# Mass tolerant and error tolerant searches

## How do they compare?

**MASCOT** : *Mass tolerant searches*   *© 2016 Matrix Science*   MATRIX SCIENCE

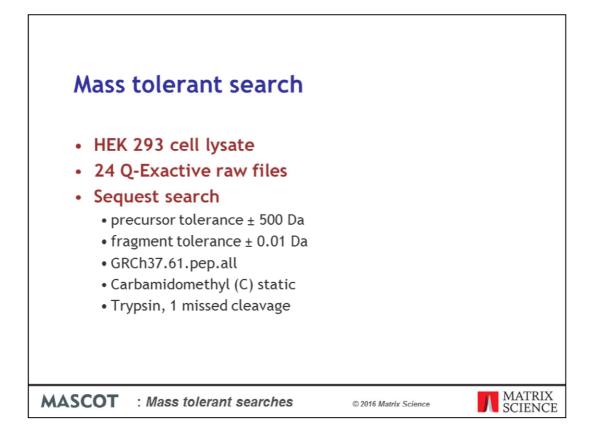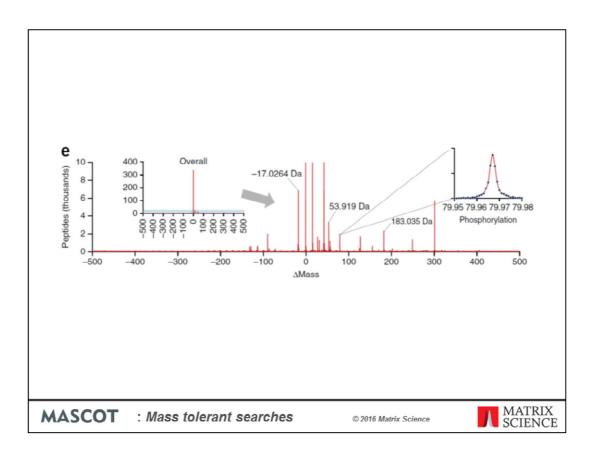MASCOT : *Mass tolerant searches* © 2016 Matrix Science MATRIX SCIENCE

Steven Gygi's lab at Harvard Medical School published this paper in Nature Biotechnology last year. It describes the use of a very wide precursor mass tolerance, +/- 500 Da, to identify modified peptides in a Sequest search. How does this approach, which the authors also call an open search, compare with a "conventional" multi-pass search, such as the Mascot error tolerant search?

The sample was a lysate of human embryonic kidney cells: 24 fractions analysed by Q-Exactive Orbitrap. Peak lists were searched against a human proteome database using Sequest. The only unusual aspect of the search was the 500 Da precursor tolerance.
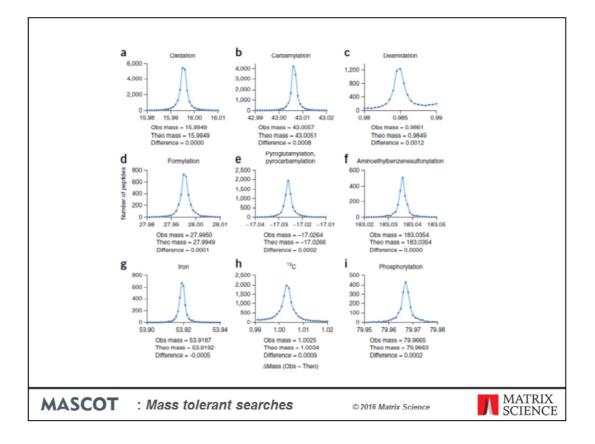
What do the results of a mass tolerant search look like? Well, it's a long list of matches, just like a regular search, except some of them have substantial differences between the calculated and observed peptide mass.

MASCOT : *Mass tolerant searches* © 2016 Matrix Science

In the paper, this is summarized using a histogram of significant PSM count against mass difference. Most of the counts are for unmodified peptides, so the y axis has been expanded to show some of the more common delta masses, many of which correspond to the 'usual suspects' – ammonia loss, oxidation, acetylation, carbamylation, phosphorylation, etc.

The search doesn't say anything about the site of the modification, which has to be determined after the search using a separate algorithm. In the paper, A-Score was used for this.

