

What's new in Mascot 2.3

MASCOT

MATRIX
SCIENCE

What's new in Mascot 2.3

- New report formats
- Searching multiple fasta files
- Support for HUPO PSI standard formats
- Support for Percolator
- Improvements to Mascot Daemon
- Mascot Distiller 2.3

MASCOT : *Version 2.3*

© 2009 Matrix Science



In this session, I'll be describing some of the changes in Mascot 2.3

I'll be describing some significant improvements to Mascot Server, Mascot Daemon and also to Mascot Distiller.

New (MS-MS) report formats

1. Improved protein grouping
2. Significantly faster to load huge reports
3. More flexible user interface
4. Minimal amount of memory required for browser on client

MASCOT : *Version 2.3*

© 2009 Matrix Science



We have put a significant amount of effort into changing the way we report results from Mascot searches.

The main improvements that I want to highlight are:

- Improved protein grouping
- Significantly faster to load huge reports
- More flexible user interface
- Minimal amount of memory required for browser on client

Protein grouping in Mascot 2.2

20: Q3UWB9 ... UDP-glucuronosyltransferase 2 family

34: Q91WH2 ... SIMILAR TO UDP- GLUCURONOSYLTRANSFERASE 2 FAMILY

51: Q3UEP4 ... similar to UDP-glucuronosyltransferase 2 family

| 20: Q3UWB9 | 34: Q91WH2 | 51: Q3UEP4 |
|--------------------|--------------------|---------------|
| GAVALNIR | | |
| | GAAVTLNIR | GAAVTLNIR |
| | AEMWLIR | AEMWLIR |
| IILDELVQR | IILDELVQR | |
| FSPGYQIEK | FSPGYQIEK | |
| DNLENFFIK | DNLENFFIK | |
| FETFPTSISK | FETFPTSISK | |
| FVDVWTYEMPR | FVDVWTYEMPR | |
| | | WTYEVPR |
| | | IILDELK |
| | | TPATLGPNTR |
| | | FSPGYLEK |
| WLPQNDLLGHPK | WLPQNDLLGHPK | WLPQNDLLGHPK |
| ANAIWALAQIPQK | | ANAIWALAQIPQK |
| GHEVTVLRPSAYYVLDPK | GHEVTVLRPSAYYVLDPK | |

MASCOT : Version 2.3

© 2009 Matrix Science



I'd like to show an example of the protein grouping in Mascot 2.2 producing less than ideal results.

If you look at the description lines for these three proteins, I think that we would all agree that they should be grouped together. However, they are reported as three separate proteins in Mascot 2.2

The first column shows the peptide matches for the hit 20.

Hit number 34 shares 7 peptide matches with hit number 20, but also matches these 2 additional peptides, which is why it is recorded as a separate match.

Hit number 51 has 4 new peptide matches, and the rest are either in hit 34, hit number 20 or both.

If we were to swap the order of hits 34 and 51, we can see that hit 34 just contains peptides in the other two – i.e. there are no new peptides.

Mascot 2.3 - protein grouping

1. Take the (next) top scoring protein
2. List all peptides above homology threshold
3. Find all other proteins with a match to one or more of the same peptides
4. Repeat from 2 until no new proteins found
5. Group using hierarchical pairwise single-linkage clustering.
6. Repeat from 1

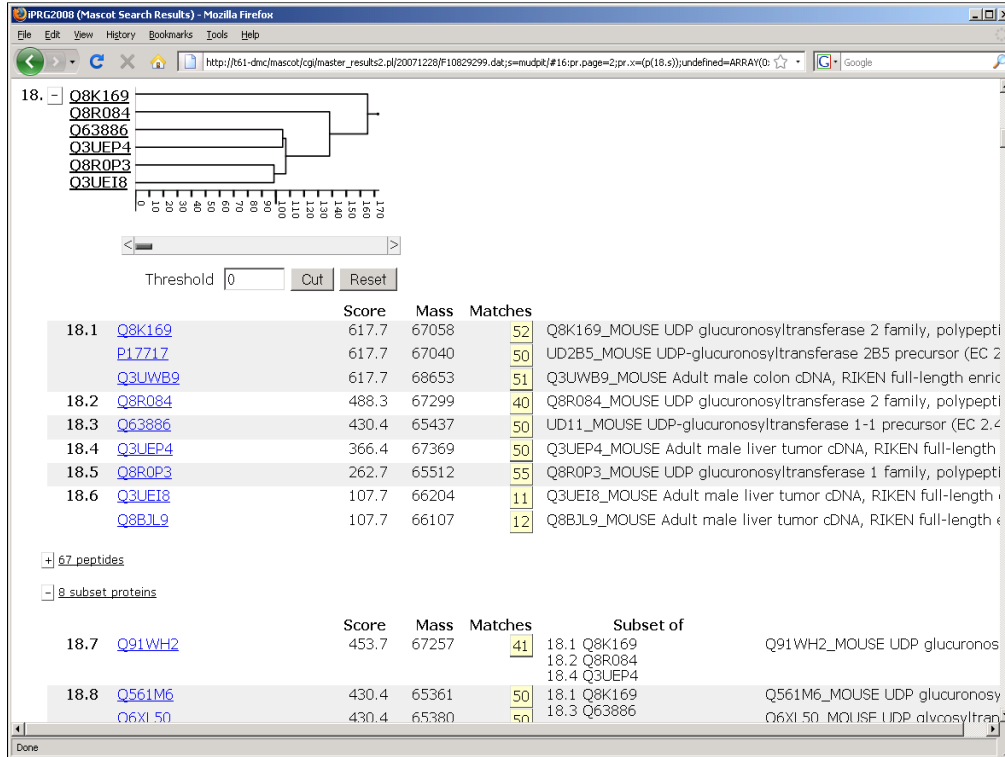
MASCOT : *Version 2.3*

© 2009 Matrix Science

MATRIX
SCIENCE

In Mascot 2.3, these proteins are grouped together.

Here's an overview of how the new protein grouping works



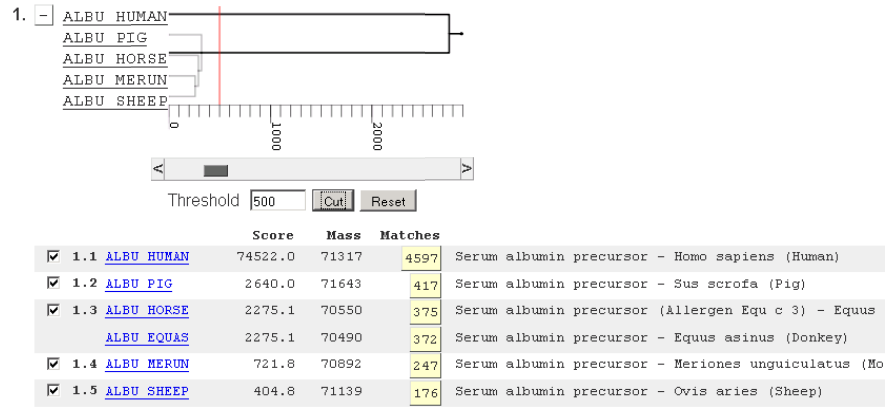
If we look at this search using the Mascot 2.3 reports, we can see that it has also grouped together a number of other glucuronosyltransferase proteins.

We can see that although they are all part of the same ‘family’, we have used hierarchical clustering to differentiate between protein groups. So, Q3UEP4 is in a separate group within the family.

As expected we can also see that Q91WH2 is a subset of Q3UEP4 and Q3UWB9

I’ll explain a little bit about the dendrogram now, but it’s easier to do this with a different example:

New reports - dendrogram



MASCOT : Version 2.3

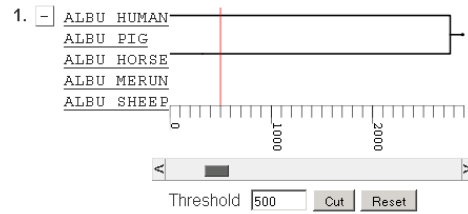
© 2009 Matrix Science



This is an example of a match to a human serum albumin protein, where we have searched all species.

The dendrogram is a pictorial representation of the hierarchical clustering. We calculate a distance tree based on the scores of significant peptides found in one protein but not the other. It's immediately obvious that the pig, horse and sheep albumins are similar. It's also obvious that there is big difference between them and the human protein. We could decide that we aren't interested in distinguishing between the different animal proteins, and cut the dendrogram at, say 500

New reports - dendrogram



| | Score | Mass | Matches | | |
|---|----------------------------|---------|---------|------|--|
| <input checked="" type="checkbox"/> 1.1 | ALBU HUMAN | 74522.0 | 71317 | 4597 | Serum albumin precursor - Homo sapiens (Human) |
| <input checked="" type="checkbox"/> 1.2 | ALBU PIG | 2640.0 | 71643 | 417 | Serum albumin precursor - Sus scrofa (Pig) |
| | ALBU HORSE | 2275.1 | 70550 | 375 | Serum albumin precursor (Allergen Equ c 3) - Equus |
| | ALBU MERUN | 721.8 | 70892 | 247 | Serum albumin precursor - Meriones unguiculatus (M |
| | ALBU SHEEP | 404.8 | 71139 | 176 | Serum albumin precursor - Ovis aries (Sheep) |
| | ALBU EQUAS | 2275.1 | 70490 | 372 | Serum albumin precursor - Equus asinus (Donkey) |

