

Sneak preview of new features in Mascot 2.2

MASCOT

{MATRIX}
{SCIENCE}

Quantitation - Strategy

- Some quantitation methods will be supported in Mascot 2.2
- Other methods will be added ASAP in patch releases
- Some methods will only require Mascot Server, e.g. iTRAQ, emPAI
- Others will require the raw data to be processed using Mascot Distiller, e.g. Silac, ICAT, ^{18}O

MASCOT : Mascot 2.2 - Sneak Preview

© 2006 Matrix Science



The major new piece of work in Mascot 2.2 is quantitation. It seems like there is a strong need for more systematic and reliable quantitation software. Many people are working on developing chemistries and protocols, but the software to report the numbers seems to be a little neglected.

Since identification is an essential aspect of most MS based quantitation, it seems natural to try and integrate this functionality more closely into Mascot. So, quantitation will be one of the main components in the next release of Mascot. There is still more work to be done, but I wanted to give you a preview of what's in the pipeline.

We intend to support as many of the popular methods as possible. Only some of these will be finished in time for the initial release of Mascot 2.2. Other will be added as soon as possible in patch releases.

For some methods, like iTRAQ and emPAI, all of the information required for quantitation is contained in the peak list.

Most methods require additional information from the raw data file, either because it is necessary to integrate the elution profile of each peptide or because information is required for multiple peaks in the survey scan. These methods will require that the raw data files are processed using Mascot Distiller, and the Mascot search results will be imported into Distiller to generate a quantitation report.

Quantitation

Simplicity
for user

MASCOT : Mascot 2.2 - Sneak Preview

© 2006 Matrix Science



We wanted to keep the user interface simple. Quantitation adds a huge number of choices and parameters, but there is no point in exposing all of these in the search form.

The approach we have chosen is encapsulate these choices and parameters into named quantitation methods. This means that the search form has just a single new control, which replaces the old ICAT checkbox.

Methods that have [MD] at the end are the ones that require Mascot Distiller

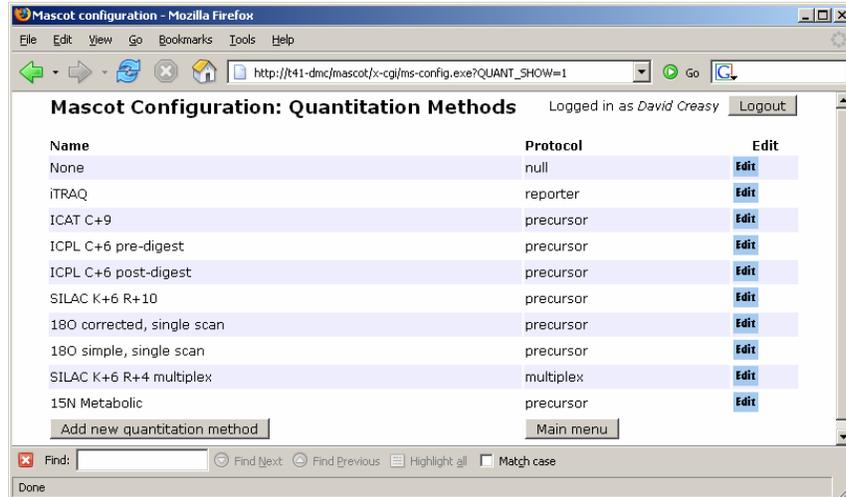
Quantitation

The quantitation methods are defined in a single XML configuration file

- quantitation.xml
- Browser based editor
- Add new methods as required
- Used by Mascot Server and Mascot Distiller

The configuration file that encapsulates the choices and parameters for each quantitation method is called quantitation.xml. This is an XML file, and there is a browser based editor for modifying methods and creating new ones. quantitation.xml lives on the Mascot server and is read by both the search engine, and Mascot Distiller. Other applications can also read the file remotely using Mascot Parser.