The authors used Gaussian fit analysis to divide the matches into 523 delta mass bins. This is Figure 2 from the publication, showing narrow, symmetric distributions in the delta mass distributions for selected modifications. Half widths are typically 0.005

## Error tolerant search

- **24 Q-Exactive raw files downloaded from PRIDE**
- **Peak picked by Mascot Distiller**
- **Mascot error tolerant search**
    - Enzyme                        : Trypsin/P
    - Fixed modifications           : Carbamidomethyl (C)
    - Variable modifications        : Oxidation (M)
    - Peptide mass tolerance        : ± 5 ppm (# 13C = 1)
    - Fragment mass tolerance       : ± 15 ppm
    - Max missed cleavages          : 2
    - Instrument type               : ESI-TRAP
    - Database                      : GRCh37.61.pep.all

**MASCOT** : *Mass tolerant searches*      © 2016 Matrix Science      MATRIX SCIENCE

To make a comparison with the Mascot error tolerant search, we downloaded the raw files from PRIDE and processed them into peak lists using Mascot Distiller. The same database was searched using very standard settings. Target/decoy was used to set the false discovery rate for PSMs for the first pass search to 1%

MASCOT : *Mass tolerant searches*          © 2016 Matrix Science

In the automatic error tolerant search, every protein containing one or more significant matches from the first pass search is selected for a second pass search, which uses a much wider search space: all the modifications in the Unimod database, non-specific cleavage at one peptide terminus, and all possible single amino acid substitutions. For each peptide, these possibilities are tested serially. That is, we don't look for two unsuspected modifications on the same peptide, or an unsuspected modification plus a SNP, etc.

In the result report, these additional matches are displayed with a mass delta and a tooltip showing the modifications or SNPs that fit to the delta within the specified mass tolerance. For accurate data like this, where the precursor tolerance is 5ppm, there is usually just one possibility.

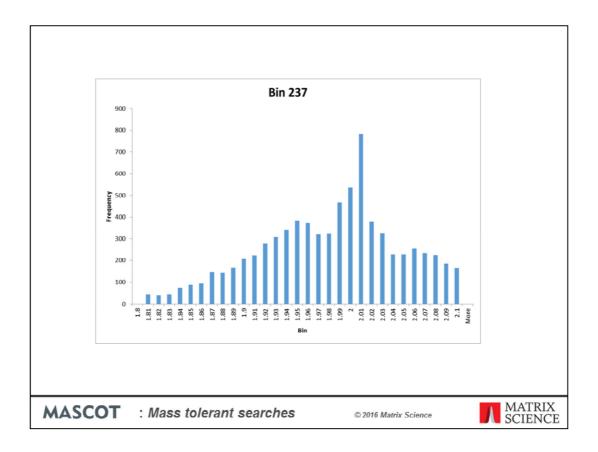We can't claim that the way the matches are reported is infallible. Sometimes, the exact site of the modification will be uncertain. Other times, the error tolerant match has a score that is only slightly higher than an unmodified peptide, and we might prefer to take the simpler explanation. But, in general, for high accuracy data, the displayed modifications represent a reasonable interpretation.

| Mass-tolerant (open) Search | | | |
|---|---|---|---|
| Bin | Delta | Count | Assignment |
| 234 | -0.0002 | 339578 | (unmodified) |
| 252 | 15.9944 | 21171 | Oxidation |
| 277 | 43.0059 | 13660 | Carbamyl |
| 236 | 1.0259 | 12741 | 13C |
| 235 | 0.9608 | 11747 | Deamidated |
| 237 | 1.9755 | 7614 | Should be 2.01, 13C2? |
| 216 | -17.0255 | 6627 | Ammonia-loss, Gln->pyro-Glu |
| 399 | 301.9864 | 5600 | ? |
| 233 | -0.9464 | 4521 | artefact |
| 287 | 53.9190 | 3326 | Cation:Fe[II] |
| 264 | 27.9946 | 3285 | Formyl |
| 232 | -1.0281 | 3185 | artefact |
| 230 | -2.0534 | 2599 | artefact |
| 269 | 31.9893 | 2561 | Dioxidation |
| 333 | 183.0367 | 2290 | AEBS |
| 254 | 16.9961 | 2030 | Oxidation+13C? |
| 189 | -89.0305 | 1934 | Met-loss+Acetyl |
| 305 | 79.9666 | 1866 | Phospho |
| 318 | 128.0964 | 1588 | Lys |
| 231 | -1.9276 | 1573 | artefact |
| 239 | 3.0216 | 1514 | 13C3? |
| 238 | 2.9008 | 1272 | artefact |
| 369 | 249.9803 | 1254 | ? |
| 292 | 57.0227 | 1108 | Carbamidomethyl |