MASCOT : Version 2.3

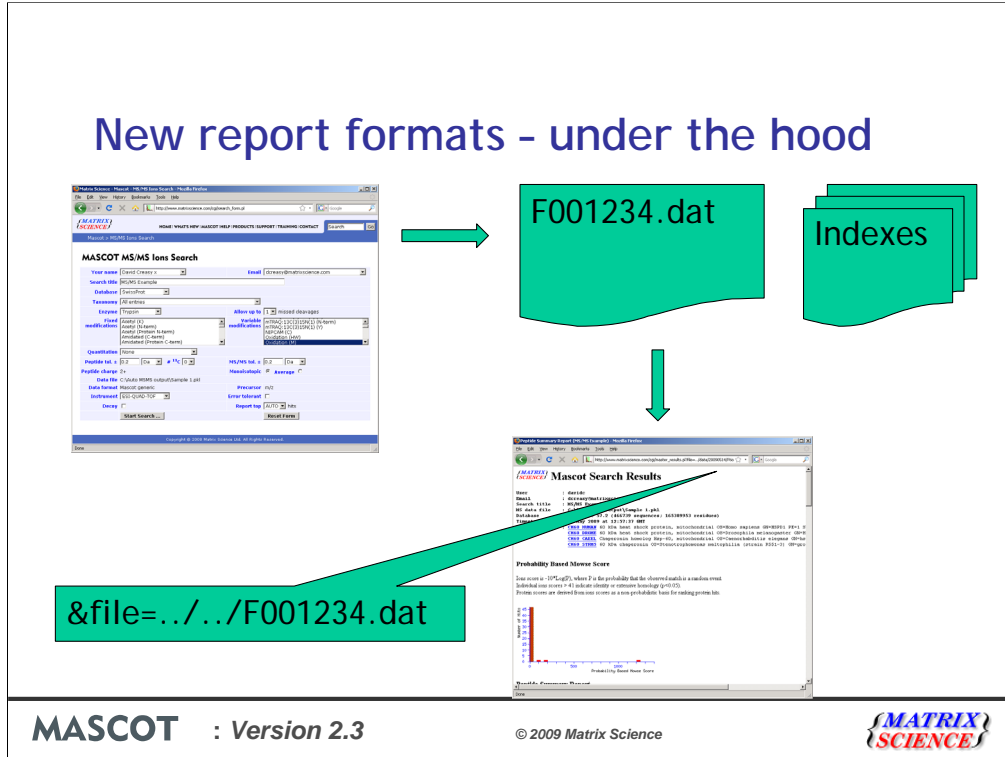
© 2009 Matrix Science

MATRIX
SCIENCE

And we now just have two family members. You can see here that the horse, donkey and sheep albumin proteins are all now considered as 'samesets' of the pig protein.

We aren't sure how useful this feature is, so we'd appreciate any feedback.

New report formats - under the hood



I will now describe some of the other changes to the reports, but first a bit of background information... when a search is submitted, it runs a cgi program called `nph-mascot.exe`. This program saves a file with all the search results onto the Mascot server. It automatically then loads another program which displays the results in a friendly manner.

It's the changes to this second piece of software that I am describing. Because of this architecture, it means that old searches can be viewed using the new protein grouping and with the other enhancements that I'm about to describe. To help speed things up, we've also added some indexes.

New reports - load much more rapidly

Search of 114,943 ms-ms spectra against a database with 1,319,480 sequences

Mascot 2.2:

- 35 minutes to load report
- No progress reports
- 430Mb memory for 1487 proteins

Mascot 2.3 - dramatic improvement...

MASCOT : Version 2.3

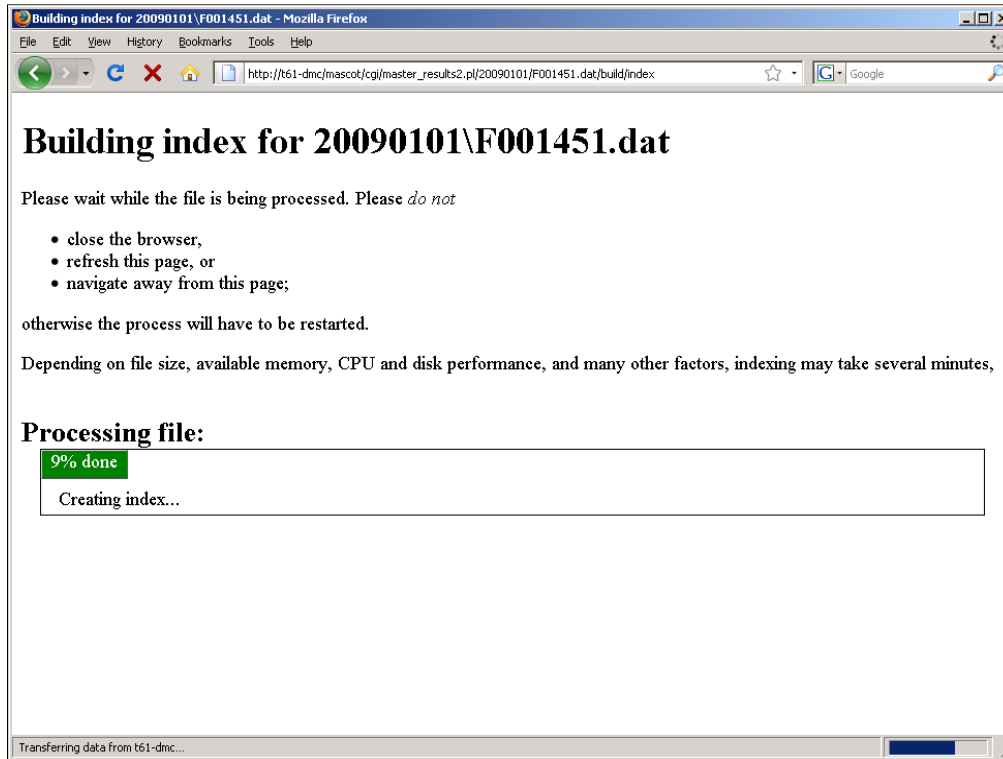
© 2009 Matrix Science



The first thing that you will notice is that the reports load much more rapidly. As an example, I'm going to show a search of 114 thousand spectra against a sequence database with 1.3 million entries. The resulting Fxxxx.dat file is just over a GB in size.

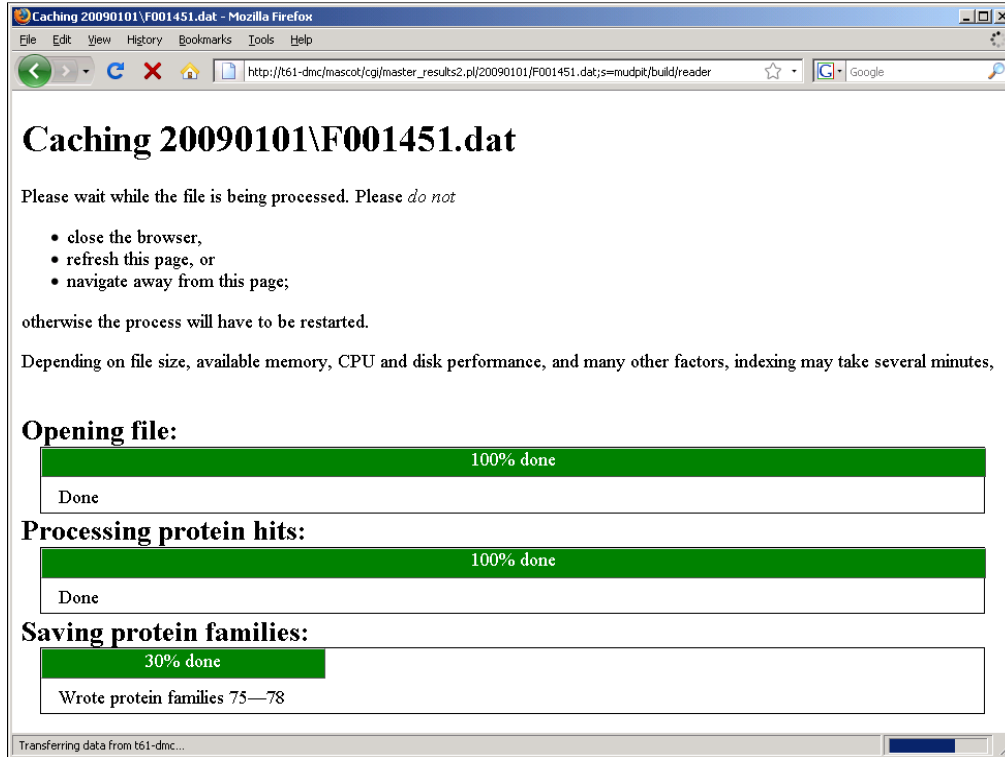
Mascot 2.2 takes 35 minutes to load the report which contains 1487 proteins. One of the problems is that there is no indication as to how much longer you need to wait for the report to be displayed.

In Mascot 2.3 ...



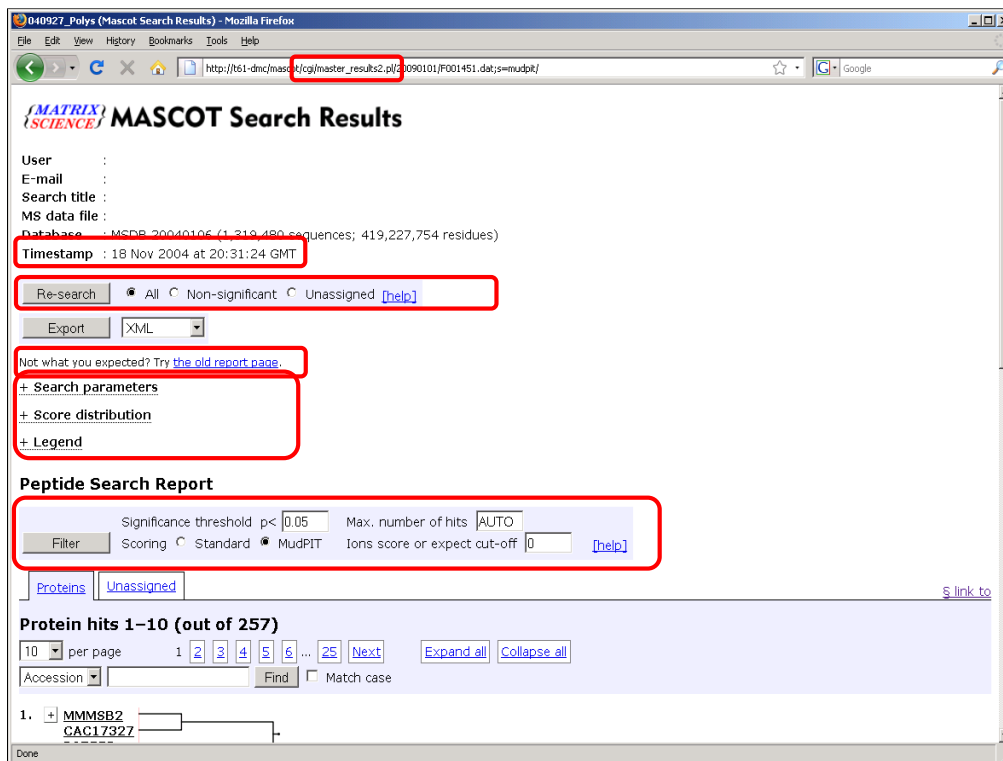
you instantly get some feedback. The first time that you load a report, various indexes get created. You may not see this page if you are running searches from Mascot Daemon, because the first index gets created as part of running the search. This part for our 1Gb file takes less than 3 minutes. I'm running this on a 2GHz laptop with a fairly slow disk, so you'd expect any recent server to be faster.

