

# Making sense of complex data sets with Mascot Insight

An in depth guide to reporting

**MASCOT** : Mascot Insight Reports

© 2014 Matrix Science

 **MATRIX  
SCIENCE**

This presentation follows on from the talk presented by Patrick Emery at the beginning of this workshop but I will be focusing on the data analysis and reporting functions available in Mascot Insight.

## Reports

- **Supports Mascot, mzIdentML, protXML based results**
- **Ships with over 30 reports covering**
  - Quantitation comparison, qc and clustering
  - Gene ontology
  - Interactions database analysis
  - Shared protein and peptide comparisons
  - Scatter plots, histograms etc
- **Can use filters with the reports**

As mentioned previously Mascot Insight Can support Mascot, mzIdentML, protXML based results so these reports are not just relevant to Mascot Server results.

Insight ships with over 30 reports designed specifically for proteomics data. These cover a wide range of areas such as result comparison, quantitation and quantitation clustering/grouping, Gene Ontology analysis, Interactions database analysis and general graphing reports such as scatter plots. It is possible to stack some of the filters and reports and use the results as inputs in to later analysis. For example filtering out the proteins identified in the contaminates database used in the initial search so that the reporting only works with the matches to the target database.

**Mascot Insight**  
User: Patrick Emery | 107 unread notifications | Preferences | Help | Logout

Searches: Xeno092-2.raw, Xeno092-3.raw, Xeno441-1.raw, Xeno441-2.raw, Xeno441-3.raw

Root  
 Quantitation datasets  
 Dundee iTRAQ dataset  
 Gent dilution series  
 Stress-induced stimulation in A.T.  
 Super-SILAC analysis of human l...  
 7 Lung cancer Xeno092  
 Xeno092-1.raw  
 Xeno092-2.raw  
 Xeno092-3.raw  
 9 Lung cancer Xeno441  
 Xeno441-1.raw  
 Xeno441-2.raw  
 Xeno441-3.raw  
 10 Lung cancer Xeno...  
 11 Lung cancer Xeno...  
 Human oral cancer brus...  
 YABBY  
 THINKPAD-TS20  
 T500-PAE  
 T440-PAE  
 paemery  
 Patrick Emery  
 Unassociated Mascot search...  
 Trash

Select  
 Expand  
 Collapse  
 Clear search  
 Add notes  
 Edit notes  
 Adjust  
 Import  
 Import  
 Open search  
 View search  
 Merge  
 Run report  
 Cut notes  
 Copy notes  
 Clone notes  
 Delete

| Hit rank | Accession   | Protein                                      |
|----------|-------------|--|
| 1.1      | MIM3_HUMAN  | Myosin-9 OS=Homo sapiens GI=MIM3 PE=1        |
| 1.2      | MIM10_HUMAN | Myosin-10 OS=Homo sapiens GI=MIM10 PE=1      |
| 1.3      | MIM11_HUMAN | Myosin-11 OS=Homo sapiens GI=MIM11 PE=1      |
| 1.4      | MIM14_HUMAN | Myosin-14 OS=Homo sapiens GI=MIM14 PE=1      |
| 1.5      | MIM1_HUMAN  | Myosin-1 OS=Homo sapiens GI=MIM1 PE=1        |
| 1.6      | MIM2_HUMAN  | Myosin-2 OS=Homo sapiens GI=MIM2 PE=1        |
| 1.7      | MIM4_HUMAN  | Myosin-4 OS=Homo sapiens GI=MIM4 PE=1        |
| 1.8      | TBM44_HUMAN | Mitochondrial import inner membrane transloc |
| 1.9      | MIM3_HUMAN  | Myosin-3 OS=Homo sapiens GI=MIM3 PE=1        |
| 2.1      | ENOC_HUMAN  | Alpha-enolase OS=Homo sapiens GI=ENOC1       |
| 2.2      | ENOC_HUMAN  | Beta-enolase OS=Homo sapiens GI=ENOC2        |
| 2.3      | ENOC_HUMAN  | Gamma-enolase OS=Homo sapiens GI=ENOC3       |
| 3.1      | FLNA_HUMAN  | Filamin-A OS=Homo sapiens GI=FLNA PE=1       |