Mascot 2.2 - Quantitation editor

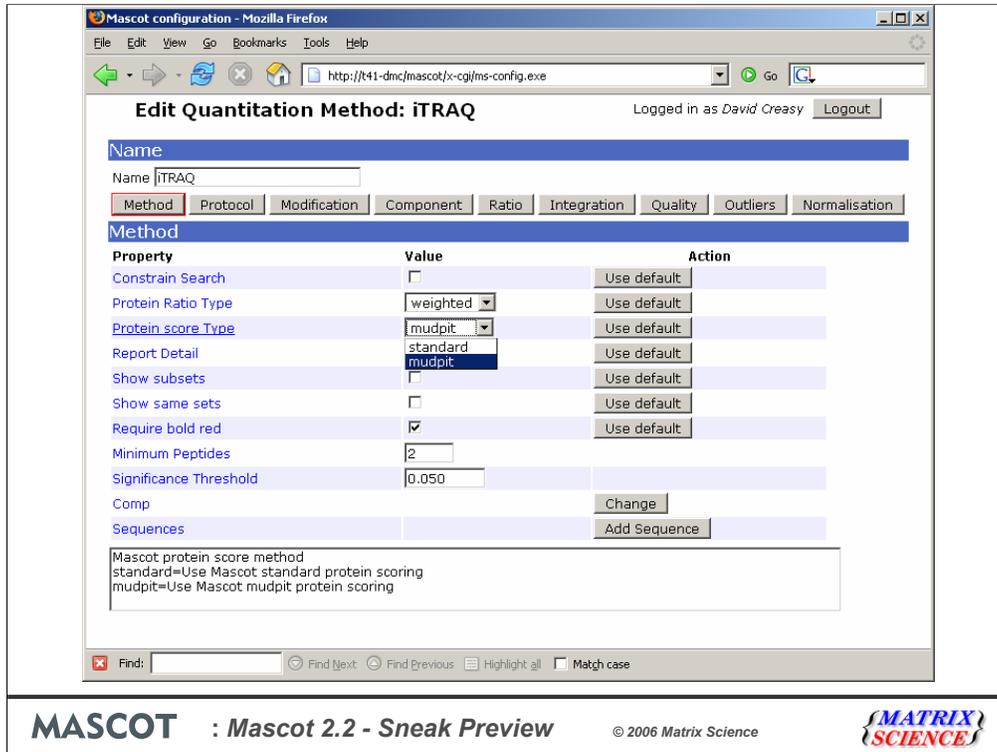


MASCOT : Mascot 2.2 - Sneak Preview

© 2006 Matrix Science



When you first open the quantitation editor, there will be a list of all the currently defined modifications. You can also of course add a new one from here. I'll show an example with ITRAQ. Clicking on the Edit button brings up this page:

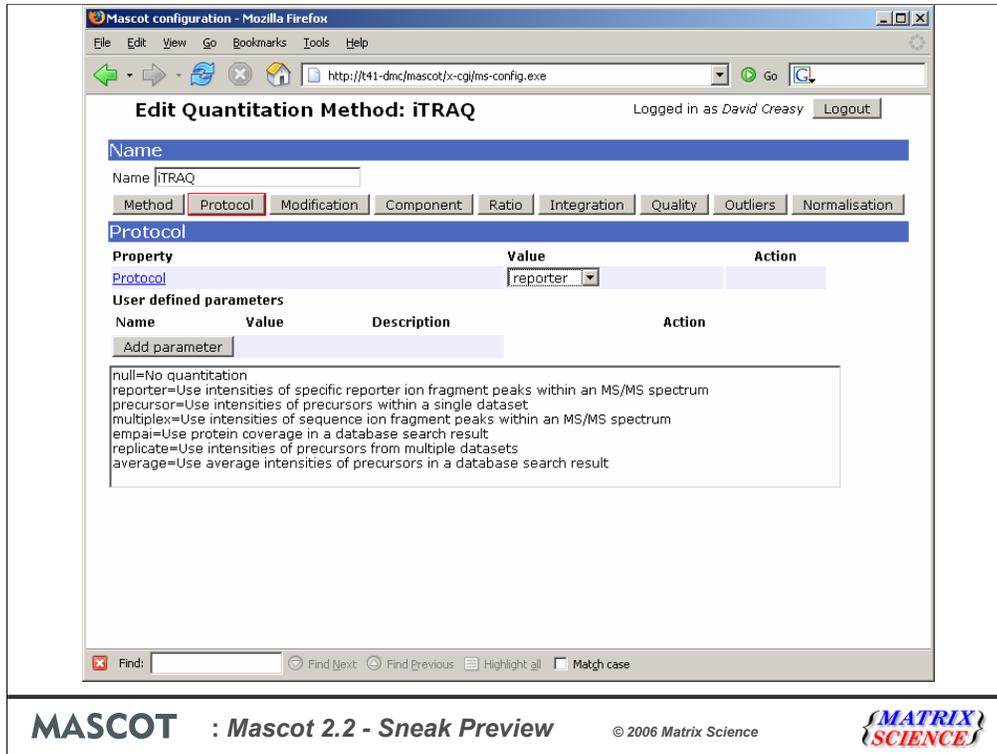


You'll see that there are a lot of different parameters, split into groups. So, there are settings for the method, protocol, modifications, etc. Some techniques won't have entries for some of these settings.