Apart from our own reports, any 3rd party programs that use Mascot Parser can also take advantage of these indexes without any change to their code.



Once the main index has been created, the file and creates an initial list of candidate proteins from the peptide hits and then groups the proteins.

It then saves the results to cache files. These 3 steps takes a further 3 and a half minutes on my laptop.



And about 5 seconds later the report loads. If you decide to re-load the report later this will take just 15 seconds because it doesn't need to re-create the cache files. So, I hope you'll agree that this is a rather dramatic improvement – 35 minutes to 15 seconds. Furthermore, it only uses 50Mb memory in the browser compared with 430Mb

As you can see, this is an old search from 2004. As I mentioned earlier, the new reports can be used with old searches – there should be no need to repeat the search unless, for example, you want to search against a new database.

The next thing to point out is that it is now much easier to repeat a search, either for all ms-ms spectra, or just ones below the significance threshold.

Some people will find that the new reports take a bit of getting used to. There will always be a link to the old reports, or for anyone that really doesn't like the new report format, there's a configuration option available to load the old report by default. Also, those with 20 vision will be able to see that the new report is called master_results2.pl – this means that any third party software that automatically loads the reports and then parses the html may still work because the original master_results.pl is almost unchanged. I do of course need to point out that this is a bad thing to do, and you really should be using the Export function for this type of task.

The filtering options are much the same as in previous versions of Mascot. If you change any of these values, new cache files need to be created, so there will be a delay. You'll notice that there is no longer a 'require bold red' option. This was useful in earlier versions of Mascot where the grouping was less sophisticated. There is also no show subsets option because we can now show or hide these for a particular protein 'on demand'.

In fact, many parts of the new reports can be expanded/contracted on demand. We are using 'AJAX' so this is generally very fast. If I click here, for example,

040927_Polys (Mascot Search Results) - Mozilla Firefox

File Edit View History Bookmarks Tools Help

http://t61-dnc/mascot/cgi/master_results2.pl(20090101/F001451.dat;s=mudpit/

MASCOT Search Results

User :
 E-mail :
 Search title :
 MS data file :
 Database : MSDB 20040106 (1,319,480 sequences; 419,227,754 residues)
 Timestamp : 18 Nov 2004 at 20:31:24 GMT

Re-search All Non-significant Unassigned [\[help\]](#)

Export XML

Not what you expected? Try [the old report page](#).

- Search parameters

Type of search : MS/MS Ion Search
 Enzyme : Trypsin
 Fixed modifications : Carbamidomethyl (C)
 Mass values : Average
 Protein mass : Unrestricted
 Peptide mass tolerance : ± 2 Da
 Fragment mass tolerance : ± 0.8 Da
 Max missed cleavages : 1
 Instrument type : ESI-TRAP
 Number of queries : 114,943
 Error tolerant search : 0

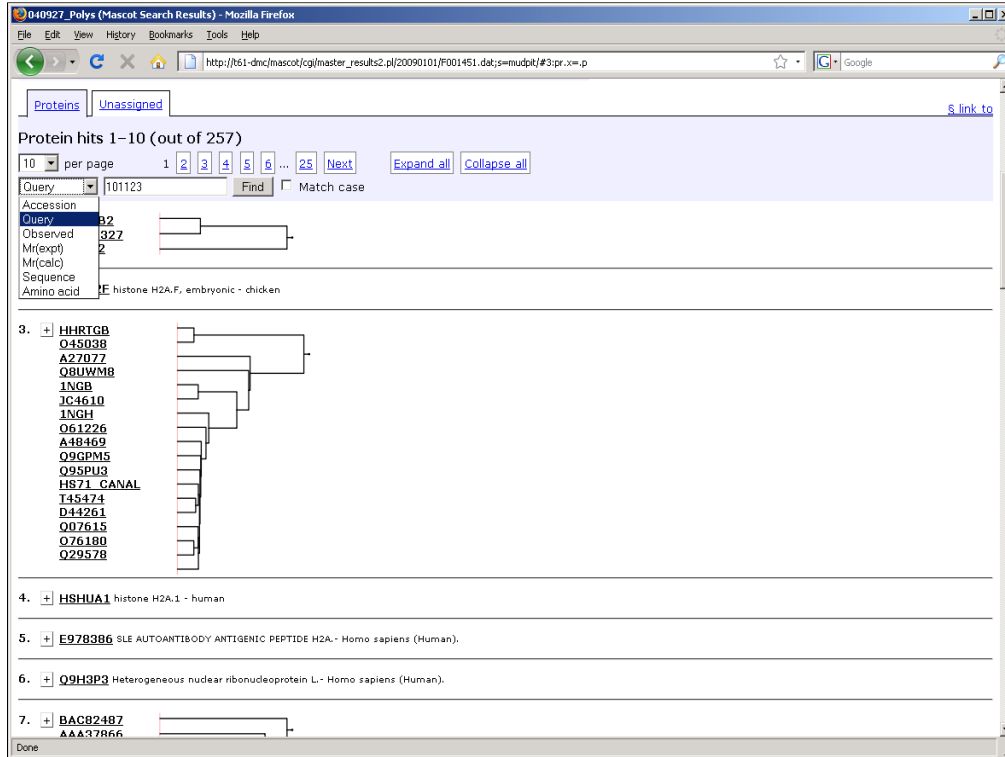
+ Score distribution

- Legend

| Expect | Rank | 1 | 2 | Peptide | |
|----------|------|---|---|-------------|-----------------------------|
| 0.037 | 2 | | | GAYLSLR | significant |
| 9 | 1 | | | GFFLPVEGGR | top ranking |
| 6.4e-005 | 1 | | | GSSIFGLAPGK | significant and top ranking |
| 0.010 | 1 | | | MAPLFFGLNPK | duplicate peptide |

Done

...we can see the search parameters.



The next thing that you will notice is that by default, we only show 10 hits per page, although this can easily be changed from the drop down list here. The significant reduction in time to load the report is due to this being a paged report where not all hits are shown at once.

There's also an option to search for a particular query, accession or even precursor match – very important with a paged report.

Anybody who is really wide awake will notice that I said that there were 1487 protein matches in Mascot 2.2, but here you can see that there are only 257 protein families. This is because the protein grouping is so much more efficient in Mascot 2.3

Mascot MS/MS Search Results (20080424/F002926.dat) - Mozilla Firefox

ALBU EQUAS 2275.1 70490 372 Serum albumin precursor - Equus asinus (Donkey)