| Error tolerant Search | | | | |
|---|---|---|---|---|
| Modification | Site | Delta | Count | Notes |
| Carbamidomethyl | C | 57.0214 | 136316 | Fixed mod in search |
| Oxidation | M | 15.9949 | 79590 | Variable mod in search |
| Non-specific cleavage | - | - | 16836 | |
| Carbamyl | N-term | 43.0058 | 13056 | |
| Gln->pyro-Glu | N-term | -17.0265 | 8094 | |
| Deamidated | N | 0.9840 | 7295 | |
| AEBS | Y | 183.0354 | 4472 | |
| Dioxidation | W | 31.9898 | 3984 | |
| Formyl | S | 27.9949 | 3761 | |
| Ammonia-loss | N-term | -17.0265 | 2919 | pyro-carbamidomethyl |
| Phospho | S | 79.9663 | 2669 | |
| AEBS | K | 183.0354 | 2529 | |
| Acetyl | N-term | 42.0106 | 2510 | |
| Formyl | T | 27.9949 | 2153 | |
| Oxidation | W | 15.9949 | 2117 | |
| Deamidated | Q | 0.9840 | 1848 | |
| Carbamyl | K | 43.0058 | 1699 | |
| Glu->Gln | E | -0.9840 | 1514 | same as amidation |
| Arg | N-term | 156.1011 | 1275 | ISD / non-specific cleavage |
| Carbamyl | T | 43.0058 | 1224 | |
| Cation:Fe[II] | D | 53.9193 | 1172 | |
| Iodo | Y | 125.8966 | 1138 | |
| Cation:Fe[II] | E | 53.9193 | 1132 | |
| Delta:H(2)C(2) | N-term | 26.01565 | 1121 | |
| Carbamyl | S | 43.0058 | 1091 | |
| Ammonia-loss | N | -17.0265 | 1030 | |

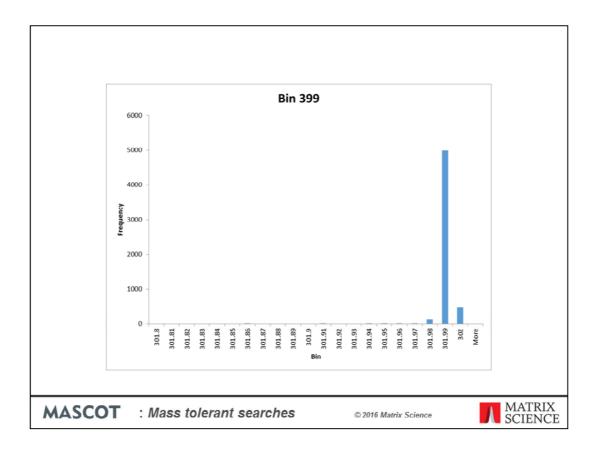MASCOT : *Mass tolerant searches*  © 2016 Matrix Science  MATRIX SCIENCE

These tables list the most abundant matches from the two types of search with an arbitrary cut-off of 1000 instances. There are some differences in the way the results are reported that are not important. For example, the table for the error tolerant search includes the fixed and variable modifications while the table for the mass tolerant search includes unmodified peptides and 13C matches.

For the mass tolerant search, the counts are independent of specificity. For example, the carbamyl count includes carbamylation of N-term, S, T, and any other sites that are susceptible to this modification. The error tolerant search reports separate counts for each specificity, although this isn't always going to be meaningful. When alternative sites are close together or when the spectrum is noisy, there may be little difference in score between two alternatives. And, of course, if there is a choice of modifications within the precursor mass tolerance, the very identity of the mod may be uncertain, although such cases will be rare for this particular search because the tolerance was 5ppm.
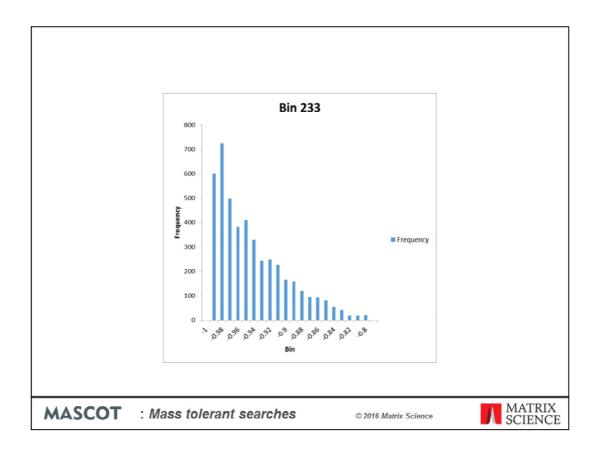
The really important point is that several of the most abundant modifications from the mass tolerant search have a question mark against them or are labelled artefact. Let's look at the first three of these: bins 237, 399, and 233