As with the database maintenance and security administration tool, holding the mouse over any of these items gives further help. So, we can see, for example that the help for constrain search says that "If any modification group specifies exclusive mode, then apply this constraint during the search so that only matches that can be used for quantitation will be returned." This isn't a sensible option for iTRAQ where the modifications would normally be specified as Fixed

The protein ration type is the method of calculating the ratio for the protein from the individual peptide ratios. You can specify whether to use the median, the average or a wighted average.

For the protein scoring, either standard or MudPIT scoring can be chosen.



The protocol selection is perhaps the most important. It's a simple drop down list with the following options:

Quantitation - Protocols

Protocol	Description	Examples
null	No quantitation	
reporter	Specific reporter ion peaks within a single MS/MS spectrum	iTRAQ
precursor	Extracted ion chromatograms for related precursors within a single dataset	ICAT, SILAC, ¹⁸ O, ICPL, AQUA, Metabolic
Multiplex (Neubert et. al.)	Pairs of sequence ion fragment peaks within a single MS/MS spectrum	SILAC, ¹⁸ O
Empai (Ishihama et. al.)	Protein coverage in a database search result	Label-free
replicate	Extracted ion chromatograms for identical precursors across multiple datasets	Label-free
average	Extracted ion chromatograms for selected peptides per protein within a single dataset	Label-free

MASCOT : Mascot 2.2 - Sneak Preview

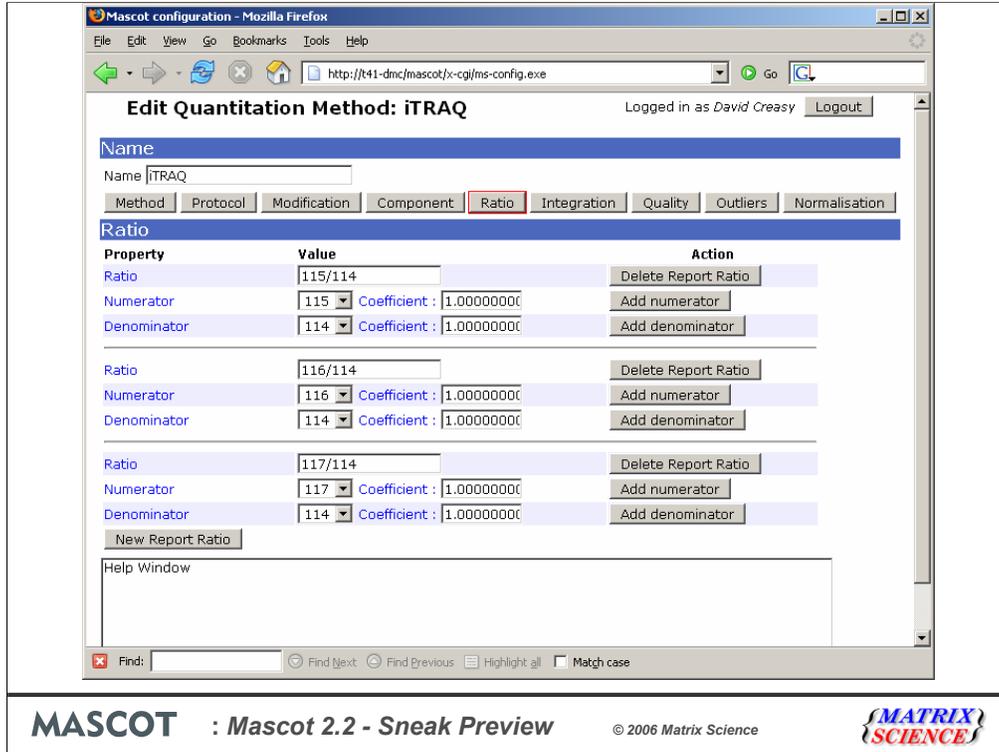
© 2006 Matrix Science



At this time, we have identified 6 protocols plus null, which means no quantitation.