Show selected rows only? | 3,171 / 3,171 proteins selected |

| Query No | %I-term rest | Peptide  | C-term rest | Query Title     | Charge | Observed | Intensity |
|----------|--------------|----------|-------------|-----------------|--------|----------|-----------|
| 134      | K            | QLAADR   | L           | 1733: Scan ...  | 2      | 401.21   | 1.1067    |
| 673      | R            | FVQGLR   | E           | 8230: Scan ...  | 2      | 404.79   | 859020.90 |
| 684      | K            | QLAADR   | L           | 1747: Scan ...  | 2      | 406.22   | 5.7785    |
| 1042     | K            | TLEERAK  | F           | 3265: Scan ...  | 2      | 410.31   | 7.4686    |
| 1274     | K            | TLEERAK  | F           | 2739: Sum ...   | 2      | 414.22   | 1.4527    |
| 1393     | K            | SALARRK  | R           | 11099: Scan ... | 2      | 415.73   | 5.5326    |
| 1663     | R            | NCAATLK  | L           | 3452: Scan ...  | 2      | 420.31   | 8.1826    |
| 1736     | K            | SYMELEK  | S           | 1504: Scan ...  | 2      | 421.22   | 9.4085    |
| 1925     | K            | DQGLER   | Q           | 3675: Scan ...  | 2      | 423.70   | 3.6056    |
| 2049     | K            | SYMELEK  | S           | 1500: Sum ...   | 2      | 425.23   | 6.8935    |
| 3052     | R            | SQAMARK  | K           | 22087: Sum ...  | 2      | 425.23   | 3.2485    |
| 2781     | -            | MAQQAAQK | V           | 16128: Sum ...  | 2      | 435.72   | 275709.77 |
| 2847     | D            | hencosr  | V           | 6111: Scan ...  | 2      | 436.26   | 4.2162    |

Show selected rows only? | 237/274 peptides selected | Selection: Significant peptides

| Quantitation | Protein View | Taxonomy | Family dendrogram | Family | Annotation | Spectrum | Ions matched | Error distribution |
|--------------|--------------|----------|-------------------|--------|------------|----------|--------------|--------------------|
| 2            | QLAARF       |          |                   |        |            |          | 0.12         | 0.22               |
| 1            | QLAAR        |          |                   |        |            |          | 0.11         | 0.40               |
| 2            | GALALE       |          |                   |        |            |          | 0.14         | 0.25               |
| 1            | GALAL        |          |                   |        |            |          | 0.17         | 0.65               |
| 2            | NCAATK       |          |                   |        |            |          | 0.18         | 0.23               |
| 1            | NCAATLK      |          |                   |        |            |          | 0.13         | 0.63               |
| 2            | SYMELEK      |          |                   |        |            |          | 0.09         | 0.60               |
| 1            | SYMELEK      |          |                   |        |            |          | 0.26         | 0.91               |
| 2            | DQGLER       |          |                   |        |            |          | 0.15         | 0.58               |
| 1            | DQGLER       |          |                   |        |            |          | 0.83         | 0.42               |
| 2            | ADEVLAK      |          |                   |        |            |          | 0.53         | 0.46               |
| 1            | ADEVLAK      |          |                   |        |            |          | 0.88         | 0.13               |
| 2            | ASTALEAK     |          |                   |        |            |          | 0.56         | 0.57               |
| 1            | ASTALEAK     |          |                   |        |            |          | 0.81         | 0.83               |

MASCOT : Mascot Insight Reports © 2014 Matrix Science MATRIX SCIENCE

You can initiate reports from the data tree or within MIRA at the Search, Protein, Peptide or quantitation level.

# Report selection

|  |        |                                     |                |  |
|--|--------|-------------------------------------|----------------|--|
| Scatter plot                               | Report | One or more Search results          | Matrix Science | Plot Scattergraph<br><a href="#">Detailed description</a>  |
| Gene ontology barchart                     | Report | One or more Search results          | Matrix Science | Gene Ontology barchart<br><a href="#">Detailed description</a>   |
| De novo sequence homology report           | Report | Two Search results                  | Matrix Science | Protein sequence Homology search for de novo sequencing results<br><a href="#">Detailed description</a>                              |
| Gene ontology proteome comparison          | Report | One or more Search results          | Matrix Science | GO proteome comparison report<br><a href="#">Detailed description</a>  |
| Unique peptides vs all peptides cross plot | Report | One or more Search results <b>Q</b> | Matrix Science | Protein ratios based on unique peptides plotted against Protein ratios based on all peptides<br><a href="#">Detailed description</a> |

