# From Search Results to Publication in Nine Mouse Clicks

(sequence shortened)

**MASCOT**

*MATRIX SCIENCE*

I'd like to show you how the new features in the protein family report make it easy to generate the figures and tables needed for a publication.

To illustrate, I'm going to use data from this year's ABRF Proteome Informatics Workgroup study - iPRG2012

The sample was a yeast lysate that had some additional non-yeast proteins spiked in. It was analysed on an AB Sciex 5600 tripleTOF and both raw data and peak lists were provided. Participants were asked to search a specified database and use target decoy to report peptide matches at 1% FDR. They were also asked to characterize modifications with special emphasis on modifications not introduced by sample handling

MASCOT : *From Search Results to Publication* © *2012 Matrix Science*

I'm not going to show you all the steps required to participate in the study because it was a peptide-centric study and most of the work was in formatting the results to fit the spreadsheet template.

First of all, we need to make the Fasta database available for searching in Mascot. As you've seen in the earlier presentation, Database Manager makes this very easy. If you wanted to use Mascot's automatic target/decoy function, you would download the target only database, which contains SwissProt entries.

MASCOT : *From Search Results to Publication* © 2012 Matrix Science

We know from the sample description that the cysteine alkylation is carbamidomethyl. Usually, the only other modification I would select for a first, trial search is Met-Ox. The other settings are guesses which we will refine by looking at the results.

In particular, we can use Peptide View and Protein View to estimate mass accuracy.

| 2 | 173.0921 | 87.0497 | | | 155.0815 | 78.0444 | A | 1733.9443 | 877.4758 | 1736.9177 | 868.9625 | 1733.9357 | 868.4703 | 17 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 3 | 230.1135 | 115.5604 | | | 212.1030 | 106.5551 | G | 1682.9072 | 841.9572 | 1665.8806 | 833.4440 | 1664.8966 | 832.9519 | 16 |
| 4 | 343.1976 | 172.1024 | | | 325.1870 | 163.0972 | I | 1625.8857 | 813.4465 | 1608.8592 | 804.9332 | 1607.8751 | 804.4412 | 15 |
| 5 | 471.2562 | 236.1317 | 454.2296 | 227.6185 | 453.2456 | 227.1264 | Q | 1512.8016 | 756.9045 | 1495.7751 | 748.3912 | 1494.7911 | 747.8992 | 14 |
| 6 | 584.3402 | 292.6738 | 567.3137 | 284.1605 | 566.3297 | 283.6685 | I | 1384.7431 | 692.8752 | 1367.7165 | 684.3619 | 1366.7325 | 683.8699 | 13 |
| 7 | 683.4087 | 342.2080 | 666.3821 | 333.6947 | 665.3981 | 333.2027 | V | 1271.6590 | 636.3331 | 1254.6325 | 627.8199 | 1253.6484 | 627.3279 | 12 |
| 8 | 754.4458 | 377.7265 | 737.4192 | 369.2132 | 736.4352 | 368.7212 | A | 1172.5906 | 586.7989 | 1155.5640 | 578.2857 | 1154.5800 | 577.7937 | 11 |
| 9 | 869.4727 | 435.2400 | 852.4462 | 426.7267 | 851.4621 | 426.2347 | D | 1101.5535 | 551.2804 | 1084.5269 | 542.7671 | 1083.5429 | 542.2751 | 10 |
| 10 | 984.4997 | 492.7535 | 967.4731 | 484.2402 | 966.4891 | 483.7482 | D | 986.5265 | 493.7669 | 969.5000 | 485.2536 | 968.5160 | 484.7616 | 9 |
| 11 | 1097.5837 | 549.2955 | 1080.5572 | 540.7822 | 1079.5732 | 540.2902 | L | 871.4996 | 436.2534 | 854.4730 | 427.7402 | 853.4890 | 427.2482 | 8 |
| 12 | 1198.6314 | 599.8193 | 1181.6048 | 591.3061 | 1180.6208 | 590.8141 | T | 758.4155 | 379.7114 | 741.3890 | 371.1981 | 740.4050 | 370.7061 | 7 |
| 13 | 1297.6998 | 649.3535 | 1280.6733 | 640.8403 | 1279.6892 | 640.3483 | V | 657.3679 | 329.1876 | 640.3413 | 320.6743 | 639.3573 | 320.1823 | 6 |
| 14 | 1398.7475 | 699.8774 | 1381.7209 | 691.3641 | 1380.7369 | 690.8721 | T | 558.2994 | 279.6534 | 541.2729 | 271.1401 | 540.2889 | 270.6481 | 5 |
| 15 | 1512.7904 | 756.8988 | 1495.7639 | 748.3856 | 1494.7799 | 747.8936 | N | 457.2518 | 229.1295 | 440.2252 | 220.6162 | | | 4 |
| 16 | 1609.8432 | 805.4252 | 1592.8166 | 796.9120 | 1591.8326 | 796.4199 | P | 343.2088 | 172.1081 | 326.1823 | 163.5948 | | | 3 |
| 17 | 1680.8803 | 840.9438 | 1663.8537 | 832.4305 | 1662.8697 | 831.9385 | A | 246.1561 | 123.5817 | 229.1295 | 115.0684 | | | 2 |
| 18 | | | | | | | R | 175.1190 | 88.0631 | 158.0924 | 79.5498 | | | 1 |

