# De novo sequencing in Mascot Distiller

**Why carry out de novo?**

- **Unsequenced genomes**
  - No protein sequences available for database search
- **Unmatched spectra**
  - Unexpected PTMs
  - SNPs
  - Database errors
  - Missing sequences

MASCOT : De novo sequencing © 2013 Matrix Science    MATRIX SCIENCE

A typical peptide identification protocol works by carrying out a search of the uninterpreted MS/MS spectra against a protein (or translated Nucleic acid) sequence database. This can be a very efficient and effective technique, but for this approach to work the matching peptide sequence must be available in the sequence database searched.

De novo sequencing interprets peptide sequence information directly from the MS/MS spectrum, with no sequence database required. This can be a useful approach if you are working on an organism which does not have a sequenced genome. If the sequences from a related species is available, you can then use the interpreted de novo sequence results in a database homology search, such as an MS-BLAST search.

From a standard MS/MS database search, you'll often get a large number of unmatched spectra. Database searching of MS/MS spectra will often be caused to fail because of an unexpected post translational modification, or because of differences between the database sequence and the sequence in your dataset (often caused by sequence variations such Single Nucleotide Polymorphisms, sequencing errors in the database or sequences missing from the searched database). De novo sequencing is one approach that can help to match some of these spectra.

In general, we'd recommend carrying out a two pass error tolerant search for matching unexpected PTMs and SNPs before carrying out a de novo search as this is a much more

efficient way to pick up these types of matches.  If you are looking for SNPs or unexpected PTMs, only carry out a de novo search if you're unable to get matches from a standard, first pass search – if you're looking at endogenous peptides for example.

## De novo sequencing considerations

- **Leucine and Isoleucine have the same mass**
- **Near isobaric mass differences**
  - Glutamine and Lysine
  - Phenylalanine and oxidised methionine
- **Some amino acids have the same or nearly the same mass as pairs of others**
  - Mass of A+G=128.058578
  - Mass of Q = 128.058578

**MASCOT** : De novo sequencing © 2013 Matrix Science — MATRIX SCIENCE

When carrying out de novo sequencing there are some considerations which need to be taken which need to be taken into account and which can lead to errors and ambiguity in the generated sequence.

For example, Leucine and Isoleucine have the same mass and you cannot tell the difference between them in a standard MS/MS experiment.

Certain other amino acid pairs, such as glutamine and lysine, Phenylalanine and Oxidised methionine have near isobaric mass differences, and without a high accuracy instrument you're not going to be able to tell the difference between them.

Some amino acids have nearly the same mass as pairs of others. For example the mass of alanine plus glycine is the same as that of Glutamine. You may have enough information in the spectrum to confidently tell the difference between the options, but often you won't.

## De novo in Mascot Distiller

- Requires the optional 'Search toolbox'
- Novel algorithm
- Scores approximate to Mascot scores
- Single spectrum, unassigned list or entire dataset

MASCOT : De novo sequencing © 2013 Matrix Science    MATRIX SCIENCE

The Mascot Distiller search toolbox includes a powerful de novo sequencing module.

The score assigned to a de novo solution is similar to a Mascot score. In general, these scores will be higher than you expect to see in a Mascot database search, because the algorithm has selected the best matching sequence from all possible sequences, rather than the limited number of sequences found in any database. So, you should not judge the quality of the match by applying any rule of thumb or significance threshold to the score. However, if you get the same solution by de novo and by database search, using identical parameters, you should find the Mascot scores are very close.

You can de novo sequence a single spectrum, the 'unassigned list' from a Mascot search result or all the MS/MS spectra in the file