**Protein**

| ID   | Type   | Available for                 | Author         | Description   |
|--|--------|-------------------------------|----------------|---|
| Protein interaction network                | Report | One Protein                   | Matrix Science | Protein interaction network report<br><a href="#">Detailed description</a>  |
| Protein interaction shortest path          | Report | Two Proteins                  | Matrix Science | Calculate shortest interaction path between two proteins<br><a href="#">Detailed description</a>  |
| K-means clustering                         | Report | One or more Proteins <b>Q</b> | Matrix Science | K-means clustering report quantitation results<br><a href="#">Detailed description</a>  |
| Histogram                                  | Report | One or more Proteins          | Matrix Science | Plot Histogram<br><a href="#">Detailed description</a>  |
| Scatter plot                               | Report | One or more Proteins          | Matrix Science | Plot Scattergraph<br><a href="#">Detailed description</a>   |
| Gene ontology barchart                     | Report | One or more Proteins          | Matrix Science | Gene Ontology barchart<br><a href="#">Detailed description</a>  |
| SPlot                                      | Report | One or more Proteins          | Matrix Science | Plot S-plot<br><a href="#">Detailed description</a>   |
| Unique peptides vs all peptides cross plot | Report | One or more Proteins <b>Q</b> | Matrix Science | Protein ratios based on unique peptides plotted against Protein ratios based on all peptides<br><a href="#">Detailed description</a>  |
| Protein coefficient of variation histogram | Report | One or more Proteins <b>Q</b> | Matrix Science | This report generates histograms showing the coefficient of variation (CV) for the protein quantitation ratios for the selected protein hits (in log space)<br><a href="#">Detailed description</a> |

**Peptide**

| ID                                      | Type   | Available for                 | Author         | Description   |
|---|--------|-------------------------------|----------------|---|
| Peptide quantitation ratio distribution | Report | One or more Peptides <b>Q</b> | Matrix Science | Quantitation peptide ratio distribution for selected peptides<br><a href="#">Detailed description</a> |
| K-means clustering                      | Report | One or more Peptides <b>Q</b> | Matrix Science | K-means clustering report quantitation results<br><a href="#">Detailed description</a>                |
| Histogram                               | Report | One or more Peptides          | Matrix Science | Plot Histogram<br><a href="#">Detailed description</a>  |
| Scatter plot                            | Report | One or more Peptides          | Matrix Science | Plot Scattergraph<br><a href="#">Detailed description</a>   |
| Normal probability plot                 | Report | One or more Peptides <b>Q</b> | Matrix Science | Quantitation Normal Probability Plot<br><a href="#">Detailed description</a>                          |
| SPlot                                   | Report | One or more Peptides          | Matrix Science | Plot S-plot<br><a href="#">Detailed description</a>   |

**MASCOT** : Mascot Insight Reports

© 2014 Matrix Science

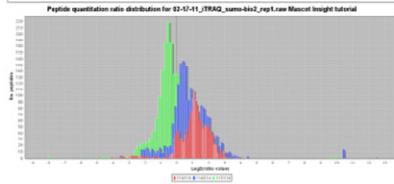


The report selection window appears. Only the reports that are suitable for the data selected are active in the list, shown in blue.

Quantitation reports are marked with a Q symbol.

# Report selection

| Id                                 | Quantitation distribution histogram  |  |
|------------------------------------|--|--|
| <a href="#">Statistics export</a>  |  | delineated format for use in statistical software (e.g., R, Perseus) |
| <a href="#">JTreeView export</a>   | <b>Available for</b><br>Single quantitation result   | compatible pdl file format   |
| <a href="#">Standard report</a>    | <b>Report description</b><br>This report plots the protein (anchor proteins) and peptide (rank 1 matches) quantitation ratio value distributions for the selected search, enabling you to view the overall quantitation ratio distribution. Since quantitation ratio distributions are expected to be log normal, the quantitation ratio values are always log2 transformed prior to being added to the histogram. | report for selected search in Excel format                           |
| <a href="#">Publication report</a> |  | based on the MCP guidelines for publication                          |



**MASCOT** : Mascot Insight Reports

© 2014 Matrix Science



There is a one line description for each report and a more detailed description a click away.