201 peptides

| Query | Dups | Observed | Mr (expt) | Mr (calc) | Delta | Miss | Score | Expect | Rank | 1 2 | Peptide |
|--------|------|-----------|-----------|-----------|---------|------|-----------|---------|------|-----|---|
| 179721 | 23 | 1043.2390 | 2084.4634 | 2085.8303 | -1.3668 | 97 | 4.9e-07 | 1 | 1 | 1 | K.VHTECHGDLLECADDR.A |
| 179823 | 26 | 696.1702 | 2085.4888 | 2085.8303 | -0.3415 | 65 | 0.00081 | 1 | 1 | 1 | K.VHTECHGDLLECADDR.A |
| 180937 | | 699.7211 | 2096.1415 | 2094.9746 | 1.1669 | 6 | 6.7e+02 | 6 | 6 | 6 | R.DELPADLNPLEHDFVEDK.E |
| 182923 | | 705.3275 | 2112.9607 | 2112.8775 | 0.0831 | 1 | 2.1e+03 | 9 | 9 | 9 | K.VHKECHGDLLECADDR.A |
| 183017 | | 1058.0150 | 2114.0154 | 2112.8775 | 1.1379 | 1 | 0.23e+03 | 1 | 1 | 1 | K.VHKECHGDLLECADDR.A |
| 183784 | | 708.1685 | 2121.4837 | 2123.9802 | -2.4965 | 1 | 19 | 29 | 1 | 1 | K.AAFTECCQADKACLPEK.L |
| 184508 | | 1064.2730 | 2126.5314 | 2130.1877 | -3.6563 | 1 | 3.1e+03 | 2 | 2 | 2 | M.RWVTFISLLLLFSSAYS.R.G |
| 186016 | | 1069.5870 | 2137.1594 | 2132.8085 | 4.3509 | 1 | 2.15e+03 | 6 | 6 | 6 | R.ETTYGMADCCERQEPER.N + Oxidation (M) |
| 188763 | | 721.5112 | 2161.5118 | 2164.1721 | -2.6603 | 1 | 6.58e+02 | 9 | 9 | 9 | KWVTFISLLFLSSAYS.R |
| 200372 | 359 | 1130.2270 | 2258.4394 | 2259.0154 | -0.5759 | 129 | 2.7e-10 | 1 | 1 | 1 | K.EFNAETTFHADICTLSEK.E |
| 200508 | 206 | 754.1395 | 2259.3967 | 2259.0154 | 0.3813 | 70 | 0.00023 | 1 | 1 | 1 | K.EFNAETTFHADICTLSEK.E |
| 207858 | | 776.4385 | 2326.2937 | 2330.0705 | -3.7768 | 1 | 7.4e+02 | 10 | 10 | 10 | K.LCTVATLRATYGELADCCCK.Q |
| 208116 | | 777.2828 | 2328.8266 | 2328.1062 | 0.7203 | 1 | 5.6e+02 | 9 | 9 | 9 | K.NYQEAQDVFLGSLFYYSR.R |
| 218180 | | 813.7062 | 2438.0968 | 2441.2148 | -3.1180 | 1 | 4.88e+02 | 8 | 8 | 8 | R.MSQTFPNADFIEITKRLATLTK.V |
| 223059 | 90 | 830.6788 | 2489.0146 | 2489.2777 | -0.2631 | 102 | 1.3e-07 | 1 | 1 | 1 | K.ALVLIAFAQYLQCCPFEDHVK.L |
| 223217 | 49 | 1246.0610 | 2490.1074 | 2489.2777 | 0.8290 | 91 | 1.4e-06 | 1 | 1 | 1 | K.ALVLIAFAQYLQCCPFEDHVK.L |
| 223347 | | 831.4108 | 2491.2106 | 2495.2631 | -4.0525 | 4 | 7.8e+02 | 6 | 6 | 6 | K.GLVLIAPSOHLQCCPYEERVK.L |
| 223474 | | 831.8135 | 2492.4187 | 2492.1893 | 0.2293 | 1 | 9.23e+02 | 4 | 4 | 4 | K.AETTFHADICTLPEDEKQIK.K |
| 225560 | 6 | 1259.5270 | 2517.0394 | 2517.2066 | -0.1671 | 87 | 3.4e-06 | 1 | 1 | 1 | R.MPCAEDYLSVVLNQLCVLHEK.T |
| 225590 | 21 | 840.1052 | 2517.2938 | 2517.2066 | 0.0872 | 107 | 4.2e-08 | 1 | 1 | 1 | R.MPCAEDYLSVVLNQLCVLHEK.T |
| 225705 | | 840.4322 | 2518.2748 | 2517.0610 | 1.2138 | 1 | 15 | 55 | 1 | 1 | K.TYETTLERCCAADPHECYAK.V |
| 226938 | 23 | 1267.2640 | 2532.5134 | 2533.2015 | -0.6880 | 114 | 7.3e-09 | 1 | 1 | 1 | R.MPCAEDYLSVVLNQLCVLHEK.T + Oxidation (M) |
| 227008 | 39 | 845.3395 | 2532.9967 | 2533.2015 | -0.2048 | 102 | 1.2e-07 | 1 | 1 | 1 | R.MPCAEDYLSVVLNQLCVLHEK.T + Oxidation (M) |
| 229970 | | 855.0638 | 2562.1696 | 2565.2097 | -3.0402 | 1 | 5.67e+02 | 1 | 1 | 1 | K.EDPPACYATVDFKQPLVDEPK.N |
| 232089 | 20 | 862.2865 | 2583.8377 | 2584.1105 | -0.2728 | 1 | 65 | 0.00054 | 1 | 1 | K.VHTECHGDLLECADDRDLAK.Y |
| 232114 | 7 | 1293.0030 | 2583.9914 | 2584.1105 | -0.1190 | 1 | 70 | 0.00018 | 1 | 1 | K.VHTECHGDLLECADDRDLAK.Y |
| 233455 | 1 | 867.4965 | 2599.4677 | 2603.2910 | -3.8233 | 1 | 11.14e+02 | 2 | 2 | 2 | R.MPCTEDYLSLILNRLCVLHEK.T |
| 235680 | | 877.8392 | 2630.4958 | 2628.1683 | 2.3274 | 1 | 19 | 24 | 1 | 1 | K.LVNEVTEFAKTCVAESAENCDK.S |
| 237131 | 7 | 1325.5300 | 2649.0434 | 2649.2567 | -0.2113 | 103 | 8.4e-08 | 1 | 1 | 1 | R.LVREVDVNTAFHDEETFLK.K |

Transferring data from frill...

Scrolling down for this match, we can see a summary of the 201 peptide matches to the albumin proteins. You'll see on the left there are a lot of cases where more than one spectrum matched to the same peptide.

We found that many people believed that bold red meant a significant match, so we've now given in and that's what it means in Mascot 2.3.

Much of this should be familiar. You can still click on the query number to get the peptide view. However, there is no yellow popup window when you hover the mouse over the query number. To get the same information, click on the rank number

Mascot MS/MS Search Results (20080424/F002926.dat) - Mozilla Firefox

ALBU EQUAS 2275.1 70490 372 Serum albumin precursor - Equus asinus (Donkey)

201 peptides

| Query Dups | Observed | Mr(expt) | Mr(calc) | Delta Miss Score | Expect | Rank | 1 2 Peptide |
|-------------------------------|----------|-----------|-----------|------------------|---------|------|--|
| CBS00155-02-02-12.3491.2.dta | | | | | | | |
| Score > 46 indicates identity | | | | | | | |
| Score > 33 indicates homology | | | | | | | |
| 179721 | 23 | 1043.2390 | 2084.4634 | 2085.8303 | -1.3668 | 97 | 4.9e-07 1 ■ K.VHTECCGDLLECADDR.A |
| | | | | | -1.6022 | 14 | 96 2 SYSPMETIGGGIILIDVPVFK + Oxidation (M) |
| | | | | | -1.4899 | 13 | 1.4e+02 3 CVDTSMGFPLEGLVMGTR + Oxidation (M) |
| | | | | | -1.6312 | 10 | 2.3e+02 4 LISEEVGPVGDIIITNFDIR |
| | | | | | -1.5631 | 9 | 3.1e+02 5 LDEFELERHITQOAR |
| | | | | | -1.6094 | 8 | 4e+02 6 ESTIATMGTSLTDHVKILR |
| | | | | | 3.4242 | 7 | 4.9e+02 7 NNIEHMLIGGLVLAAMTK + 2 Oxidation (M) |
| | | | | | -1.5631 | 7 | 5e+02 8 LESPGGMVHGVLAAQLQR + Oxidation (M) |
| | | | | | -1.6617 | 7 | 5.4e+02 9 SPSIFHQVHLTLQYLLK |
| | | | | | -1.5189 | 6 | 6.3e+02 10 MGAQLMVDHLMIEVDTR + Oxidation (M) |
| 179823 | 26 | 696.1702 | 2085.4080 | 2085.8303 | -0.3415 | 65 | 0.00081 1 ■ K.VHTECCGDLLECADDR.A |
| 180937 | | 899.7211 | 2096.1415 | 2094.9746 | 1.1669 | 6 | 6.7e+02 6 ■ R.DELPADLNPLEHDFVDEK.E |
| 182923 | | 705.3275 | 2112.9607 | 2112.8775 | 0.0831 | 1 | 2 1.4e+03 9 ■ K.VHTECCGDLLECADDR.A |
| 183017 | | 1058.0150 | 2114.0154 | 2112.8775 | 1.1379 | 1 | 0 2.3e+03 1 ■ K.VHTECCGDLLECADDR.A |
| 183784 | | 708.1685 | 2121.4837 | 2123.9802 | -2.4965 | 1 | 19 29 1 ■ K.AAFTECCQAADKACLPLK.L |
| 184508 | | 1064.2730 | 2126.5314 | 2130.1877 | -3.6563 | 1 | 3 1.1e+03 2 ■ M.KWVTFISLILLFSSAYS.R |
| 186016 | | 1069.5870 | 2137.1594 | 2132.8085 | 4.3509 | 1 | 2 1.5e+03 6 ■ R.ETYGMDADCEKQEPER.N + Oxidation (M) |
| 188763 | | 721.5112 | 2161.5118 | 2164.1721 | -2.6603 | 1 | 6 5.8e+02 9 ■ KWVTFISLILLFSSAYS.R |
| 200372 | 359 | 1130.2270 | 2258.4394 | 2259.0154 | -0.5759 | 129 | 2.7e-10 1 ■ K.EFHAETFFHADICTLSEK.E |
| 200508 | 206 | 754.1395 | 2259.3967 | 2259.0154 | 0.3813 | 70 | 0.00023 1 ■ K.EFHAETFFHADICTLSEK.E |
| 207858 | | 776.4385 | 2326.2937 | 2330.0705 | -3.7768 | 1 | 7 4.4e+02 10 ■ K.LCTVATLRATYGLADCCCK.Q |
| 208116 | | 777.2828 | 2328.8266 | 2328.1062 | 0.7203 | 1 | 5 6.4e+02 9 ■ R.NYQEARQVFLGSLFYEYS.R |
| 218180 | | 813.7062 | 2438.0968 | 2441.2148 | -3.1180 | 1 | 4 8.8e+02 8 ■ R.MSQTFPNADFAEITRLATDLTK.V |
| 223059 | 90 | 830.6788 | 2489.0146 | 2489.2777 | -0.2631 | 102 | 1.3e-07 1 ■ K.ALVLIAFAQYLQCPFEDHVK.L |
| 223217 | 49 | 1246.0610 | 2490.1074 | 2489.2777 | 0.8298 | 91 | 1.4e-06 1 ■ K.ALVLIAFAQYLQCPFEDHVK.L |
| 223347 | | 831.4108 | 2491.2106 | 2495.2631 | -4.0525 | 4 | 7.8e+02 6 ■ K.GLVLIAFASHLQCCPYEEHVK.L |
| 223474 | | 831.8135 | 2492.4187 | 2492.1893 | 0.2293 | 1 | 9 2.3e+02 4 ■ K.AETFFHADICTLPEDEKIQK.K |

+7 subset proteins

http://frill/mascot_2_3_beta/cgi/master_results2.pl?20080424/F002926.dat;ismudpk/proteins/1/queries/2/1/182923/9

You will see that this information is now included ‘in-line’ which means that you can display more than one at a time – something that was impossible with the popup windows. Again, we are using AJAX (shorthand for asynchronous JavaScript and XML) to just load these parts on demand.