De novo settings
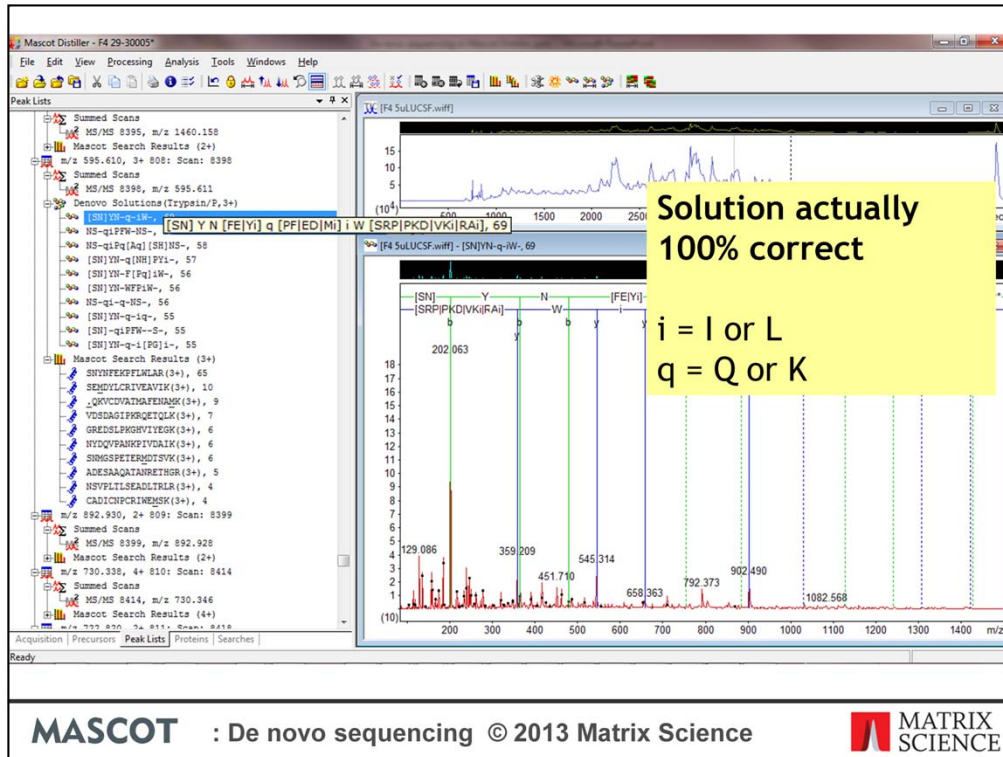
MASCOT : De novo sequencing © 2013 Matrix Science

As with a Mascot database search, there are various search settings which need to be set for a de novo search before carrying out the search. These can be found on the preferences dialog in Mascot Distiller, accessed from the 'Tools' menu.

The controls are very similar to those on the Mascot search settings form that you use when submitting a database search and it is important to set the correct de novo settings before initiating a de novo search, otherwise the default settings will be used which will probably not be appropriate for your dataset.

The starting point for de novo can be any MS/MS scan where Distiller has been used to create a peak list. Select the 'Peak list' tab in the dataset explorer in Distiller. Right click the peak list node and choose 'de novo Search', or choose the de novo button from the toolbar when a Summed Scans node is selected.

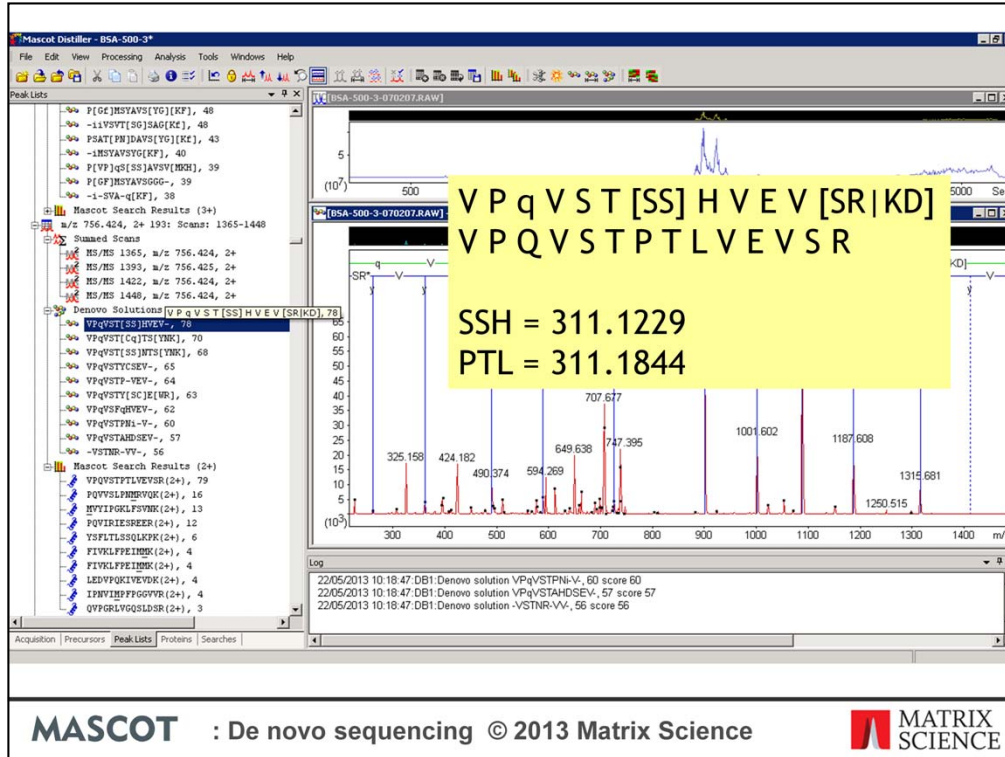When the de novo search has completed, the top ten solutions are added to the 'Peak lists' tab.