MASCOT : Mascot 2.2 - Sneak Preview © 2006 Matrix Science **MATRIX SCIENCE**

The contents of the component tab depend on the protocol. So, with iTRAQ for example, since we are using reporter ions, we have to specify the m/z value for each of the four reporter ions. Also, we have the option to enter some correction information that is supplied with the reagents.



And for the ratios, we are simply choosing to report the ratio of the 115, 116 and 117 peaks to the 114 peak. Obviously, it's easy for people to select other ratios if desired.

Edit Quantitation Method: iTRAQ Logged in as *David Creasy* [Logout](#)

Name:

Method | Protocol | Modification | Component | Ratio | Integration | **Quality** | Outliers | Normalisation

Quality

Property	Value	Action
Minimum precursor charge	<input type="text" value="0"/>	
Isolated precursor	<input type="checkbox"/>	
Minimum A1	<input type="text" value="0.000"/>	
Precursor Threshold Type	<input type="text" value="maximum expect"/>	
Precursor Threshold Value	<input type="text" value="0.0500"/>	

User defined parameters

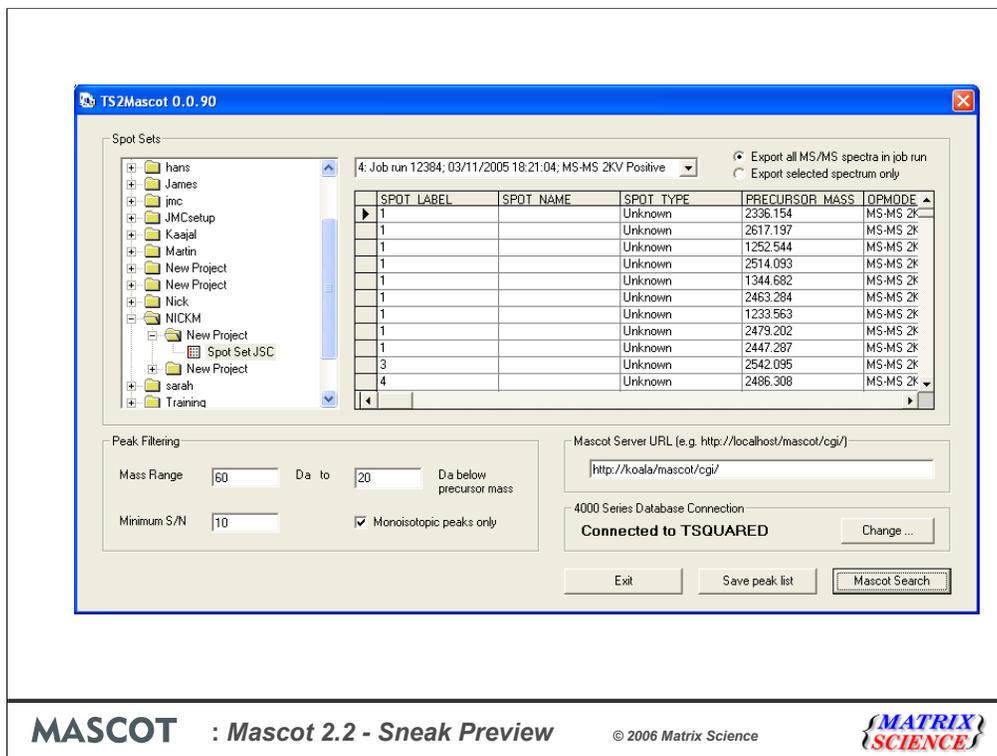
Name	Value	Description	Action
<input type="button" value="Add parameter"/>			

Significance threshold for peptide matches
 minimum score=Significance threshold for peptide matches is a minimum score
 maximum expect=Significance threshold for peptide matches is a maximum expectation value
 at least identity=Significance threshold for peptide matches is score at or above identity threshold
 at least homology=Significance threshold for peptide matches is score at or above homology threshold

Find: Match case

MASCOT : Mascot 2.2 - Sneak Preview © 2006 Matrix Science

The quality page allows you to select which peptides will be used for calculating ratios - so it would be normal to disregard anything with a 1 in 20 chance of it being a random match.