## Super-SILAC analysis of human lung primary tumor-derived xenografts LC-MS/MS

- **Data set available from PRIDE PXD000438**
  - Zhang W, et al, Proteomics, 14(6):795-803(2014)
- **Lys +8, Arg +10 labeling of standard vs unlabeled tumor**
- **2 major histological subtypes analysed:**
  - Adenocarcinoma (ADC)
  - Squamous cell carcinoma (SCC)
- **2 biological samples for each subtype, 3 technical replicates**
- **Aim was to use protein expression to find candidate proteins that differentiate the two tumor types**

Proteomic profiles of human lung adeno and squamous cell carcinoma using super-SILAC and label-free quantification approaches. Zhang W., et al, Proteomics. 2014 Mar;14(6):795-803.

**MASCOT** : Mascot Insight Reports

© 2014 Matrix Science

**MATRIX**  
SCIENCE

The data set that I am going to use for most of the examples is available at the EBI PRIDE repository.

It is a Super-SILAC experiment using a lung cell culture labelled with Lys +8 and Arg +10 as a standard and comparing it to unlabeled human lung primary tumor-derived xenografts and analysed by LCMS/MS.

2 major histological subtypes were analysed, ADC and SCC with 2 biological samples for each subtype and 3 technical replicates

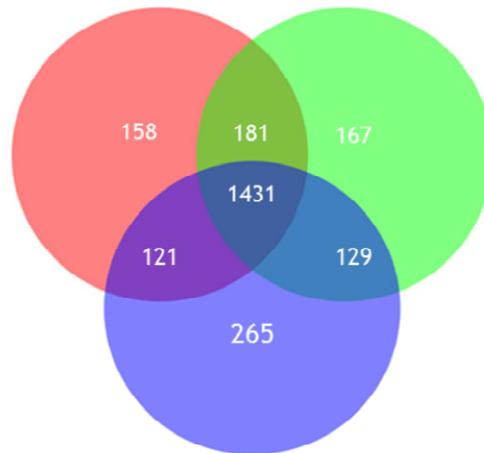
for each sample for a total of 12 MS/MS datasets. The data was analysed using the same search settings described in the publication.

The goal of the experiment was to compare the two lung tumor types and to use protein expression to differentiate them.

Reference: Proteomics. 2014 Mar;14(6):795-803. doi: 10.1002/pmic.201300382, Proteomic profiles of human lung adeno and squamous cell carcinoma using super-SILAC and label-free quantification approaches, Zhang W1, Wei Y, Ignatchenko V, Li L, Sakashita S, Pham NA, Taylor P, Tsao MS, Kislinger T, Moran MF.

## Venn diagram of identified proteins for three technical replicates of sample Xeno 441/ADC

■ Run 20140221\_000011\_000001\_1\_000  
■ Run 20140221\_000011\_000002\_1\_000  
■ Run 20140221\_000011\_000003\_1\_000  
Only protein/peptide protein hits counted