| Score | Mass | Matches | |
|-------|-------------|---------|------------|
| 1.1 | ALBU HUMAN | 74522.0 | 71317 4597 |
| 1.2 | ALBU PIC | 2640.0 | 71643 417 |
| | ALBU HORSE | 2275.1 | 70550 375 |
| | ALBU GERBIL | 721.8 | 70892 247 |
| | ALBU SHEEP | 404.8 | 71139 176 |
| | ALBU EQUAS | 2275.1 | 70490 372 |

| Query Dups | Observed | Mr(expt) | Mr(calc) | Delta Miss Score | Expect Rank | 1 | 2 | Peptide | | |
|------------|----------|-----------|-----------|------------------|-------------|----|---------|---------|---|-----------------|
| 43279 | 5 | 1128.7790 | 1127.7717 | 1127.6914 | 0.0803 | 1 | 47 | 0.092 | 1 | K.KQTALVELVK.H |
| 44140 | | 569.0507 | 1136.0868 | 1136.4808 | -0.3940 | 11 | 3.1e+02 | 2 | 2 | K.EACFAEEGPK.L |
| 44207 | 14 | 1137.8360 | 1136.8287 | 1137.4907 | -0.6620 | 38 | 0.72 | 1 | 1 | CCTESLVNR |
| 44326 | 77 | 570.1287 | 1138.2428 | 1137.4907 | 0.7522 | 57 | 0.0083 | 1 | 1 | CCTESLVNR |
| 44573 | | 1141.6930 | 1140.6857 | 1140.6866 | -0.0009 | 1 | 49 | 0.049 | 1 | K.KLVAAQAALGL.- |
| 44611 | 16 | 571.4742 | 1140.9338 | 1140.6866 | 0.2472 | 1 | 90 | 3.9e-06 | 1 | K.KLVAAQAALGL.- |
| 45411 | 33 | 575.1782 | 1148.3418 | 1148.6077 | -0.2659 | 71 | 0.00036 | 1 | 1 | LVNEVTEPAK |
| 45415 | | 575.1957 | 1148.3768 | 1148.5686 | -0.1918 | 1 | 11 | 3e+02 | 2 | R.DAHKSEVAHR.F |
| 45441 | 23 | 1149.6160 | 1148.6087 | 1148.6077 | 0.0010 | 36 | 1 | 1 | 1 | LVNEVTEPAK |
| 49914 | 2 | 597.3976 | 1192.7806 | 1190.6659 | 2.1147 | 1 | 18 | 58 | 2 | K.LVREVFPAK.T |
| 50374 | 15 | 599.9357 | 1197.8568 | 1197.5336 | 0.3233 | 1 | 54 | 0.017 | 1 | K.ETCFAEEGKK.L |
| 50639 | | 400.7912 | 1199.3518 | 1197.5336 | 1.8182 | 1 | 12 | 2.4e+02 | 2 | K.ETCFAEEGKK.L |
| 54765 | 13 | 409.5202 | 1225.5388 | 1225.5979 | -0.0591 | 1 | 37 | 0.69 | 1 | R.FKDLGEENFK.A |
| 54774 | 6 | 1226.5920 | 1225.5847 | 1225.5979 | -0.0132 | 1 | 57 | 0.0071 | 1 | R.FKDLGEENFK.A |
| 54862 | 47 | 614.1812 | 1226.3478 | 1225.5979 | 0.7500 | 1 | 63 | 0.0018 | 1 | R.FKDLGEENFK.A |
| 55952 | | 412.6049 | 1234.7929 | 1231.6053 | 3.1876 | 10 | 3.5e+02 | 8 | 8 | K.ACCDKPLLQK.S |
| 58500 | 1 | 627.3412 | 1252.6678 | 1251.6499 | 1.0179 | 1 | 20 | 36 | 4 | R.FPKAFAEVSK.L |
| 64884 | | 432.8069 | 1295.3989 | 1295.6972 | -0.2984 | 1 | 9 | 5e+02 | 3 | R.LAKTYTTELEK.C |
| 65001 | 1 | 648.8746 | 1295.7346 | 1295.6972 | 0.0374 | 1 | 22 | 22 | 1 | R.LAKTYTTELEK.C |
| 66534 | 5 | 435.8438 | 1304.5096 | 1304.6104 | -0.1009 | 30 | 3.5 | 1 | 1 | K.ECCEKPLEK.S |
| 66549 | 1 | 1305.5660 | 1304.5587 | 1304.6104 | -0.0517 | 38 | 0.56 | 1 | 1 | K.ECCEKPLEK.S |
| 66688 | 28 | 653.7496 | 1305.4846 | 1304.6104 | 0.8742 | 51 | 0.029 | 1 | 1 | K.ECCEKPLEK.S |

The little square boxes need some explanation. The number at the top of the column signifies the protein number – so in this case 1 is the human albumin and 2 is the pig, horse, donkey albumin. For query 50374, you can see that this peptide was just found in the human protein. A grey box shows that a match was found in some, but not all of the proteins in the family member.

As you can see, this new report is a replacement for both the peptide and select summary. There is no longer a requirement for separate reports for small and large searches.

Select multiple fasta files for searching

Why:

- Best to concatenate a few of your own sequences onto the end of SwissProt or NCBIInr
- Want to search SwissProt and Trembl, or a species database and contaminants.

Concatenating fasta files not easy because:

- Files are often huge
- May need to also update the reference file
- Different accession formats

In Mascot 2.3, you can select multiple databases for searching in a single search.

It's not advisable to search databases with a single, or just a few entries. We recommend that you add your own sequences to the end of a reasonable size database such as SwissProt or NCBIInr so that the statistics are more reasonable. Also, a common requirement is to search a species specific database and some common contaminants.

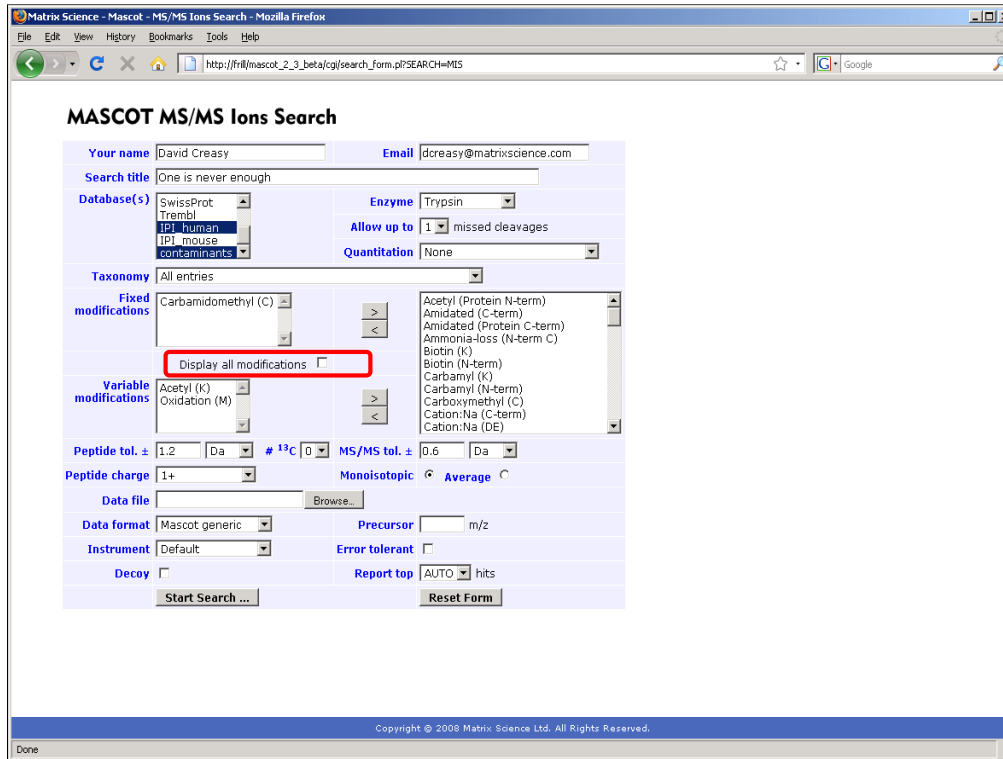
These points illustrate why the new feature will make it easier to maintain the databases.

Dealing with multi GB files is never easy.

It's also difficult to add a sequence to a database with a reference file, such as SwissProt, because the reference file needs to be updated too.

You need to make sure that the accession string format is similar in both databases and that there are no duplicate accessions.

Finally, if you update the large public databases such as NCBIInr, then you need to concatenate additional sequences to the end of the file after every update.



You can now simply select more than one database at a time.

We've also changed the way that you select modifications. It's now very easy to see which modifications have been selected, and there's a checkbox to show all the modifications.

Select multiple fasta files for searching

Databases : **1:** SwissProt 51.6
2: IPI_human 20081014
(in total 331,981 sequences; 125,104,298 residues)

| | | | | | |
|------------|-------------------------------|-------|-------|----|-----|
| 1.1 | 1:PPB1_HUMAN | 786.8 | 58259 | 17 | All |
| | 2:IPI00007289 | 786.8 | 59938 | 17 | Tæ |
| 1.2 | 1:PPBN_HUMAN | 605.8 | 57656 | 12 | All |
| | 2:IPI00290380 | 605.8 | 57626 | 12 | Tæ |
| 1.3 | 1:PPBI_HUMAN | 100.9 | 57119 | 2 | In |
| | 2:IPI00298622 | 100.9 | 57119 | 2 | Tæ |

MASCOT : Version 2.3

© 2009 Matrix Science

MATRIX
SCIENCE

At the top of the report, there is a list of the databases that were searched and the total number of sequences and residues.

In this case, the search was just against SwissProt and IPI_human. The databases are numbered sequentially, and these numbers are used to refer to the database in the body of the report.

Support for new HUPO PSI formats

HUPO Proteomics Standards Initiative

Founded in April 2002

<http://psidev.info/>

Academic and industrial support

Aims:

- Create minimum reporting standards
- Enable easier transfer of proteomics data

MASCOT : *Version 2.3*

© 2009 Matrix Science

MATRIX
SCIENCE

The HUPO Proteomics standard initiative was founded in April 2002 and has had consistent support from a number of academic groups and mass spec instrument and software vendors. Matrix Science has been involved throughout this period.

The main aims are to create standards which recommend the minimum information about a proteomics experiment that should be reported. From these, XML schema have been developed which should allow data to be more easily moved between different proteomics applications and repositories.