Even the simpler methods can become complicated. With iTRAQ, for example, we found that the peak list exported from the 4000 series data system or submitted to Mascot from GPS Explorer did not have the correct peak areas for the reporter ions. The numbers are different from those used within GPS Explorer for quantitation. We have had to write our own application to export a suitable peak list from the Oracle database.

So, for iTRAQ, you might choose Mascot Search ...

Quantitation

Matrix Science - Mascot - MS/MS Ions Search - Microsoft Internet Explorer

Address: http://k41-jsc/mascot/cgi/search_form.pFORMVER=2&SEARCH=MS

MASCOT MS/MS Ions Search

Your name: John Cottrell Email: jcottrell@matrixscience.com

Search title: 41-42-43-43.5; Job run 12384; MS-MS 2KV Positive

Database: NCBIInr

Taxonomy: Proteobacteria (purple bacteria)

Enzyme: Trypsin/P Allow up to: 2 missed cleavages

Fixed modifications: Acetyl (K), Acetyl (N-term), Amide (C-term), Biotin (K), Biotin (N-term)

Variable modifications: Acetyl (K), Acetyl (N-term), Amide (C-term), Biotin (K), Biotin (N-term)

Protein mass: kDa Quantitation: iTRAQ all ratio to 114

Peptide tol. ±: 0.3 Da MS/MS tol. ±: 0.3 Da

Peptide charge: 1+ Monoisotopic: Average

Data file: nc\LOCALS~1\Temp\ts230.tmp

Data format: Mascot generic Precursor: m/z

Instrument: MALDI-TOF-TOF

Overview: Report top: AUTO hits

Start Search ... Reset Form

Copyright © 2005 Matrix Science Ltd. All Rights Reserved.

Local intranet

MASCOT : Mascot 2.2 - Sneak Preview

© 2006 Matrix Science



Which brings up the search form. We choose an appropriate quantitation method. We don't need to specify the iTRAQ modifications, because these are contained in the quantitation method. Submit the search...

Select Summary Report (Phi/Caulo 41 42 43 43.5; Job run 12384; MS MS 2KV Positive) - Microsoft Internet Explorer

Address: http://k11-jsc/mascot/cgi/master_results.pl?file=../data/20060510/F001747.dat

1. [q1113423172](#) Mass: 72428 Score: 1952 Queries matched: 51
TonB-dependent receptor [Caulobacter crescentus CB15]

Quantitation:	Ratio	Value	StdDev	Weighted SD
	115/114	1.01	0.15	0.03
	116/114	0.98	0.17	0.03
Settings	117/114	1.31	0.33	0.06