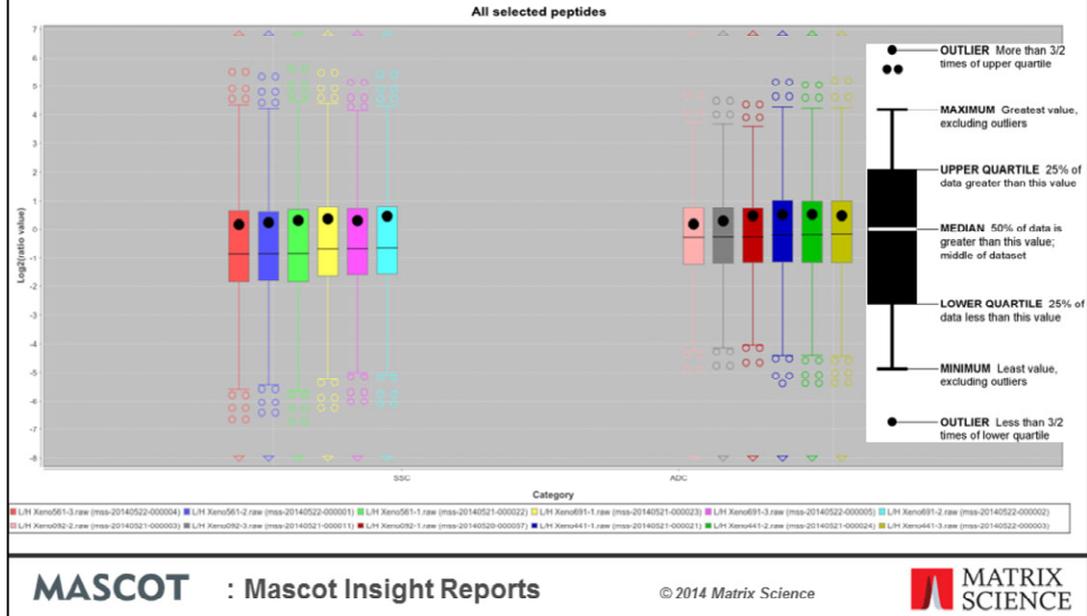
**MASCOT** : Mascot Insight Reports

© 2014 Matrix Science

**MATRIX**  
SCIENCE

First I took a look at the technical replicates for each biological sample to see how consistent they are. These are the three replicates for the ADC sample Xeno 441. As you can see there is good overlap between the different analysis and each individual run identifies no more than 15% unique proteins.

## Box and whisker plot



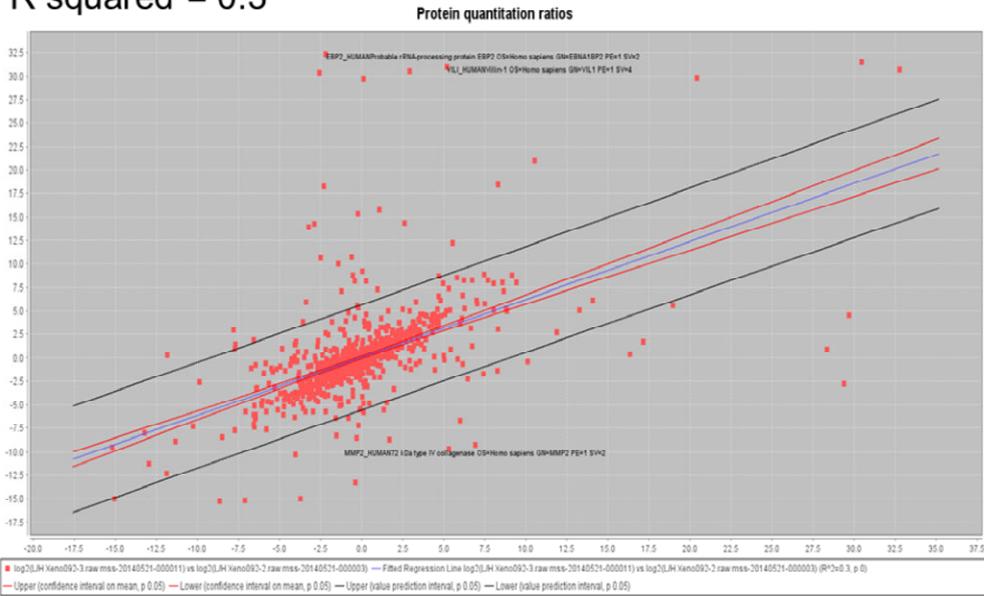
Next I generated box and whisker plot. For those not familiar with box and whisker plots here is a quick guide. The plot will give us global information for each ratio from each replicate, showing us things like the global mean and median ratios, the data spread, presence of outliers etc. The black circle is the average.

You can see the overall distributions of peptide ratios for each technical replicate, and also each tumour type, are pretty close to each other.

Once a report is generated, you can easily copy and paste it out of MIRA, or export it in a number of formats – including exporting the underlying data for the dataset as a CSV file, in case you want to get the data into another package such as Excel or R.