: De novo sequencing © 2013 Matrix Science

Let's take a look at a couple of de novo results where we already have strong database matches in order to illustrate some of the differences between a de novo result and a Mascot database search result. Of course you wouldn't normally carry out a de novo search of spectra where you have strong database matches as this would normally be a waste of time.

Good signal to noise and good mass accuracy are critical for successful de novo sequencing; much more so than in database searching. GIGO (garbage in - garbage out) is guaranteed.

In a de novo solution, i always represents I or L, and q represent Q or K when the mass tolerance does not allow these residues to be distinguished. However, K is assumed at the C terminus of a peptide when tryptic specificity applies. Also remember that phenylalanine may be oxidised methionine, and vice versa. Ambiguity is indicated by a dash in the sequence. The tooltip shows details of the ambiguity in square brackets, using pipe symbols to separate alternatives. Note that the order of any pairs and triplets is undefined, so that RAi could also be iAR.
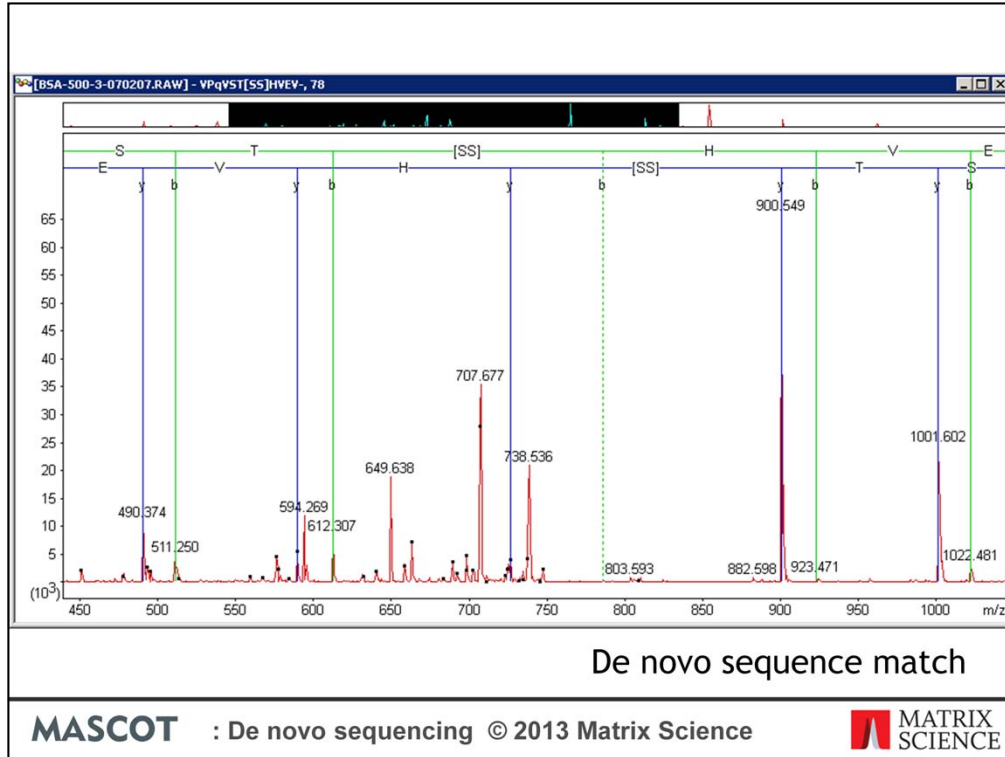
Although at first glance, the example shown here looks very different to the Mascot database match, they are actually in perfectly agreement. Some uncertainty is unavoidable in de novo, because the search space is so very much larger.
