say that the Protein Summary is deprecated?

*MATRIX SCIENCE*

*Why does it say that the Protein Summary is deprecated?*

- **The Protein Summary is intended for peptide mass fingerprint results**
  - Limited to 50 hits
  - Protein score and expectation value for a search containing MS/MS data may be misleading because the matches of the precursor masses are being scored as a PMF
  - Not available for large searches (> 1000 queries)

*ASMS 2005*

*MATRIX SCIENCE*

Peptide Summary is the default for any search containing MS/MS data. Unfortunately, some older, third party software specifies a Protein Summary report when submitting an MS/MS search. This means that you have to switch formats yourself. A lot of people didn't realise this, so we have made the warnings more prominent.

# *What is MudPIT scoring?*

- Standard protein score
  - the sum of the ions scores
  - excluding the scores for duplicate matches, which are shown in parentheses
  - correction to reduce the contribution of low-scoring random matches

```
183.  IPI00141647                        Mass: 3011421  Score: 47    Queries matched: 5
      Tax_Id=9606 titin isoform N2-B
   ☐ Check to include this hit in error tolerant search or archive report


   Query  Observed   Mr(expt)    Mr(calc)   Delta Miss Score Expect Rank Peptide
     76  424.6860   847.3575    845.4355    1.9220  1   17     27    8  NDGGSRIK
    118  446.7500   891.4854    889.4691    2.0164  0   17     25    4  GGIQIMAGK
    358  366.3330  1095.9773   1092.6179    3.3594  0   23    5.7    8  YISSLEILR
    569  439.3649  1315.0730   1313.6615    1.4115  0   26    2.9    1  EPVLYDTHVNK
   1182  870.8864  1739.7583   1741.8886   -2.1303  0   15     27    3  VTAVNEYGPGVPTDVPK
```
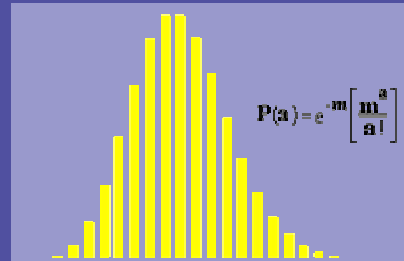
**ASMS 2005**

*MATRIX SCIENCE*

With standard peptide summary scoring, the protein score is the sum of the ions scores of all the non-duplicate peptides. Where there are duplicate peptides, the highest scoring peptide is used.

This example shows how we can get a protein score of 47 even though none of the peptide matches are significant

10

*What is MudPIT scoring?*

- – Even if you only have random matches, you can still get multiple matches to a protein.
- – The distribution of random matches depends on the ratio between the number of spectra and the number of entries in the database
- – Poisson distribution

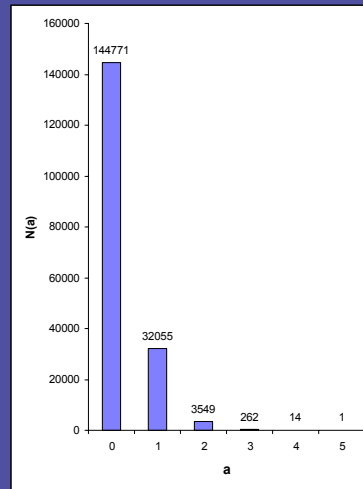$$P(a) = e^{-m} \left[ \frac{m^a}{a!} \right]$$

ASMS 2005

MATRIX SCIENCE

Even if peptide matches are random, you can still get multiple matches to a single protein. How likely this is depends on the ratio between the number of spectra and the number of entries in the database. We can predict whether this will be a serious problem or not using a function called a Poisson distribution.

If average number of events per interval is $m$, then the Poisson distribution gives us the probability of observing $a$ events in a particular interval.

## What is MudPIT scoring?

- **Shotgun / MudPIT**
  - 20 SCX fractions
  - 160,000 scans total
  - 80,000 after processing
  - 40,000 random matches in search of Swiss-Prot (180652 entries)

MATRIX SCIENCE

For this MudPIT search, 262 proteins are expected to pick up 3 random matches by chance. 1 protein will pick up 5

12

*What is MudPIT scoring?*

- Small database
  - 30 minute run
  - 1500 scans total
  - 1200 after processing
  - 550 random matches in search of Swiss-Prot using drosophila taxonomy filter (2727 entries)

*ASMS 2005*

The problem isn't limited to large searches. It is the ratio between the number of spectra and the number of entries in the database that matters. So, a small search against a small database can give similar numbers

**What is MudPIT scoring?**

- MudPIT protein score
  - The sum of the excess of the ions score over the identity or homology threshold for each query
  - Plus 1 x the average threshold