Min no peptides = 1

R squared = 0.3



**MASCOT** : Mascot Insight Reports

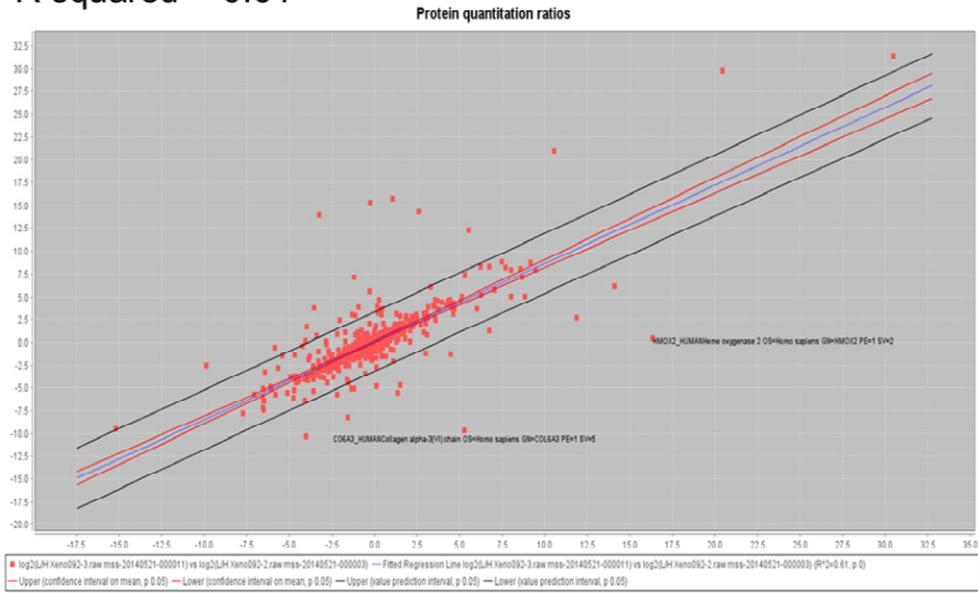
© 2014 Matrix Science



We can explore the quality of the data set further looking at the effect of the minimum number of quantitation values used per a protein on the ratio correlation. The question is how many peptide ratio measurements do I need for reliable quantitation results. For these examples I used two of Xeno 092 ADC replicates. Prior to creating the report I filtered out matches to the contaminants database.

I then plotted the ratios from two replicates against each other using different numbers of minimum required peptide ratios per a protein (1-4), getting better and better correlations. Here is the plot for using 1 peptide ratio and it does not give me much confidence.