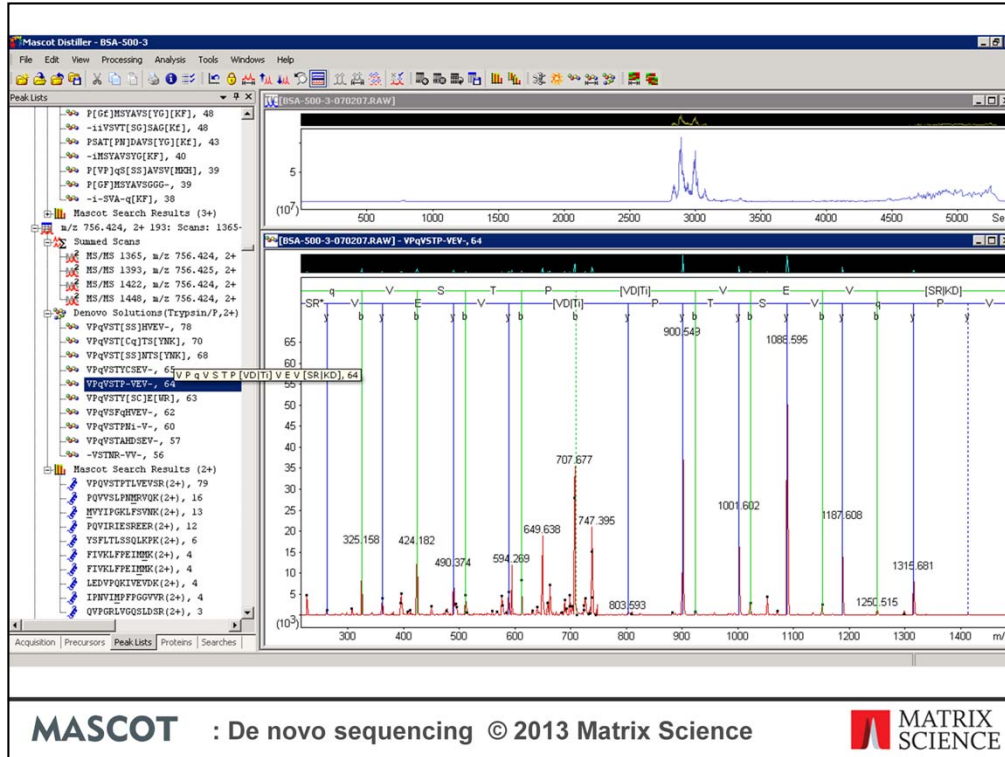
Lets take a look at another example. Here, we have a Mascot database match and a de novo match with practically the same score. The two sequences are almost identical – the difference being the de novo solution suggests SSH in the middle of the peptide where the database sequence is PTL, the masses of which are almost identical, and if you take a look at the matches in this region
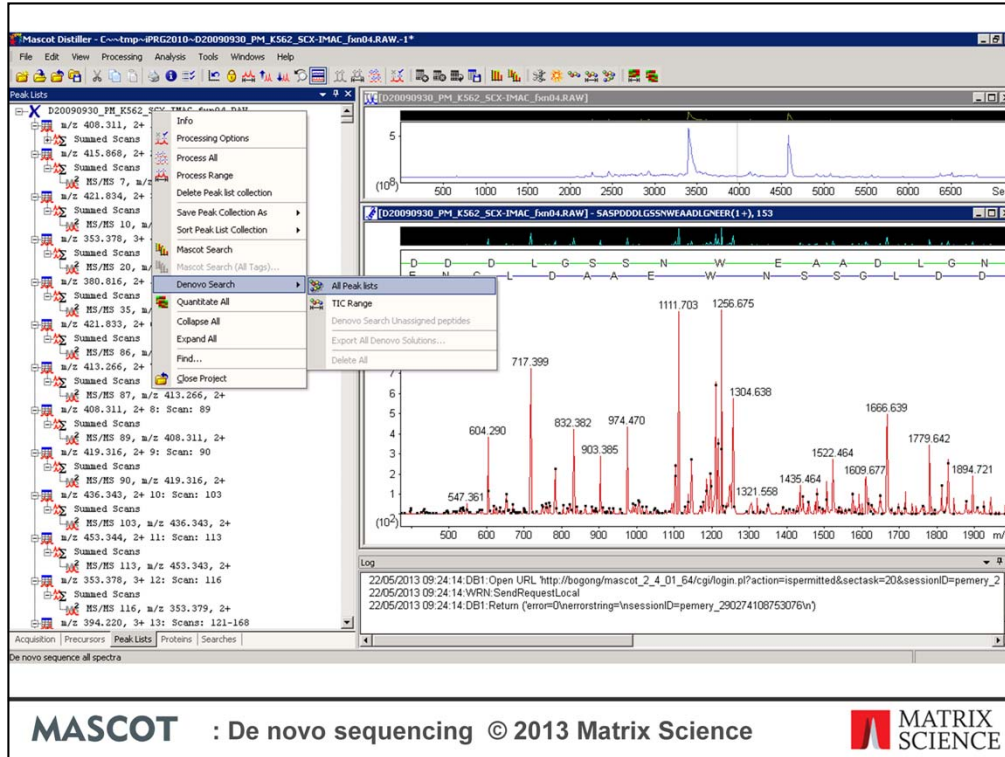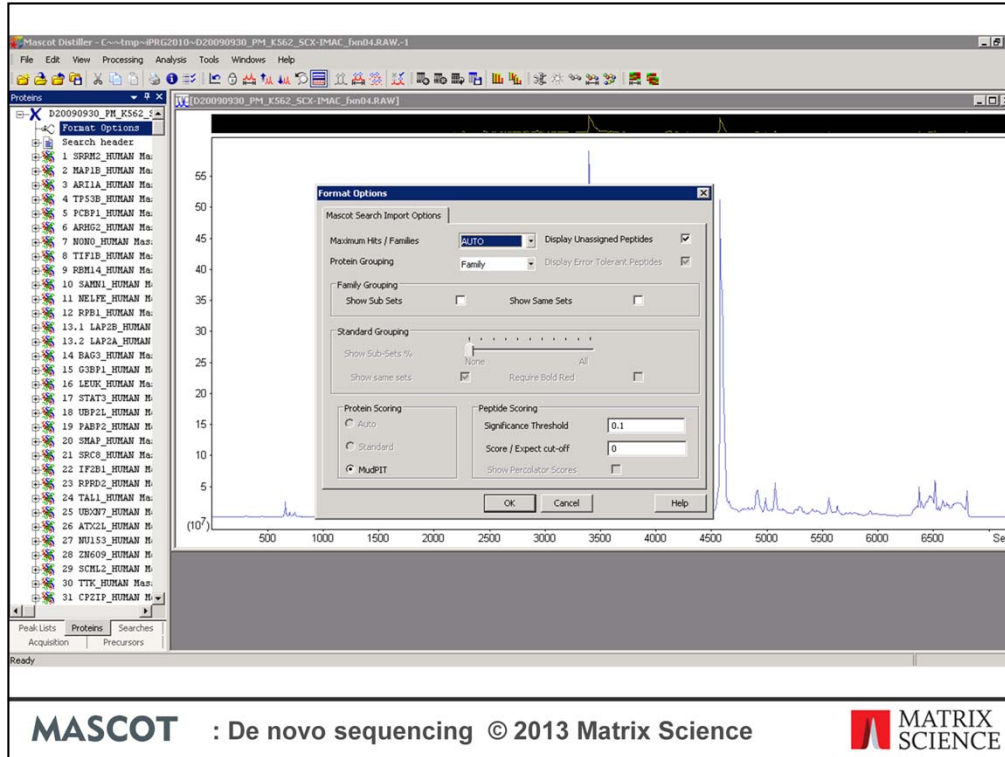
We can see that neither the Mascot database search match