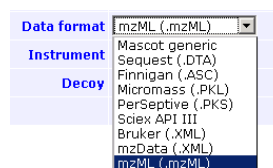
HUPO PSI mzML format

Mascot 2.2 can search mzData files

Mascot Distiller can export mzData file

Mascot Distiller can import mzXML files

Mascot 2.3 has support for mzML 1.1



MASCOT : Version 2.3

© 2009 Matrix Science



The first mass spectrometry based schema to be developed by the PSI was the mzData format. This has been used by a number of vendors and Mascot has included support for this format for many years.

However, the ISB has had a 'competing' standard called mzXML. A couple of years ago, it was agreed to merge the mzXML and mzData standard, and this has now been achieved. Almost. Version 1.0 had a few shortcomings, and it was agreed to fix these in version 1.1. The review period for 1.1 is almost over, and we've written the parsers for mzML 1.1 assuming that there will be no further changes. See the PSI web site for tools that can output mzML files

Unfortunately, there's not much to show for quite a lot of work, just an option on the list of file formats. However, we've tested with a number of files and we do pass a lot of the additional information provided in the mzML file through to the Mascot results file – perfect for developers building their own pipelines.

We'll be adding support for mzML to Mascot Distiller later in the year.

HUPO PSI mzIdentML

Export format

Current formats include pepXML, protXML, Mascot XML, X!Tandem XML, .out ...

Single format should help

PSI started developing mzIdent in 2004

More 'general' format analysisXML - failed

| | |
|---------------------------|-----------|
| Export format | mzIdentML |
| Significance threshold p< | XML |
| Max. number of hits | CSV |
| Protein scoring | mzIdentML |
| | pepXML |
| | DTASelect |
| | standard |

MASCOT : Version 2.3

© 2009 Matrix Science



The other PSI format that is relevant for Mascot is the mzIdentML format.

This is an export format and there is currently no real standard. All the search engines output in different formats.

A single format would certainly help for repositories such as Pride and for submission to journals.

The PSI started developing mzIdent in early 2004 when it was claimed that it would be “functional by the year's end”.

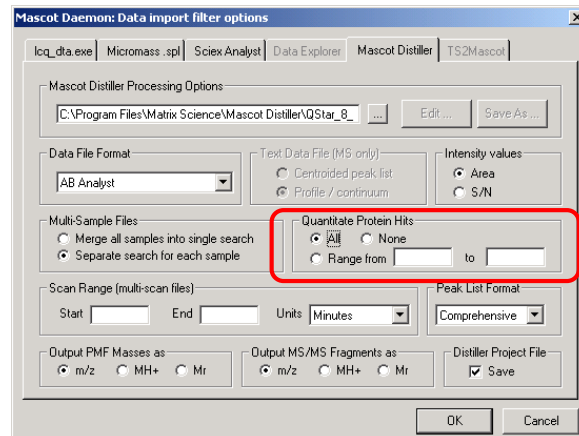
However, a few people thought that this was an opportunity to create something more general and useable by other proteomics processing such as 2D gels, and even chromatography. For this reason, the name was changed to analysisXML

However, after a couple of years, it became clear that this wasn't going to succeed, so there was a change back to reporting for just mass spec protein identification. A 60 day review period has just ended, and minor changes are currently being implemented. One of the requests for change was to the name analysisXML which was considered to be too general, so it is now called mzIdentML. The 'ML' incidentally stands for modelling language.

Again, there's not much to say or show, but it is in Mascot 2.3 and a good number of the example instance documents on the PSI web site are from Mascot.

Mascot Daemon

Add support for Distiller quantitation



MASCOT : Version 2.3

© 2009 Matrix Science



One of the major changes for Daemon 2.3 is to enable the automation of quantitation using Mascot Distiller. Those people who currently use Mascot Daemon with the Mascot Distiller import filter will know that you can automate the peak processing using the Distiller libraries and the search the data automatically. Daemon saves the Distiller project file which can then be opened in Mascot Distiller. If you are performing a quantitation experiment, you then need to open the project file in Distiller and press the quantitate button – fine for one or two files, but tedious otherwise.

So, it is now possible to choose to perform the quantitation automatically from here to automate the whole process

Mascot Daemon

- Automatically remove additional header lines after the first occurrence
- Daemon should specify username/password separate from URL
- Mascot.dll memory leak
- Daemon now supports an https (ssl enabled) Mascot server.

MASCOT : Version 2.3

© 2009 Matrix Science



There's a few other important changes in Daemon.

The first occurs when you chose to produce mgf files from an instrument data system and then want Daemon to merge all the files together. If the data system outputs header lines, then when these are merged, the header lines appear in the middle of the mgf file which causes an error. Daemon now strips the header lines after the first file.

If you are required to enter a username and password to access your Mascot server, you used to have to enter this on the url, which is then easily visible to anyone else. There's now a separate place to enter this.

A number of people have had problems with a memory leak in Mascot.dll for Analyst files. Unfortunately, it appears that the problem was not in Mascot.dll, but in Analyst itself. It seemed that there was unlikely to be a fix for this in the near future, so we've played some clever tricks to work around the problem in Mascot Daemon

Finally, Mascot Daemon now supports Mascot running on https.

Semi-supervised learning for peptide identification from shotgun proteomics datasets

Lukas Kall¹, Jesse D Canterbury¹, Jason Weston², William Stafford Noble^{1,3} & Michael J MacCoss¹

Shotgun proteomics uses liquid chromatography–tandem mass spectrometry to identify proteins in complex biological samples. We describe an algorithm, called Percolator, for improving the rate of confident peptide identifications from a collection of tandem mass spectra. Percolator uses semi-supervised machine learning to discriminate between correct and decoy spectrum identifications, correctly assigning peptides to 17% more spectra

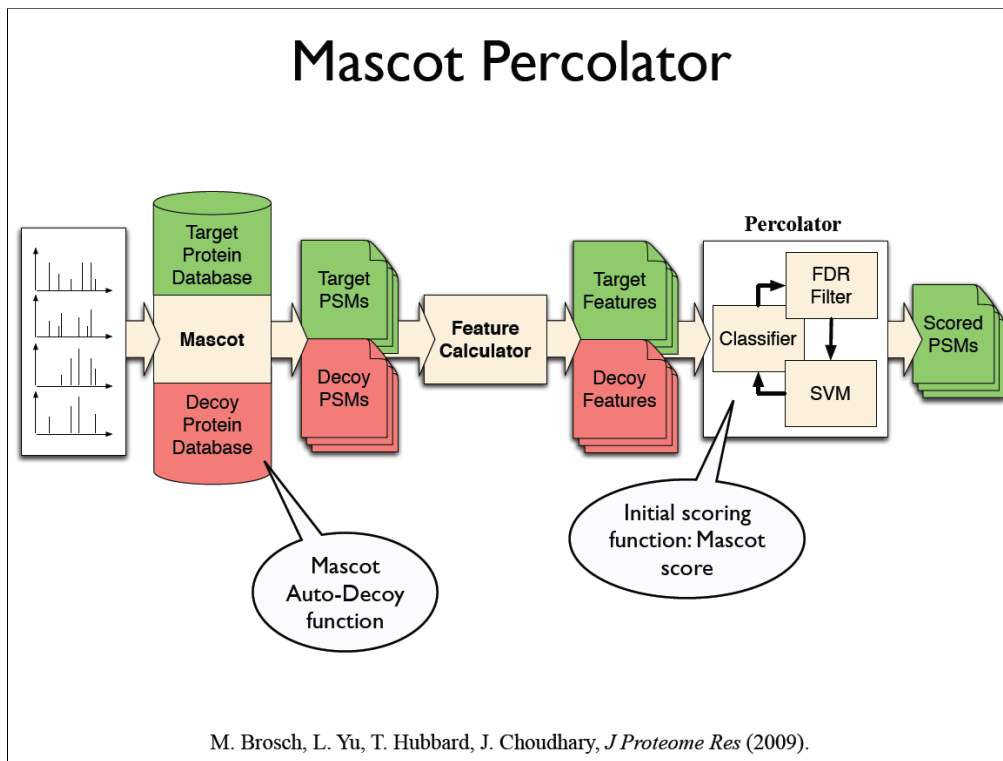
matches (PSMs) that exceed a given threshold⁶. This approach allows the user to adjust the score threshold to obtain a target false discovery rate.

Because most database search algorithms return multiple scores (for example, $XCorr$, Sp , and DC_n for SEQUEST), most proteomics studies apply separate thresholds to each score. Using multiple orthogonal score criteria is useful for eliminating false discoveries that might exceed one threshold but not another. However, in most cases these orthogonal scores are considered independently, ignoring the benefits that can be obtained if the features are considered jointly.

An alternative approach is to use machine learning methods to re-rank the PSMs and then set a threshold automatically in the re-ranked list^{4,7}. This approach uses a supervised classification algorithm to discriminate between correct and incorrect PSMs. Each PSM is characterized by a fixed-length vector of features, and the relative weights of the individual features are learned from a training set of manually curated PSMs. This approach provides substantially greater confidence in peptide identification than using SEQUEST alone; however, obtaining a high-quality training set is complicated.

ing Group <http://www.nature.com/naturemethods>

We've heard earlier from Markus about Percolator which was developed by Lukas Kall and in Mike MacCoss.



Markus described how Mascot Percolator provides a neat interface between Percolator and Mascot version 2.2. This provides alternative and improved scoring in cases where you perform a decoy search and have sufficient data for the results to be reliable.

To use Mascot Percolator with Mascot 2.2, you'll need to download the appropriate files and run Mascot percolator. You then have a choice of how to view the results. As Markus has shown, you can run a script to modify the Mascot results file and then view the scores in the standard Mascot results, or you can combine the results using spreadsheets.

In Mascot 2.3, we will simply be providing an easier to use interface for this. We will be shipping the percolator binaries and providing the infrastructure to run percolator automatically.

Mascot Percolator

Configuration in mascot.dat:

```
DisplayValues MS ME  
  DisplayPercolatorValues MS PEP
```

PS = Percolator score

PQ = Percolator q value

PEP = Percolator PEP.

MASCOT : Version 2.3

© 2009 Matrix Science



In mascot.dat, you'll be able to select what values to display in the standard Mascot reports. In this case, the Mascot score and Percolator PEP (probability that individual match with this Percolator score is random match) will be displayed. If percolator 'fails' in any way, for example if you didn't search against a decoy database, then the DisplayValues will be used instead.