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| 178. | IPI00001639 | | | Mass: 98420 | | Score: 46 | | Queries matched: 3 |

Tax_Id=9606 Importin beta-1 subunit

☐ Check to include this hit in error tolerant search or archive report

| Query | Observed | Mr(expt) | Mr(calc) | Delta | Miss | Score | Expect | Rank | Peptide |
|---|---|---|---|---|---|---|---|---|---|
| 22 | 386.4956 | 770.9767 | 770.4286 | 0.5481 | 0 | 22 | 8.5 | 3 | DPSVVVR |
| 914 | 779.7214 | 1557.4282 | 1555.8205 | 1.6077 | 1 | 23 | 4.9 | 2 | TVSPDRLELEAAQK |
| 1359 | 918.3068 | 1834.5991 | 1832.8839 | 1.7152 | 0 | 46 | 0.024 | 1 | GDQENVHPDVMLVQPR |

**ASMS 2005**

**MATRIX SCIENCE**

For MudPIT scoring, the score for each peptide is not its absolute score, but the amount that it is above the threshold. Therefore, peptides with a score below the threshold do not contribute to the score. Finally, the average of the thresholds used is added to the score. For each peptide, the "threshold" is the homology threshold if it exists, otherwise it is the identity threshold.

You shouldn't see proteins with a large number of weak peptide matches getting a good score. If there are no significant peptides, the protein score will be 0.

By default, MudPIT protein scoring is used for searches with more than 1000 spectra. You can also choose which scoring to use in the format controls at the top of the report

14

Probability based scoring tells you the probability that the match is random. This is, the probability that the match is meaningless. Many people would prefer a probability that the match is correct. Is this possible?

It is certainly possible if you are analysing a known protein or standard mixture of proteins. If you know what the sequences are, or think you know, then the matches to the known sequences are defined to be correct and those to any other sequence are defined to be wrong.

If the sample is an unknown, then it is difficult to define what is meant by a correct match.

Is the correct match the best match in the database? Certainly not ... this would be a false positive if the correct sequence was not in the database.

What about the best match out of all possible peptides. Yes, a reasonable definition, but not a very practical one. This is what we try to find in de novo sequencing. The reason for searching a database is that the data quality are not good enough for reliable de novo, so we reduce the size of the search space to the content of the chosen sequence database.

How about the peptide sequence that is uniquely and completely defined by the MS data? This is equally impractical. One rarely, if ever, sees a mass spectrum perfect enough to meet this criterion
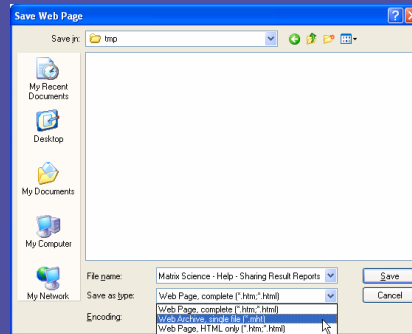
This is a typical MS/MS search result, where we see a series of high scoring homologous peptides. The sequences of the top four matches are very similar, and their expectation values vary from random through to very unlikely to be random. The best match has an expectation value of 2E-5. However, we cannot be sure that this is an identity match to the analyte peptide. It is simply the best match we could find in the database. There is always the possibility that a better match exists, that is not in the database, so to call it the correct match can be misleading.

The important thing is that we can recognise and discard matches that are nothing more than random matches. I guess we aren't even sure how to define correct, never mind calculate a probability for a particular match being correct

16

*How can I send a result report to a colleague?*

- Save a single report as web page complete or web archive

ASMS 2005

Saving as "Web page, HTML only" is no good because graphics like the score histogram will be missing.
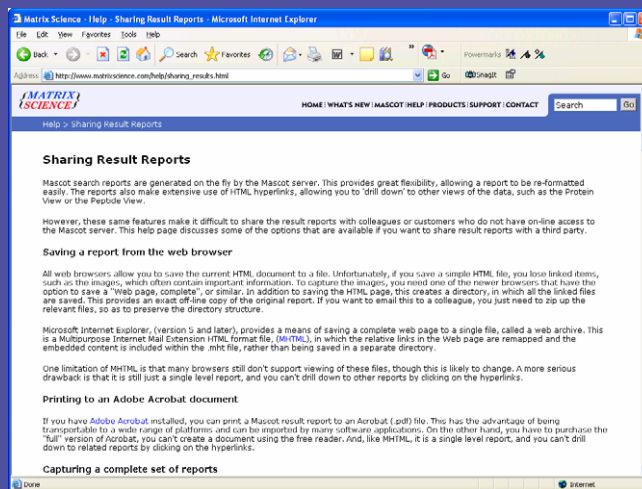
*How can I send a result report to a colleague?*