This is what the distribution looks like for bin 237. The paper reports this bin as a mass of 1.9755, which is the mean, but I suspect the spike at 2.01 is a better representative value, and can be assigned as 13C2. But what are all the other matches? It is essentially a continuum. The paper claims a peptide FDR of 0.12%, with just 625 modified peptides in total. If this is correct, and these matches are real, then each of these bins requires a different elemental composition, which is very hard to believe.
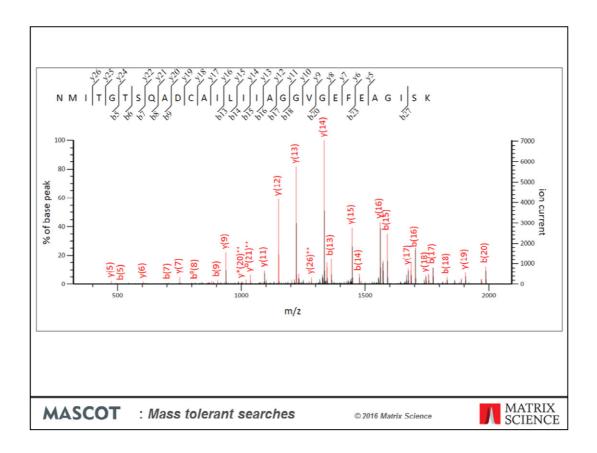
**Bin 399**

Bin 399 is a nice, clean peak. The problem is trying to figure out an assignment. It could be a combination of mods, so one approach is to look at combinations of the other high abundance modifications, but I haven't been able to come up with an assignment. The paper says nothing about this peak, even though it is the 7th most abundant modification. Does anyone have any ideas? Even with good mass accuracy, there are many possible elemental compositions for a mass of 302, and I haven't found any standard utility for listing possible formulae that includes negative counts for some elements, as may be required for a delta.

Bin 233

Bin 233 is another continuum. In the paper, it is assigned a mass of -0.95, but maybe -0.98 would be a better choice. As with the earlier example, even if you can come up with a composition for one or two of these channels, what are all the rest? These are not low level features, hidden in the grass.

The paper doesn't say anything about this issue, but we can make a good guess as to the likely cause if we consider exactly how a modification is found in a mass tolerant search

MASCOT : Mass tolerant searches    © 2016 Matrix Science    MATRIX SCIENCE

The calculated fragment masses used to test for a match are always those for the unmodified peptide. If you have a nice spectrum like this, with a good balance of b and y ions, and there is an unsuspected modification somewhere in the middle, this will take out roughly half the fragment matches. That is, the match is only based on those fragments that do not include the modified residue. If the modification was at or near a terminus, it would take out one complete series. For a modification on the amino terminus, you lose all the b ion matches and for a modification on the carboxy terminus, you lose all the y ion matches. If you have a good balance of b and y ions, this is much the same as having a modification in the middle – you lose half your matches - but if you only have one series it will give a bias. For example, if you only have y ions, then the closer the modification is to the C-terminus, the less likely you are to get any kind of match.

The critical weakness of the mass tolerant approach is that the mass of the modification comes solely from the difference between the calculated mass of the peptide and the observed mass of the precursor; the fragment masses play no part in determining the modification mass.

Let's take one of the strong matches from the continuum of bin 237. Observed m/z 629.3194 and a match to AQAALAVNISAAR. The difference between the observed mass and the calculated mass is 1.92 Da, which doesn't fit to anything in Unimod and is outside the 'allowed' range of mass defects for peptide-like molecules.

If we locate this scan in the original raw file and take a look at the precursor region of the survey scan …

MASCOT : *Mass tolerant searches*  © 2016 Matrix Science

MATRIX SCIENCE

This is what we find. The precursor with an m/z of 629.3199 is in the middle . But, we can see two other precursors, equally strong. Notice the difference between the first two: 0.96 m/z at charge 2+, corresponding to a mass difference of 1.92. It seems pretty clear that the mass tolerant search hasn't really discovered a modified peptide. The instrument was targeting 629.32 but the fragments in the MS/MS spectrum that gave the strongest match came from the precursor at 628.36. In effect, the precursor mass was 'wrong'. Since there is nothing to tie the fragment masses to the precursor mass, the error goes undetected, and a spurious modification is reported.