De novo sequence match

MASCOT : De novo sequencing © 2013 Matrix Science

Nor the Distiller de novo search match has a complete set of fragment matches across this region.

MASCOT : De novo sequencing © 2013 Matrix Science

In fact the de novo search did find the same match as the Mascot database search, but it has a slightly lower score and has come out as the fifth ranked match. This peptide came from BSA, and this stretch of sequence is highly conserved across serum albumins from many different species, so it seems unlikely that the alternative de novo match is correct. In fact all ten reported de novo matches have reasonable scores and are, to a greater or lesser extent, a variation on the 'correct' (database) sequence. With such a large possible search space, this type of ambiguity is a fact of life with de novo. A sequence database search is much more constraining on the possible sequence matches.

To de novo sequence a complete peak list collection, or the peak lists in the currently displayed TIC range, use the context menu obtained by right-clicking the root (world) node of the peak lists tree

MASCOT : De novo sequencing © 2013 Matrix Science

One common reason for carrying out a de novo search is to try and find additional matches to spectra in the 'Unassigned' list in a Mascot search result – these are the spectra that failed to give decent matches. The most efficient way to de novo these spectra is to switch to the proteins tree, click on Format Options, and choose to load the unassigned queries.

Scroll down to the bottom of the proteins tree and use the context menu obtained by right-clicking the unassigned node to de novo just the unassigned queries. Once completed, the de novo solutions are added to the 'Peak list' tree in the dataset explorer. From there, we can scroll down the tree and look for spectra where we have a high scoring de novo solution, with no reliable Mascot database search match.

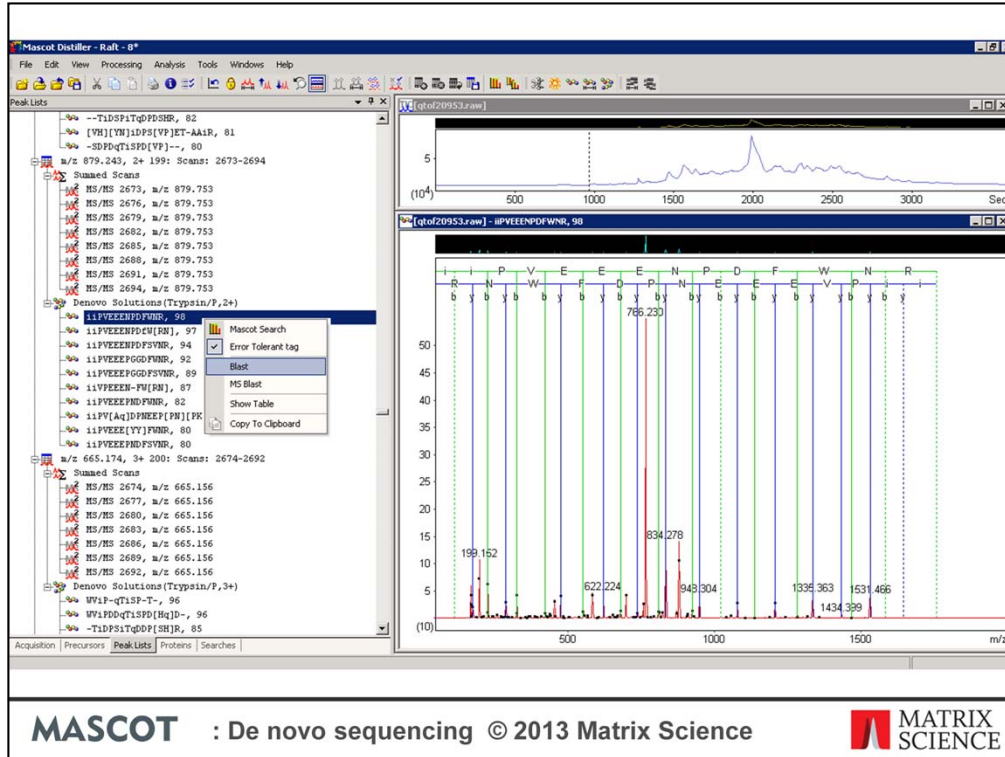MASCOT : De novo sequencing © 2013 Matrix Science

For a large dataset, an alternative is to export the de novo solutions as a CSV file and then viewing and filtering the results in Microsoft Excel.  To do this, select the 'Export All De novo Solutions' option from the 'File' menu, or from the context menu obtained by right clicking the root (world) node of the peak lists tree.  Once you've exported the results, you can load them directly into Excel and use the filtering functions to screen the de novo matches.

These data are from a QTof dataset.  Here, I've filtered on the 'Score' column in the CSV file for values greater than or equal to 80.

With high quality data, we can get something close to a complete peptide sequence, such as the match highlighted here.  To take a look at this match in Distiller, note the query m/z value of 879

Return to Distiller, select the 'Precursors' tab, and scroll down or use 'Find' function in Distiller to find the precursor. The tabs in the Distiller dataset explorer are, as far as is possible, kept synchronised – so now select the 'Peak Lists' tab

And here is our match. To see if we can get a match into a protein sequence, with a sequence this clean we could carry out a sequence homology search. To do this, right click over the de novo peptide match to bring up the context menu. This gives us various options to submit a Mascot tag search (we'll be taking a look at that in a moment), and to submit the peptide to a BLAST or MS-BLAST server. In this case we'll submit the peptide for an MS_BLAST search.

MASCOT : De novo sequencing © 2013 Matrix Science

If you submit a BLAST search, Distiller will open up submission windows – one with the forward peptide sequence and one with the reverse. Also, if there is ambiguity in the sequence only a single version of the peptide will be submitted for each direction. When submitting to MS Blast however, any ambiguity in the sequence will be expanded and formatted using MS Blast ambiguity codes, and all the generated the forward and reverse peptide sequences will be submitted together.