It will be possible to display the Percolator score, the percolator q value or the percolator PEP value.

64 Bit Windows support

- 64 bit Parser support in 2.2 for Windows
- Full 64 bit support for Linux in Mascot 2.2
- EST databases need 64 bit support
- Ability to lock / map more databases in memory.

MASCOT : *Version 2.3*

© 2009 Matrix Science

MATRIX
SCIENCE

We are also now including true support for 64 bit Windows.

In version 2.2 we added support for 64 bit Mascot Parser, which means that you can open very large results files if you install on a 64 bit platform.

Mascot 2.2 also provided full support for 64 bit Linux.

The most important issues resolved by adding 64 bit support are the ability to use the huge EST databases and the ability to lock more databases in memory.

Other changes

- Support more than 4 cores
- Change minimum ms-ms fragment tolerance in search form from 0.01 Da
- Increase limit for number of databases from 64 to 'unlimited'
- Better homology threshold
- Support higher charge states
- Changes to quantitation editor for Distiller.

MASCOT : Version 2.3

© 2009 Matrix Science



There have been a couple of patch releases to Mascot 2.2 to resolve issues with the more recent Intel processors. However, there is a limit of 4 cores per cpu in Mascot 2.2. Mascot 2.3 will provide support for processors with more than 4 cores, but you'll need one license for each 4 cores. A system with an 8 core processor, will therefore require a 2 cpu license.

A number of people have thought that Mascot can't work with a fragment tolerance less than 0.01 Da. This is simply because the search form doesn't allow you to enter a value less than this unless you use millimass units. This is changed in Mascot 2.3

We've increased the limit for the number of databases from 64 to a configurable 'unlimited' maximum value.

In some cases it wasn't possible to calculate a homology threshold, particularly with more accurate precursor masses, so we've been able to improve this.

We also now support higher charge states for the precursor which helps with top down experiments. The limit has been increased from 8+ (or 8-) to an unlimited value.

Mascot Distiller 2.3

- Label Free - replicate protocol
- Label Free - Average protocol
- Free viewer!
- Option to use 'zoom' or 'enhanced resolution' scans for standard traps
- Normalise on selected protein
- Only quantify unique peptides.

MASCOT : Version 2.3

© 2009 Matrix Science



The new version of Distiller supports two protocols for label free quantitation.

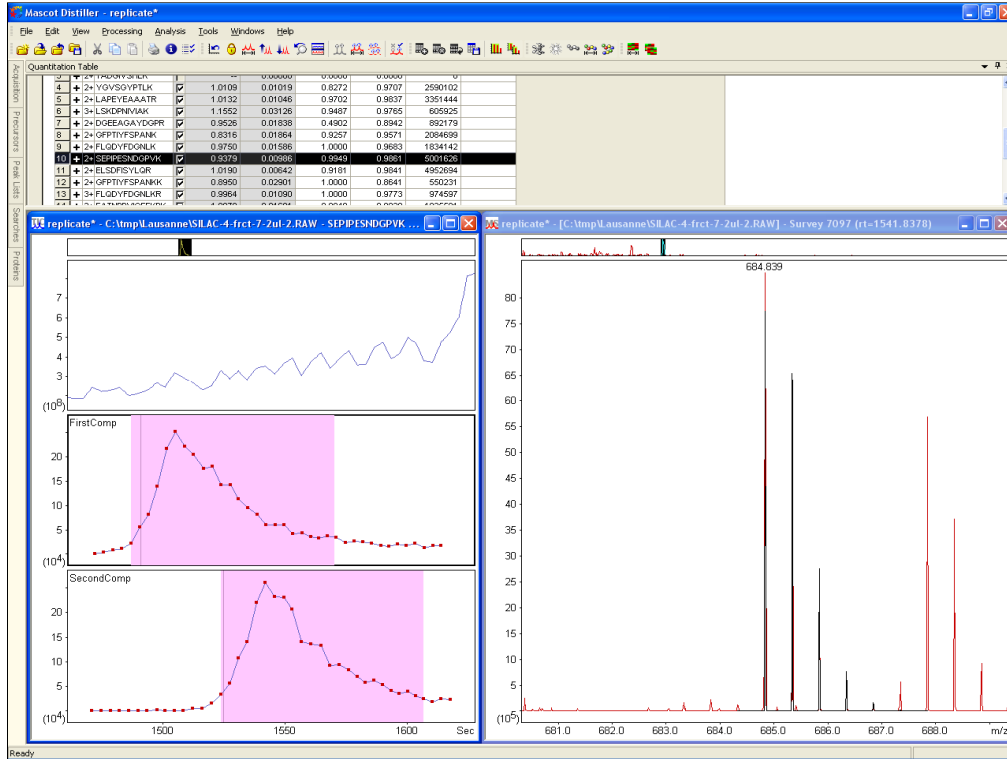
The first, is what we call the replicate protocol, and is used to determine the relative abundance of proteins between two or more samples. There will typically be a different raw file for each sample.

The Average protocol, is label-free, absolute quantitation for the proteins in a mixture in a single sample.

If you want a colleague who doesn't have a Mascot Distiller license to be able to view results, just send them the raw data, the project file and then get them to download and install Distiller. Without a license, it runs as a viewer. If you have a 30 day evaluation license, then Distiller becomes a free viewer when the license expires.

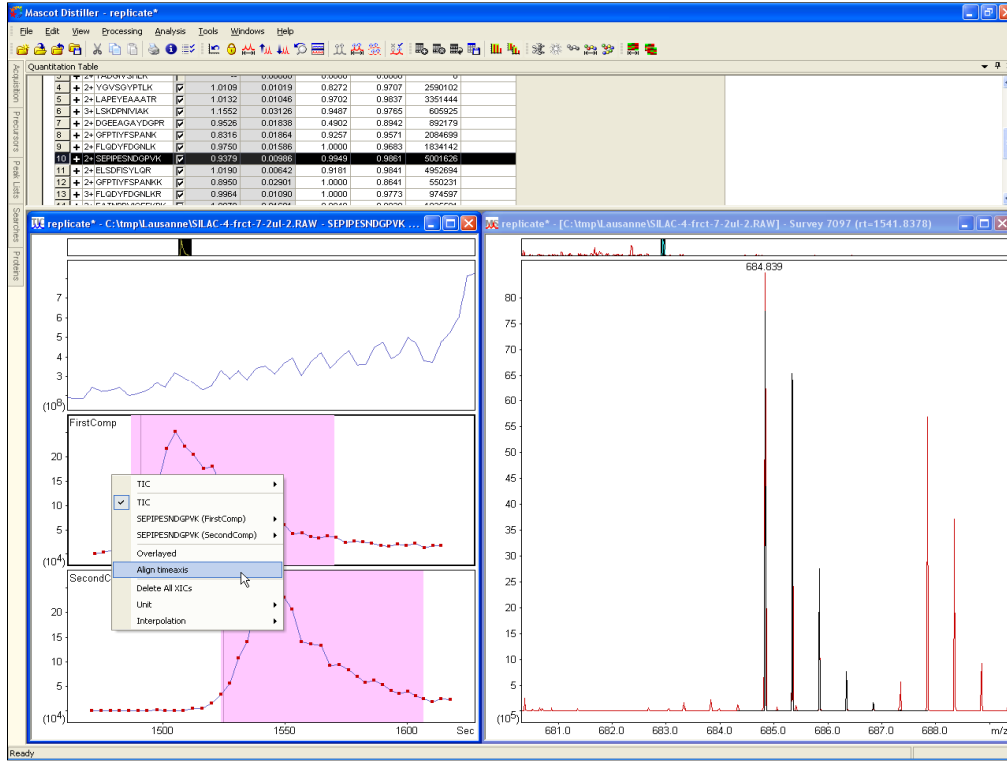
With a standard ion trap, the survey scans are often not high enough resolution for accurate quantitation. In Distiller 2.3, we've added the option to use zoom scans or enhanced resolution scans as they are called on some instruments.

We've also added the option to normalise on a selected protein

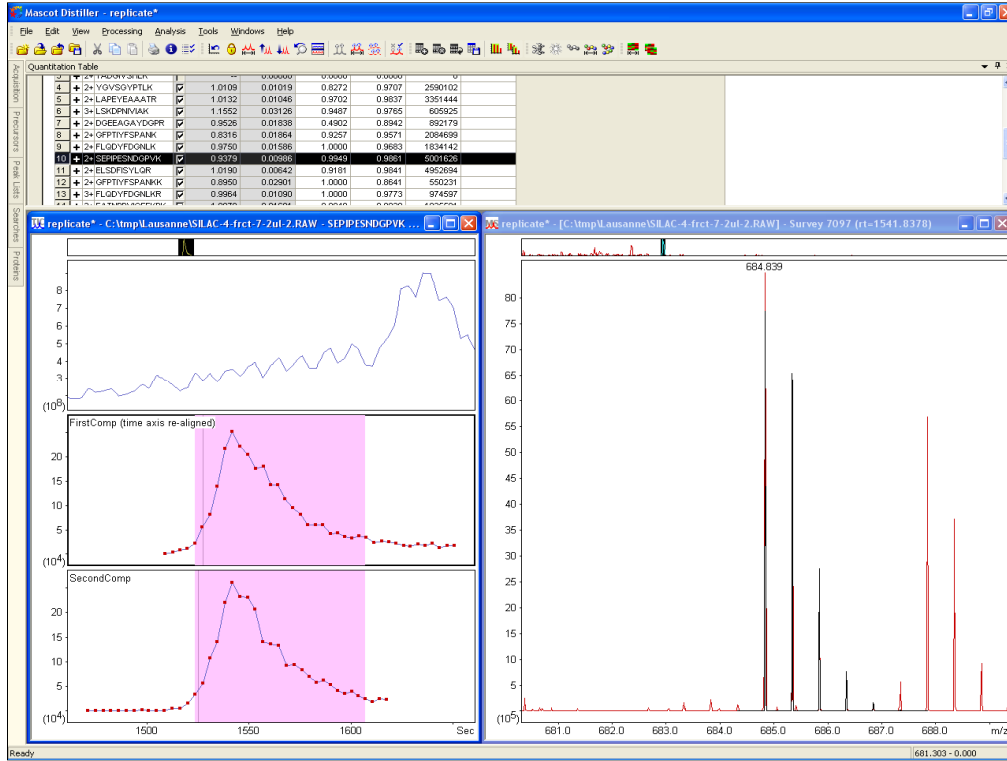


This is an example of label free using what we the replicate protocol. For clarity, there are just two raw files here, but it is of course possible to have 10 or 20 files. This method relies on getting the Mascot peptide matches first, and as with labelled methods, there is no need to get a match from both the data sets.