Min no peptides = 2  
R squared = 0.61



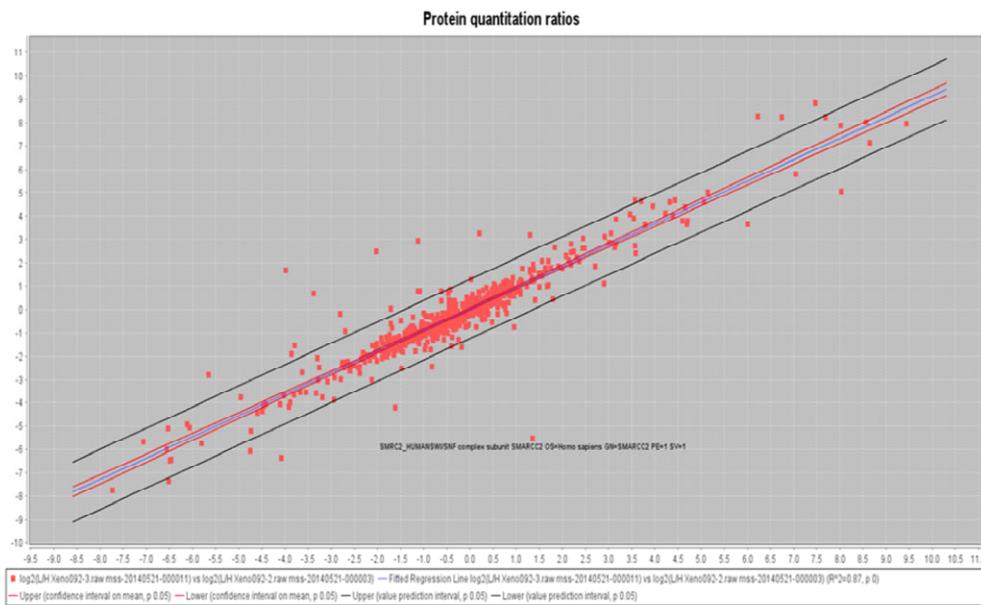
**MASCOT** : Mascot Insight Reports

© 2014 Matrix Science



Using 2 peptides the protein ratio correlation starts to improve

Min no peptides = 3  
R squared = 0.87



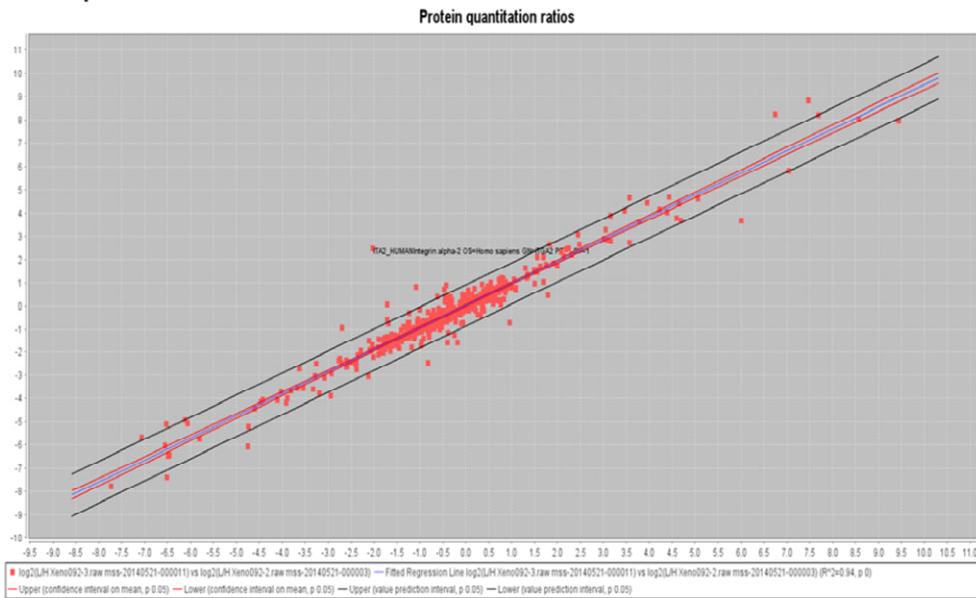
**MASCOT** : Mascot Insight Reports

© 2014 Matrix Science



With three measurements - this is more like

Min no peptides = 4  
R squared = 0.94

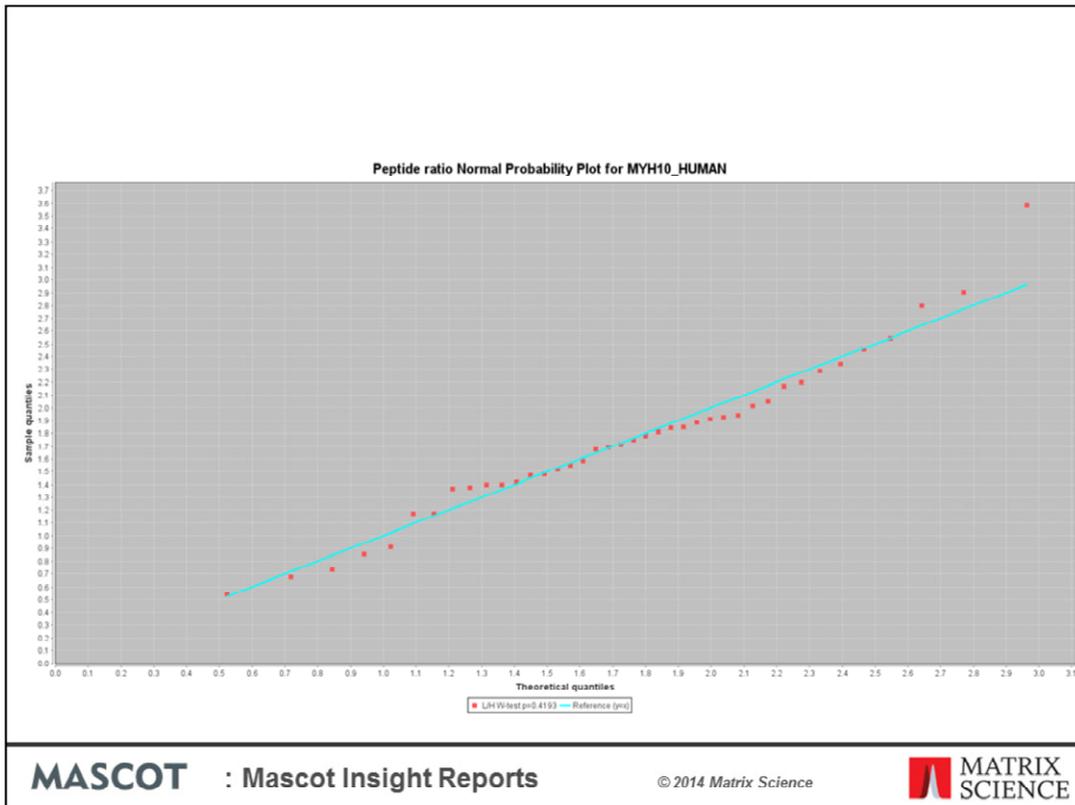


**MASCOT** : Mascot Insight Reports

© 2014 Matrix Science



Now we are there. This final plot shows that using 4 peptides gives a much higher ratio correlation between technical replicates than using just one measurement per a protein. This should be common sense but we still have queries from customers about obtaining as many protein quantitation ratios as possible even if they one have one or two peptide measurements.



We would certainly be expecting Keratins to be differently expressed in a lung tumor. Here I show a Peptide ratio Normal probability plot for the Keratin, type I C 18 peptides. The peptide quantitation data deviated from the straight line of a theoretical normal distribution. This may be due to the large number of shared peptides between different keratin isoforms. Here is a plot for MYH10\_HUMAN which is a lot closer to a straight line and of the quality we are looking for.

## Blast cluster report

- **BLASTClust is part of the BLAST package of programs**
- **Clusters either protein or nucleotide sequences**
- **Clusters protein hits of the whole protein sequence rather than identified peptides**
- **Alternative to Mascot Server protein inference algorithm**

BLASTClust is a program within the standalone BLAST package used to cluster either protein or nucleotide sequences. The program begins with pairwise matches and places a sequence in a cluster if the sequence matches at least one sequence already in the cluster. In the case of proteins, the blastp algorithm is used to compute the pairwise matches.