MASCOT : De novo sequencing © 2013 Matrix Science

When we get the results back, we find identity matches to several Alkaline phophatases, one of which PPB1_HUMAN was found in the original Mascot database search. In this case we have a nice match to a semi-tryptic peptide. An Error Tolerant Mascot database search would also have found this match.

MASCOT : De novo sequencing © 2013 Matrix Science

That is a fairly unusual case, and you could have just read the peptide sequence off. With typical data you'll often find that you only get partial sequence matches with rather more ambiguity in the sequence. Here, for example, we've got a match with where the top two de novo solutions have the same score, and both have a reasonable degree of ambiguity. The next step might be to take these partial sequences and search them as error tolerant sequence tag searches. This will allow the peptide mass to vary and allows the tags to float, potentially finding matches where there are sequence variations (from SNPs etc), none specific peptides or unsuspected modifications.

We can use the de novo solution to generate the tags automatically. To do this right click over the solution to bring up the context menu. Note that we have the 'Error tolerant tag' selected.

Distiller populates the query field with the tags taken from the non-ambiguous parts of the de novo solution.  We do this for each of the top two de novo solutions individually.  Then we submit the searches …

When the results come back, we get the same match back from the tags from the two different de novo solutions, but the match from the second de novo solution gets a better score.

**Peptide View**

MS/MS Fragmentation of **STTTGHLIYK**
Found in **EF1A_YEAST** in **SwissProt**, Elongation factor 1-alpha OS=Saccharomyces cerevisiae (strain ATCC 204508 / S288c)
GN=TEF1 PE=1 SV=1

Match to Query 1: 1133.582624 from(1134.589900,1+) index(0) etag(0.00000,STT,272.11369) etag(845.46652,TGH,550.34647)
etag(1047.53299,TTT,744.40614) etag(744.40614,GH[L|I],437.26275) etag(946.51621,TTG,687.39002)
Title: 126: Scan 3862 (rt=21.4386, p=0, c=965, e=1)

Monoisotopic mass of neutral peptide Mr(calc): 1119.5924
Unsuspected modification: 13.9902 Da, located in the region I8 to C-term
Ions Score: 106  Expect: 3.3e-09  (help)

| # | b | $b^0$ | Seq. | y | y* | $y^0$ | # |
|---|---|---|---|---|---|---|---|
| 1 | 88.0393 | 70.0287 | S | | | | 10 |
| 2 | 189.0870 | 171.0764 | T | 1033.5677 | 1016.5411 | 1015.5571 | 9 |
| 3 | 290.1347 | 272.1241 | T | 932.5200 | 915.4934 | 914.5094 | 8 |
| 4 | 391.1823 | 373.1718 | T | 831.4723 | 814.4458 | 813.4618 | 7 |
| 5 | 448.2038 | 430.1932 | G | 730.4246 | 713.3981 | | 6 |
| 6 | 585.2627 | 567.2522 | H | 673.4032 | 656.3766 | | 5 |
| 7 | 698.3468 | 680.3362 | L | 536.3443 | 519.3177 | | 4 |
| 8 | 811.4308 | 793.4203 | I | 423.2602 | 406.2336 | | 3 |
| 9 | 974.4942 | 956.4836 | Y | 310.1761 | 293.1496 | | 2 |
| 10 | | | K | 147.1128 | 130.0863 | | 1 |

**MASCOT** : De novo sequencing © 2013 Matrix Science

MATRIX SCIENCE

If we take a look at the peptide view for this match, it is suggesting an unsuspected modification of 13.9902 Da in the region I8 to the C-terminus.

Protein sequence coverage: 2%

Matched peptides shown in **bold red**.

```
  1 MGKEKSHINV VVIGHVDSGK STTTGHLIYK CGGIDKRTIE KFEKEAAELG
 51 KGSFKYAWVL DKLKAERERG ITIDIALWKF ETPKYQVTVI DAPGHRDFIK
101 NMITGTSQAD CAILIIAGGV GEFEAGISKD GQTREHALLA FTLGVRQLIV
151 AVNKMDSVKW DESRFQEIVK ETSNFIKKVG YNPKTVPFVP ISGWNGDNMI
201 EATTNAPWYK GWEKETKAGV VKGKTLLEAI DAIEQPSRPT DKPLRLPLQD
251 VYKIGGIGTV PVGRVETGVI KPGMVVTFAP AGVTTEVKSV EMHHEQLEQG
301 VPGDNVGFNV KNVSVKEIRR GNVCGDAKND PPKGCASFNA TVIVLNHPGQ
351 ISAGYSPVLD CHTAHIACRF DELLEKNDRR SGKKLEDHPK FLKSGDAALV
401 KFVPSKPMCV EAFSEYPPLG RFAVRDMRQT VAVGVIKSVD KTEKAAKVTK
451 AAQKAAKK
```