- Print a single report to an Acrobat PDF file
- Capture a complete set of reports using an off-line browser utility
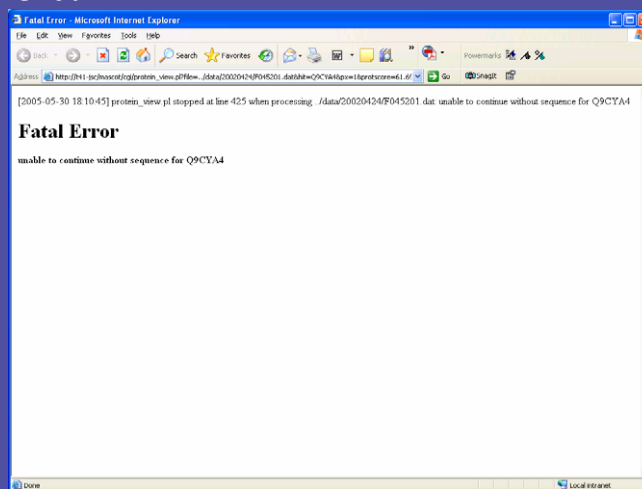
18

*How can I send a result report to a colleague?*

ASMS 2005

A detailed answer to this question can be found at
http://www.matrixscience.com/help/sharing_results.html

*Why can't I get a Protein View report for some hits from an old search?*

ASMS 2005

Mascot cannot save complete protein sequences to the result file. To do so would make each result file enormous. When you request a Protein View report, the script gets the protein sequence from the current Fasta file.

Unfortunately, database accessions are not nearly as stable as one might expect. A percentage disappear in each update because the entry is revised and gets a new accession. Even Swiss-Prot suffers from this problem

One fix would be to retain old database files on your Mascot server by creating a new database definition for each update, rather than just replacing the file. Alternatively, if the result is important, it may be easiest to repeat the search against the current Fasta file. I prefer the second route. You just have to choose "Search selected" from the Master results report

# *Why can't I get a Protein View report for some hits from an old search?*

*MATRIX SCIENCE*

The Mascot result file includes title strings and mass values for all the proteins it "expects" to display in a standard report. However, in a large search, it may miss a few. If this information is missing for a primary hit, (the first one listed for a given hit number), the report uses a utility (ms-getseq.exe) to retrieve this information from the Fasta file. For non-primary hits, it does not do this because it would greatly increase the time taken to load the report.

*Why do some protein hits in a Peptide Summary not have a mass and description?*

ProteinsInResultsFile 2

Determines the number of protein title lines saved to each results file.

1. As in Mascot 1.7 and earlier, only proteins that appear in the Summary section will appear in the Proteins section
2. Include proteins with at least one top ranking peptide match to a peptide of length greater than MinPepLengthInPepSummary
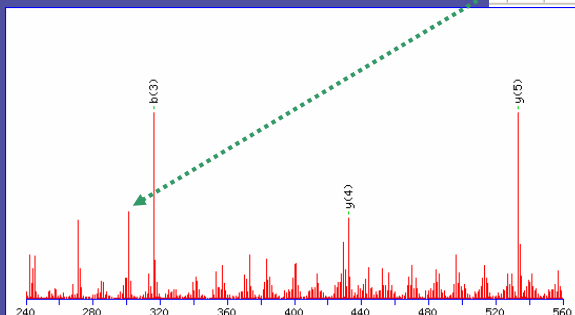3. Include all proteins

MATRIX SCIENCE

This behaviour is controlled by a setting in the Options section of mascot.dat. You can force the title and mass for all proteins to be stored in the result file, so that they are never missing. The down side is that this will cause the size of the result files to increase.

The setting is ProteinsInResultsFile in the options section of mascot.dat (see chapter 6 of the Mascot installation and setup manual).

Mascot begins by selecting a small number of experimental peaks on the basis of normalised intensity. It calculates a probability based score according to the number of matches. It then increases the number of selected peaks and recalculates the score. It continues to iterate until it is clear that the score can only get worse. It then reports the best score it found, which should correspond to an optimum selection, taking mostly real peaks and leaving behind mostly noise

Mascot is not trying to find all possible matches in the spectrum. As in this example, many spectra have "peak at every mass" noise, and can match any fragment ion from any sequence if there is no intensity discrimination. So, you may look at a peptide view report and see obvious matches that are unlabelled. However, if the peak selection was to be extended to include these additional matches, it would also have to include a number of additional noise peaks, and the score would decrease.

There are no user parameters to influence this behaviour.