RMS error 18 ppm

RMS error 18 ppm

NCBI BLAST search of TAGIQIVADDLTVTNPAR
(Parameters: blastp, nr protein database, expect=20000, no filter, PAM30)
Other BLAST web gateways

**MASCOT** : *From Search Results to Publication* © 2012 Matrix Science

{MATRIX SCIENCE}

Looks like 0.05 Da is safe for MS/MS.

For the MS errors, looks like 10 ppm is safe. We might get away with 5 ppm but, with a small database, this is going to limit the number of candidate sequences available for matching to each spectrum, so 10 ppm is a better choice.

We also see quite a few high scoring matches with 2 missed cleavages, so maybe push this up to 3.

MASCOT : *From Search Results to Publication* © 2012 Matrix Science

Repeating with the new settings, we can see that the FDR for the default setting of 5% significance threshold is approximately 3%. The iPRG study requested matches to be reported with an FDR of 1%. This is where another of the new features in Mascot 2.4 comes in useful. The 'adjust to FDR' button. Getting back to the title of the talk, let's use the first of our nine mouse clicks to obtain the required 1% FDR

In Mascot 2.3 and earlier, you had to use trial and error to adjust the FDR to a specific value, so this button is a time saver. You may also notice that the decoy sequences are reversed and not randomised. This is another new feature in Mascot 2.4. The default is reversed for MS/MS searches with enzyme specificity and randomised for no enzyme searches, but you can change these defaults if you wish.

To get a table of proteins suitable for publication, we use a second mouse click to switch to the Report Builder tab.

Lets assume we want to drop the 'one hit wonders' and only report proteins that have significant matches to at least 2 different peptide sequences

MASCOT : *From Search Results to Publication* © 2012 Matrix Science

We open up the filters section and add a suitable filter. This uses another 5 mouse clicks, so 7 in total

Click number 8 is to export as CSV and click number 9 (actually double click) is to open the CSV in excel

Now, I'm going to cheat a bit, and ignore all the keystrokes we need to use in Excel to add some formatting to the table.

MASCOT : *From Search Results to Publication* © 2012 Matrix Science

This is where the last bit of the title comes in. You may have noticed the weasel words 'sequence shortened' in technology ads. Particularly for a certain cellphone

You get the idea

And there we have it, a table of the reliably identified proteins, suitable for pasting into a publication, in just 9-ish mouse clicks

By the way, the filtering is very flexible, with lots of useful terms. Another thing that you could easily do would be to exclude proteins from the contaminants database

The columns section of Report Manager allows you to choose which columns to include and, if required, change their order

Now, the main goal of the iPRG2012 study was to characterise modifications. Quickest way to find out what modifications might be present is an error tolerant search.

MASCOT : *From Search Results to Publication* © 2012 Matrix Science

The error tolerant search discovers lots of modifications, but which ones are interesting? It would be helpful if the report included a table of the modifications that had been found together with their frequency of occurrence. I can assure you that this is on the wish list. Meanwhile, the work around is to export the results as CSV and open in Excel

Select the pep_var_mod column, containing the modifications, and choose Pivot table from the Data menu

MASCOT : *From Search Results to Publication* © 2012 Matrix Science

In the pivot table wizard, the defaults are OK

Drag and drop the pep_var_mods button to both the row fields and data items area

MASCOT : *From Search Results to Publication* © 2012 Matrix Science

And we get a table of the distinct variable modifications with a count for each. For a large search, you might want to restrict the table to just the top 50 most frequent modifications, and this is easily done (Pivot table wizard menu; Field settings; count; custom; advanced; show top 50).