MASCOT : De novo sequencing © 2013 Matrix Science  MATRIX SCIENCE

Looking at the protein view, our match runs from residues 21 to 30. Fortunately, I searched SwissProt so the protein hit has annotation data.

```
DR    SUPFAM; SSF50465; Elong_init_C; 1.
DR    SUPFAM; SSF50447; Translat_factor; 1.
DR    TIGRFAMs; TIGR00483; EF-1_alpha; 1.
DR    PROSITE; PS00301; EFACTOR_GTP; 1.
PE    1: Evidence at protein level;
KW    3D-structure; Actin-binding; Complete proteome; Cytoplasm;
KW    Cytoskeleton; Direct protein sequencing; Elongation factor;
KW    GTP-binding; Methylation; Nucleotide-binding; Phosphoprotein;
KW    Protein biosynthesis; Reference proteome.
FT    CHAIN         1    458    Elongation factor 1-alpha.
FT                               /FTId=PRO_0000090973.
FT    NP_BIND      14     21    GTP.
FT    NP_BIND      91     95    GTP.
FT    NP_BIND     153    156    GTP.
FT    SITE        298    298    Not modified.
FT    SITE        372    372    Not modified.
FT    MOD_RES       1      1    Blocked amino end (Met).
FT    MOD_RES       6      6    Phosphoserine.
```

```
FT    MOD_RES      30     30    N6-methyllysine; by EFM1 (Probable).
```

```
FT    MOD_RES      53     53    Phosphoserine.
FT    MOD_RES      72     72    Phosphothreonine.
FT    MOD_RES      79     79    N6,N6,N6-trimethyllysine.
FT    MOD_RES     163    163    Phosphoserine.
FT    MOD_RES     259    259    Phosphothreonine.
FT    MOD_RES     316    316    N6,N6-dimethyllysine; by SEE1 (Probable).
FT    MOD_RES     355    355    Phosphotyrosine.
FT    MOD_RES     356    356    Phosphoserine.
FT    MOD_RES     390    390    N6-methyllysine; by EFM1 (Probable).
FT    MOD_RES     394    394    Phosphoserine.
FT    MOD_RES     414    414    Phosphoserine.
FT    MOD_RES     430    430    Phosphothreonine.
FT    MOD_RES     458    458    Lysine methyl ester.
FT    MUTAGEN     122    122    E->K: Reduces interaction with YEF3.
FT    MUTAGEN     153    153    N->D: Increases KM for GTP to 2.7 mM.
FT    MUTAGEN     153    153    N->T: Increases KM for GTP to 6.0 mM and
FT                               reduces translation fidelity. Increases
FT                               Km for GTP to 10.3 mM and reduces
FT                               translation fidelity; when associated
```

**MASCOT** : De novo sequencing © 2013 Matrix Science     MATRIX SCIENCE

Which suggests that Lysine 30 is Methylated.

A quick trip to Unimod shows that the mass shift associated with Methylation is +14.01Da, which is well within the mass tolerance for this search for the suggested mass shift of 13.9902 Da from the error tolerant tag search – so this seems like a reasonable assignment

## Conclusion

- **Poorly represented in any database**
  - De novo followed by sequence homology search
- **Well represented in the databases**
  1. Standard Mascot database search
  2. Error tolerant search
  3. De novo search of unassigned
     - Sequence homology search
     - Error tolerant tag search

**MASCOT** : De novo sequencing © 2013 Matrix Science — MATRIX SCIENCE

If you are trying to get as much coverage as possible over your data, you might come up with a search strategy like this:

If the organism you are working on is poorly represented and there isn't a closely related species in the database, then de novo followed by a BLAST or MS-BLAST search (if you get a reasonable length region of clean sequence) is probably your best option. However, if you have protein sequences available, then you're best off starting with a standard Mascot search, followed by an error tolerant search – this will get all of the 'easy' spectra. For the unassigned spectra, carry out a de novo search followed by BLAST or MS-BLAST. You can also try an etag search to find matches to isolated peptides that have a SNP or unsuspected modification.