How often this happens is hard to say. The Gygi paper reports 185,000 modified peptides in the open search that were not found in the standard search, and it would be a mammoth task to make a forensic analysis of these. What we can say is that whenever there are overlapping precursors, there is a very real possibility of the 'wrong' mass being taken, causing the inference of a spurious modification.

Going back to the summary histogram, here we have zoomed into the central range, -20 to +100 Da, and switched to a log scale for the counts. You can see that the region that has the highest level of 'background' is the region around 0 Da. This is what you would expect for overlapping distributions of the same charge. The width of the instrument selection window is user adjustable, but I believe 4 m/z units is typical, so we are likely to see false modifications of a few da at most on unmodified peptides. However, the artefact applies equally to modified peptides. You might believe you have discovered a peptide with a delta of 40 Da or 44 Da and find it is actually a modification of 42 Da from a different precursor.

(If the overlapping distributions have different charge states, then the spurious modification could be very large, but usually these cases will fall outside the mass range studied in the paper, +/- 500 Da.)

Note that such errors are outside the scope of the FDR as estimated by target/decoy. The delta mass plays no part in the scoring and a match is counted as true or false independently of whether the delta mass is true or false.

MASCOT : *Mass tolerant searches*     © 2016 Matrix Science

The error tolerant search is much more constrained because it is looking for a fit to the modified peptide. A false modification will simply make the match worse. In this particular case, if the MS/MS spectrum was only associated with the central precursor m/z value of 629.3199, there would be no match. This is better than a false match, of course, but Distiller 2.5 and Mascot Server 2.5 introduced support for multiple precursor m/z values for a single MS/MS spectrum. This is what the Distiller peak list looks like.

The lowest m/z precursor gets the correct match to the unmodified peptide

The higher m/z value, which has charge 3+, gives a match to a completely different peptide. If you compare the fragment matches, you'll see that this is a very nice example of a chimeric spectrum. These two precursors account for all the intense fragment peaks, so it isn't surprising that we don't get a match for the middle precursor

## Mass-tolerant search limitations

- If a delta mass isn't in Unimod and isn't a combination of other abundant modifications it will be very hard to assign a composition (e.g. 249.98, 301.99)
- A continuum of spurious delta masses at low mass and near to very abundant peaks makes it hard to spot rare modifications
- Matches are weaker because only unmodified fragments are used for the match
- Matches cannot make use of known neutral loss behaviour, such as loss of 98 from phosphate

**MASCOT** : *Mass tolerant searches*      © 2016 Matrix Science      MATRIX SCIENCE

In summary, the mass tolerant search certainly shows spikes for delta masses corresponding to common modifications. But, if the mass isn't in Unimod, it will be a challenge to figure out the chemical identity.

The site of modification has to be determined separately by using a calculation such as A-Score.

A background of spurious delta masses caused by taking the wrong precursor mass makes it difficult to identify low abundance modifications.

Matches to modified peptides are weaker than in a conventional search or an error tolerant search because only half the fragments are available for matching, on average. Similarly, the matching cannot take advantage of known neutral loss behaviour

# Error tolerant search limitations

- Can only find 'known' modifications – those in Unimod
- Cannot match a peptide carrying multiple unsuspected modifications
- Requires the protein to have at least one peptide without unsuspected modifications. So, cannot match isolated peptides, e.g. endogenous peptides

**MASCOT** : *Mass tolerant searches*  © 2016 Matrix Science  MATRIX SCIENCE

One of the limitations in an error tolerant search are that it can only find 'known' modifications. As Unimod becomes more comprehensive, this becomes less of a concern.

For each peptide, it tests modifications serially, so it will not give a match to a peptide with multiple unsuspected modifications, such as might be found in a histone. In practice, the same limitation applies to the mass tolerant search; each modification takes out potential fragment matches, so having two or more makes getting a match very unlikely unless they are on adjacent residues.

I think the final limitation is the most serious. You can't use a two pass approach on endogenous peptides.

Maybe this is the most appropriate application for the mass tolerant search. If the data complexity is kept low, so that chimeric spectra are very rare, then the mass tolerant search may be an easier way to find modifications on endogenous peptides than an error tolerant sequence tag search