The list still needs some interpretation. For example, note the presence of several mods with mass delta -57. These almost certainly indicate that carbamidomethylation is not 100% quantitative. For peptides where Cys is not modified, putting a -57 mod close by cancels out the mass difference well enough to get a decent match.

| Modification | Delta | Count |
|---|---|---|
| Ammonia-loss (N-term C) | -17.03 | 4 |
| Gln->pyro-Glu (N-term Q) | -17.03 | 33 |
| Deamidated (NQ) | 0.98 | 234 |
| Methyl (K) | 14.02 | 15 |
| Methyl (R) | 14.02 | 5 |
| Oxidation (M) | 15.99 | |
| Cation:Na (C-term) | 21.98 | 4 |
| Cation:Na (DE) | 21.98 | 37 |
| Formyl (S) | 27.99 | 32 |
| Formyl (T) | 27.99 | 20 |
| Dimethyl (R) | 28.03 | 4 |
| Ethyl (K) | 28.03 | 14 |
| Ethyl (N-term) | 28.03 | 9 |
| Dioxidation (P) | 31.99 | 8 |
| Acetyl (K) | 42.01 | 4 |
| Acetyl (N-term) | 42.01 | 7 |
| Acetyl (Protein N-term) | 42.01 | 10 |
| Acetyl (S) | 42.01 | 27 |
| Guanidinyl (K) | 42.02 | 18 |
| Guanidinyl (N-term) | 42.02 | 79 |
| Trimethyl (K) | 42.05 | 10 |
| Carbamyl (Protein N-term) | 43.01 | 24 |
| Carbamyl (S) | 43.01 | 10 |
| Nitro (Y) | 44.99 | 24 |
| Carbamidomethyl (C) | 57.02 | |
| Sulfo (STY) | 79.96 | 10 |
| Phospho (ST) | 79.97 | 168 |
| Phospho (Y) | 79.97 | 16 |

**Near isobaric modifications**
**(assuming 2000 Da peptide)**

| | |
|---|---|
| Acetyl (K) Guanidinyl (K) | 5.6 ppm |
| Acetyl (N-term) + nearby deamidation Carbamyl (N-term) | 5.6 ppm |
| Sulfo (STY) Phospho (STY) | 4.8 ppm |

**MASCOT** : *From Search Results to Publication* © 2012 Matrix Science    *MATRIX SCIENCE*

After further scrutiny, we end up with these as the believable modifications that occur 4 or more times. Although the mass accuracy of the data is excellent, there can still be ambiguities, such as whether we have acetyl or guanidinyl. In the case of sulfo and phospho, we can often decide which we have from differences in neutral loss behaviour. I'll come back to this later.

Where we go next depends on the goal of the experiment. In the case of the iPRG2012 study, it was to report as many matches as possible. Clearly, this is a slightly artificial case. In real life, we are more likely to be interested in a specific modification or a specific protein. But, how would one search for all of these modifications? You can't simply select them all as variable modifications; the combinatorial explosion would mean that all specificity was lost. However, it is highly unlikely that we will see two rare modifications on the same peptide. As long as we have Oxidation (M), Deamidated (NQ), Phospho (ST), and Carbamidomethyl (C) specified in the search as variable modifications, we shouldn't miss very much when the error tolerant search looks serially through all of the modifications in Unimod.

Note that the default Mascot configuration only allows 2 variable mods in an error tolerant search. You'll need to change the value of the MaxEtVarMods option to 4 or more to perform such a search.

For the iPRG study, the next step would be to export the results to Excel. I don't want to go into a lot of detail … there isn't time … so I'll just highlight a couple of points relating to modification characterisation.

Mascot 2.4 reports site localisation probabilities using the delta score method published in MCP by Bernard Kuster's group. Here, for example, there are 4 potential phosphorylation sites but, based on the score differences between the matches, it looks fairly clear that the site is S10. The four matches with scores of 30.9 are for Sulfation on each of these four sites. Because the Sulfo modification is lost quantitatively on MS/MS fragmentation, there is no preference for any particular site; the MS/MS is identical in all cases. For this peptide, we can be confident that the modification is phospho because we see extensive loss of 98 from the fragments, and matching these gives the higher score.

| 8 | 892.4047 | 446.7060 | 875.3781 | 438.1927 | 874.3941 | 437.7007 | E | 304.1615 | 152.5844 | 287.1350 | 144.0711 | 286.1510 | 143.5791 | 2 |
| 9 | | | | | | | R | 175.1190 | 88.0631 | 158.0924 | 79.5498 | | | 1 |

NCBI BLAST search of TVIDYNGER
(Parameters: blastp, nr protein database, expect=20000, no filter, PAM30)
Other BLAST web gateways