In this particular case, there is a shift of about a minute. The alignment is surprisingly good provided chromatograms don't get too far apart. This shows an XIC peak on a true time axis

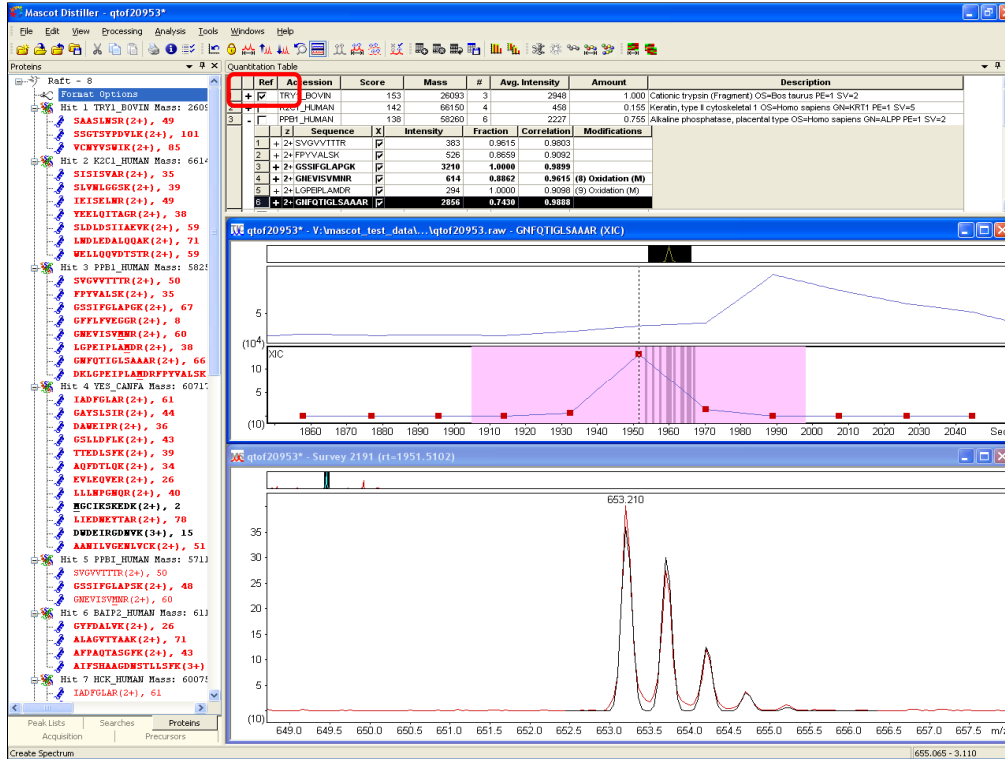


And there is a context menu to display aligned XICs



We can see here that the alignment is not bad at all.

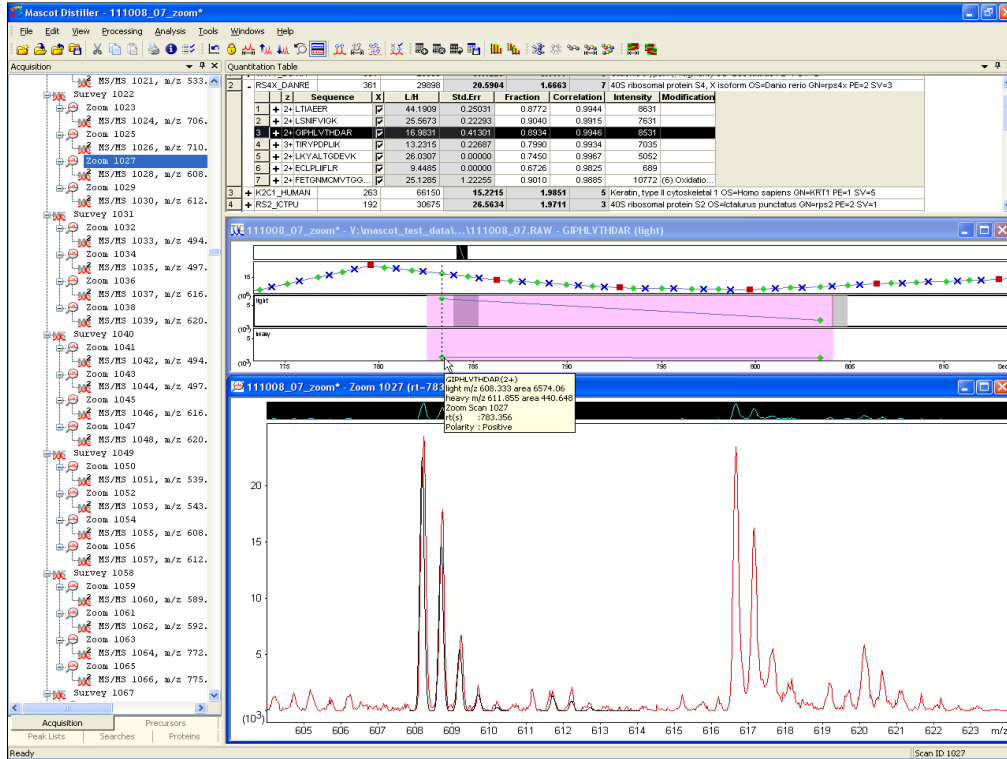
It will of course fail if XIC peaks are miles apart and there are no Mascot matches to tie them together. If there are matches to the peptides in both files, it is of course more likely to succeed.



The other label free protocol that we now support is average.

The Average protocol, is label-free, absolute quantitation for the proteins in a mixture based on the application of a rule to the intensities of extracted ion chromatograms (XICs) for the peptide matches in a database search result. The method was first described by Silva, J. C., and a group from Waters. Their observation was that the average MS signal response for the three most intense tryptic peptides per mole of protein was constant within a coefficient of variation of less than $\pm 10\%$.

Amounts are shown relative to reference protein, which can be selected using checkbox. Rule here is that amount corresponds to summed intensities of 3 most intense peptide matches per protein. Alternative to spectral counting, but should be more accurate because it is using the survey scans.



With a standard ion trap, the survey scans are often not high enough resolution for accurate quantitation. In Distiller 2.3, we've added the option to use zoom scans or enhanced resolution scans as they are called on some instruments.

This is an example of standard LTQ data off zoom scans using SILAC labelling.

Quantitation Table

| | Accession | Score | Mass | ML | SD(geo) | # | HL | SD(geo) | # | |
|---|------------|-------|-------|--------|---------|---|--------|---------|---|-----------------|
| 1 | gi 4507357 | 743 | 22549 | 1.0549 | 1.0808 | 7 | 1.0044 | 1.1013 | 7 | transgelin 2 [F |
| 2 | gi 118090 | 638 | 22786 | 0.9820 | 1.0797 | 5 | 0.9108 | 1.1343 | 5 | Peptidyl-prolyl |
| 3 | gi 4505591 | 515 | 22325 | 0.9827 | 1.0699 | 3 | 0.0841 | 1.1362 | 3 | peroxiredoxin |
| 4 | gi 2554831 | 444 | 23556 | 0.9813 | 1.0348 | 3 | 0.9684 | 1.0582 | 3 | Chain A, Crys |
| 5 | gi 2204207 | 408 | 23596 | 0.9813 | 1.0348 | 3 | 0.9684 | 1.0582 | 3 | glutathione S- |
| 6 | gi 9955007 | 337 | 21910 | 0.9812 | 1.1204 | 5 | 0.9510 | 1.1382 | 5 | Chain A, Thior |

Quantitation Method: SILAC R-6 R-10 [MD]

Method

- Constraint search:
- Protein Ratio Type: median
- Protein Score: median
- Report Detail:
- Show subsets: 0.00
- Require Bolded:
- Minimum Peptides: 2
- Significance Threshold: 0.05

Protocol

- Components
- Report Ratios
- Integration
- Quality
- Outlier Removal
- Normalisation

Normalisation Method: median

Normalisation Protein: gi|118090

Normalisation Protein
List of protein accessions

Quantitation Table [Normalised:median]

| | Accession | Score | Mass | ML | SD(geo) | # | HL | SD(geo) | # | |
|---|------------|-------|-------|--------|---------|---|--------|---------|---|-------------|
| 1 | gi 4507357 | 743 | 22549 | 1.0549 | 1.0808 | 7 | 1.0044 | 1.1013 | 7 | transgelin |
| 2 | gi 118090 | 638 | 22786 | 1.0000 | 1.1582 | 6 | 1.0000 | 1.1343 | 5 | Peptidyl-ty |
| 3 | gi 4505591 | 515 | 22325 | 1.0828 | 1.0699 | 3 | 0.0841 | 1.1362 | 3 | peroxirec |
| 4 | gi 2554831 | 444 | 23556 | 1.0011 | 1.0349 | 3 | 1.0631 | 1.0582 | 3 | Chain A, |
| 5 | gi 2204207 | 408 | 23596 | 1.0011 | 1.0349 | 3 | 1.0631 | 1.0582 | 3 | glutathior |
| 6 | gi 9955007 | 337 | 21910 | 0.9812 | 1.1204 | 5 | 0.9510 | 1.1382 | 5 | Chain A, |

MASCOT

MATRIX SCIENCE

It's now possible to normalise on selected protein(s) or peptide(s), spiked into sample.
Before (no normalisation) and after (normalise on gi|118090)

Mascot Server 2.3

- New protein grouping
- New faster reports
- Support for HUPO PSI standards
- Percolator integration

Mascot Distiller 2.3

- Two label free protocols
- Quantitation from zoom scans
- Free viewer
- Mascot Daemon integration