Query	Observed	Mr(expt)	Mr(calc)	Delta	Miss	Score	Expect	Rank	115/114	116/114	117/114	Peptide
515	1029.4802	1028.4729	1028.5161	-0.0431	0	47	0.0093	1	1.11	1.23	1.58	R.AGYSQGR.A 516 517 2339
2020	1151.5807	1150.5734	1150.5974	-0.0240	0	37	0.095	1	----	----	----	R.ENYPAR.A
3219	1267.6404	1266.6331	1266.6690	-0.0358	0	57	0.00092	1	1.03	1.11	1.59	K.DLAHLATYR.W
3690	1314.7209	1313.7136	1313.7424	-0.0288	0	54	0.002	1	1.02	1.00	1.09	K.AQGLELAAR.P
4160	1402.7280	1401.7207	1401.7486	-0.0279	0	74	2.1e-005	1	1.05	1.19	1.36	R.SPSVAHPVAFR.T
4759	1429.7885	1428.7812	1428.8210	-0.0398	0	60	0.00054	1	1.14	1.32	1.67	R.WPVGLSTVAAR.Y
4864	1439.6415	1438.6342	1438.6558	-0.0216	0	92	2.9e-007	1	1.02	1.15	1.32	K.EADGYQAGASGR.L 4863 4865 5975 5977
5081	1466.7772	1465.7699	1465.7842	-0.0143	0	66	0.00012	1	1.05	0.81	1.27	R.AAVPHTFDAAGK.N
6524	1659.8923	1658.8850	1658.8974	-0.0123	1	52	0.0032	1	0.98	0.81	1.14	R.TRSPVAHPVAFR.T 6525
7079	1728.8743	1727.8670	1727.8851	-0.0181	0	91	4.3e-007	1	1.10	1.22	1.31	R.ASPYISETLEVYGR.V 7081
7874	1850.0276	1849.0203	1849.0301	-0.0097	0	125	1.7e-010	1	1.22	1.18	1.65	R.GAESHHIVVLDGK.L 7872 7875 7876
7945	1861.8009	1860.7936	1860.8036	-0.0100	0	67	9.4e-005	1	0.78	0.90	0.94	K.NDFGHTDSPEFGR.T
7990	1869.9182	1868.9109	1868.9220	-0.0110	1	106	1.3e-008	1	1.01	0.94	1.54	K.TGKADGYQAGASGR.L 7991 7992 7993
8131	1889.0477	1888.0404	1888.0539	-0.0135	0	112	3.6e-009	1	0.99	1.04	1.02	K.ELVYQLMAALLDGR.W 8134
9372	2133.0139	2132.0066	2132.0044	0.0022	1	73	2.6e-005	1	1.21	1.12	1.48	R.VGFATDIDREYHAR.A 9370 9371 9373 93
9580	2191.0281	2190.0208	2190.0221	-0.0012	1	109	5.8e-009	1	1.00	1.00	1.44	R.GKNDGHTDSPEFGR.T 9579 9581 9582 95
9837	2262.3064	2261.2991	2261.2986	0.0005	1	41	0.044	1	1.03	0.97	0.86	R.TKELVYQLMAALLDGR.W
10998	2504.2915	2503.2842	2503.2788	0.0055	0	80	5e-006	1	0.91	1.10	0.95	K.LNDPSSIQGGFHSGLLVDIR.I
11134	2555.3291	2554.3218	2554.3270	-0.0052	1	84	1.8e-006	1	0.98	1.05	1.30	R.VENALDKYQTLILNYGTPGR.G 11133 11135
11775	2767.3257	2766.3184	2766.3240	-0.0056	2	56	0.0011	1	1.47	1.32	2.78	R.AVYSRGRKDFDGFHTDSPEFGR.T 11776
12136	3114.6531	3113.6458	3113.6590	-0.0132	1	47	0.0076	1	1.22	0.85	1.47	K.LNDPSSIQGGFHSGLLVDIR.IEVL.R.G

2. [q1113422195](#) Mass: 110633 Score: 1599 Queries matched: 35
OmpA-related protein [Caulobacter crescentus CB15]

Quantitation:	Ratio	Value	StdDev	Weighted SD
---------------	-------	-------	--------	-------------

And back comes the report. This is how the reports will look for methods that can be handled entirely within the search engine. For methods like Silac and ICAT, the quantitation report will be generated by Mascot Distiller

Quantitation - Definition example

```
<method name="ICPL Protein 13C(6) [MDJ]" constrain_search="true" protein_ratio_type="average"
report_detail="true" description="Label pre-digest, so peptide N-terms are not labelled">
  <modification_group mode="exclusive" name="light">
    <mod_file>ICPL_light (K)</mod_file>
    <mod_file>N-ICPL_light (Protein)</mod_file>
  </modification_group>
  <modification_group mode="exclusive" name="heavy">
    <mod_file>ICPL_heavy (K)</mod_file>
    <mod_file>N-ICPL_heavy (Protein)</mod_file>
  </modification_group>
  <component name="light">
    <group_name>light</group_name>
  </component>
  <component name="heavy">
    <group_name>heavy</group_name>
  </component>
  <report_ratio name="L/H">
    <numerator_component name="light"/>
    <denominator_component name="heavy"/>
  </report_ratio>
  <quality isolated_precursor="false"/>
  <integration method="trapezium" source="survey"/>
  <protocol>
    <precursor allow_mass_time_match="false"/>
  </protocol>
</method>
```

MASCOT : Mascot 2.2 - Sneak Preview

© 2006 Matrix Science



Here is an example of a quantitation method. This is for ICPL labelling of a protein prior to performing the digest.

Note that we are introducing a new way of treating modifications. Exclusive means that peptides must conform to one of these modification groups. You cannot have a peptide that has a mixture of the heavy and light labels.

It will also be possible to specify user defined modifications and isotope substitutions for metabolic labelling

The mechanics of quantitation, which is to say, where the numbers come from, is defined by the protocol. In this case, we use a protocol called precursor.