**All matches to this query**

| Score | Mr(calc) | Delta | Sequence | Site Analysis |
|---|---|---|---|---|
| 52.9 | 1145.4659 | -0.0011 | TVIDYNGER | Sulfo Y5 50.00% |
| 52.9 | 1145.4659 | -0.0011 | TVIDYNGER | Sulfo T1 50.00% |
| 47.6 | 1145.4754 | -0.0106 | TVIDYNGER | |
| 10.9 | 1145.4627 | 0.0020 | ISIQCRSCR | |
| 10.9 | 1145.4627 | 0.0020 | ISIQCRSCR | |
| 6.9 | 1145.4600 | 0.0048 | WELYNWR | |
| 6.1 | 1145.4580 | 0.0067 | VTCQSSELAK | |
| 6.1 | 1145.4580 | 0.0067 | VTCQSSELAK | |
| 6.1 | 1145.4580 | 0.0067 | VTCQSSELAK | |
| 5.6 | 1145.4722 | -0.0075 | ISIQCRSCR | |

Mascot: http://www.matrixscience.com/

**MASCOT** : *From Search Results to Publication* © 2012 Matrix Science

*MATRIX SCIENCE*

Here is a peptide that has sulfo as the top scoring match. There is simply nothing in the MS/MS to distinguish modification at T1 and Y5. The third match with the greater mass error is for Phospho on T1. Phospho on Y gets a very poor score, not even in the top 10, because it takes out most of the matching y ions

MASCOT : *From Search Results to Publication* © 2012 Matrix Science

A word of warning. Site localisation is often a function of the modifications selected for the search. Here, for example, is another peptide where the localisation looks excellent when we search with Phospho on S, T, and Y. But, in rare cases, other residues can be phosphorylated. Post translational modification of C, R, D, K and H are all documented in RESID and Unimod. If we were to perform a search where these unusual specificities were included …

| 7 | 794.3567 | 397.6820 | 776.3461 | 388.6767 | Y | *390.1425* | 195.5749 | 373.1159 | 187.0616 | | 2 |
| 8 | | | | | K | 227.0791 | 114.0432 | 210.0526 | 105.5299 | | 1 |

Error (Da): 0.03, 0.02, 0.01, 0, -0.01 — Mass (Da) 250, 500, 750
RMS error 10 ppm

Error (ppm): 30, 20, 10, 0 — Mass (Da) 250, 500, 750
RMS error 10 ppm

NCBI **BLAST** search of DISLSDYK
(Parameters: blastp, nr protein database, expect=20000, no filter, PAM30)
Other BLAST web gateways

**All matches to this query**

| Score | Mr(calc) | Delta | Sequence | Site Analysis |
|---|---|---|---|---|
| 44.7 | 1019.4212 | 0.0008 | DISLSDYK | Phospho K8 48.82% |
| 44.7 | 1019.4212 | 0.0008 | DISLSDYK | Phospho Y7 48.82% |
| 31.2 | 1019.4212 | 0.0008 | DISLSDYK | Phospho D6 2.18% |
| 20.0 | 1019.4212 | 0.0008 | DISLSDYK | Phospho S5 0.16% |
| 7.7 | 1019.4202 | 0.0019 | DISKKNR | |
| 7.7 | 1019.4202 | 0.0019 | DISKKNR | |
| 7.7 | 1019.4202 | 0.0019 | DISKKNR | |
| 7.7 | 1019.4202 | 0.0018 | DLSRLTR | |
| 7.7 | 1019.4202 | 0.0018 | DLSRLTR | |
| 7.7 | 1019.4202 | 0.0018 | DLSRLTR | |

Mascot: http://www.matrixscience.com/

**MASCOT** : *From Search Results to Publication*  © 2012 Matrix Science

{MATRIX} {SCIENCE}

Things are no longer so clear cut. In reality, this is highly likely to be Phospho on Y7 because Phospho K is very unusual. But, when we say we are confident that the phosphate is on Y7, we should really add "assuming the only possibilities are S, T, and Y"

## New Features in Mascot 2.4

- **Adjust to 1% FDR button**
- **Report Builder**
  - Filters
  - Columns
  - Export as CSV
- **Site localisation**
- **Text and Number Search**
- **Preferred Taxonomy**

**MASCOT** : *From Search Results to Publication* © 2012 Matrix Science

*MATRIX SCIENCE*

To summarise, we've seen practical examples of several of the new features in the Mascot 2.4 reports. The two that I didn't mention are the much enhanced text and number search facility. For example, you can search the protein family report for a modification or a mass value. Finally, the facility to set a preferred taxonomy. This wasn't relevant here, because the database was essentially yeast proteins, but in other searches, you might want to search a wide taxonomy, e.g. green plants, and where there are two proteins with equal scores, always choose the protein from (say) maize, because that is the particular subject of your research.