Configuration file editors

Configuration files on the Mascot server:

- enzymes
- fragmentation_rules
- quantitation
- taxonomy
- unimod (replaces mod_file and masses file)

If you want to get as many identifications as possible, as efficiently as possible, you might come up with a strategy similar to this.

Enzyme file editor

Title	Cutter	Cleavage	Restrict	Independent	Semispecific	Edit
Trypsin	C-Term	KR	P	no	no	Edit Del
Arg-C	C-Term	R	P	no	no	Edit Del
Asp-N	N-Term	BD		no	no	Edit Del
Asp-N_ambic	N-Term	DE		no	no	Edit Del
Chymotrypsin	C-Term	FLWY	P	no	no	Edit Del
CNBr	C-Term	M		no	no	Edit Del
CNBr+Trypsin	C-Term	M		no	no	Edit Del
Formic_acid	C-Term	KR	P	no	no	Edit Del
Formic_acid	C-Term	D		no	no	Edit Del
Use C	C-Term	K	D	no	no	Edit Del

MASCOT : Mascot 2.2 - Sneak Preview

© 2006 Matrix Science



If you want to get as many identifications as possible, as efficiently as possible, you might come up with a strategy similar to this.

64 Bit Support

Linux

- Mascot 2.1.03 and later supports 64 bit Perl
- Mascot 2.1.x - 32 bit apps that run in 64 bit Linux
- Mascot 2.2 will have 'true' 64 bit version

Windows

- Mascot 2.1 - 64 bit not supported because of installer
- 64 bit perl not readily available for Windows
- Mascot 2.2 - new installer, but still only 32 bit apps.

MASCOT : *Mascot 2.2 - Sneak Preview*

© 2006 Matrix Science



In Mascot 2.1.03, we added better support for 64 bit versions of Linux by adding a 64 bit version of Mascot Parser. Most 64 bit versions of Linux come with a pre-installed 64 bit perl, so it is quite difficult to install 32 bit perl on some Linux distributions. However, the search engine and all the other applications are still 32 bit which isn't really an issue. Mascot 2.2 will have the option to install 64 bit applications, and this may give a slight performance increase - possibly as much as 10%.

We currently don't support 64 bit Windows. We have confirmed that the current version of Mascot will run on 64 bit Windows XP, but it is very hard to install because we are still using an old installer.

The other issue is that we use ActiveState Perl, and they don't have a 64 bit version of perl available yet.

Mascot 2.2 will have a new installer, and 64 bit Windows will be supported, but we will only have 32 bit versions of the Mascot Parser and the applications. When a 64 bit version of ActiveState Perl is available, we will review this - (I just checked again yesterday and see that they have just released a beta version.)

Option to search randomised entries

- Check box on the search form
- Sequences based on average AA or base composition for the whole database
- For each sequence, create a randomised sequence of same length
- Option to view results from random sequences
- Summary stats:

	NCBI nr	Decoy	False Positive rate
Peptide matches above identity threshold	4970	25	0.50%
Peptide matches above homology threshold	6600	330	5.00%

MASCOT : Mascot 2.2 - Sneak Preview

© 2006 Matrix Science



We are delighted that the proteomics community has recently been taking an active interest in validating search results by searching a data set against a randomised or reversed database using all the same parameters. Since 1999 we have had the facility on our public web site to search a randomised database. We've also recently put up a help page on our public web site that discusses some of the issues, and has a link to a free script to create a randomised or reversed database. In Mascot 2.2, we are going to make this even easier.

There will be a checkbox on search form. If you tick this, the search will take twice as long.

We will take average AA or base composition for database as a whole from stats file. The new sequences that are generated will all have this average composition.

When searching each database entry, independently create and search a randomised sequence of the same length and of this average composition.

The results from these searches will be accumulated separately in two new sections in the result file. This means that it will be easy for us to add a link on the main results page to view the results from searching the randomised database.

The standard master results report will include summary stats as additional block in header, as shown here:

And lots more . . .

PMF searches now use intensity information

Module for top down searches from FT

Support for ETD/CID searches etc:

```
BEGIN IONS  
INSTRUMENT=ETD  
.  
.  
.  
END IONS  
BEGIN IONS  
INSTRUMENT=CID  
.  
.  
.
```

MASCOT : *Mascot 2.2 - Sneak Preview*

© 2006 Matrix Science



If you want to get as many identifications as possible, as efficiently as possible, you might come up with a strategy similar to this.