BLASTClust clusters protein hits on the whole protein sequence rather than just the identified peptides. This means that isoforms and protein families will be clustered together while proteins that only share one peptide sequence are unlikely to be grouped together. This makes it an alternative protein grouping algorithm to those used by Mascot Server.

## Blast cluster results

| Accession           | Protein Name                           | Length     | Score        | Identity        | Diagram    |
|---------------------|--|------------|--------------|-----------------|------------|
| <b>K1C19_HUMAN</b>  | <b>Keratin, type I cytoskeletal 19</b> | <b>440</b> | <b>21336</b> | <b>440/79.0</b> | <b>638</b> |
| msa-20140521-000024 | OS=Homo sapiens GN=KRT19 PE=1 SV=4     | 440        | 21336        | 440/79.0        | 638        |
| K1C19_HUMAN         | msa-20140521-000021                    | 20890      | 44079.0      | 591             | 400        |
| K1C19_HUMAN         | msa-20140522-000003                    | 20753      | 44079.0      | 579             | 400        |
| K1C19_HUMAN         | msa-20140521-000003                    | 7374       | 44079.0      | 201             | 400        |
| K1C19_HUMAN         | msa-20140521-000011                    | 6713       | 44079.0      | 204             | 400        |
| K1C19_HUMAN         | msa-20140520-000057                    | 6119       | 44079.0      | 198             | 400        |
| <b>K1C18_HUMAN</b>  | <b>Keratin, type I cytoskeletal 18</b> | <b>422</b> | <b>16670</b> | <b>48029.0</b>  | <b>422</b> |
| msa-20140522-000003 | OS=Homo sapiens GN=KRT18 PE=1 SV=2     | 422        | 16670        | 48029.0         | 422        |
| K1C18_HUMAN         | msa-20140521-000021                    | 16578      | 48029.0      | 440             | 400        |
| K1C18_HUMAN         | msa-20140521-000024                    | 15659      | 48029.0      | 466             | 400        |
| K1C18_HUMAN         | msa-20140520-000057                    | 6481       | 48029.0      | 184             | 400        |
| K1C18_HUMAN         | msa-20140521-000003                    | 5660       | 48029.0      | 138             | 400        |
| K1C18_HUMAN         | msa-20140521-000011                    | 4854       | 48029.0      | 114             | 400        |
| <b>LMNA_HUMAN</b>   | <b>Profilin A/C OS=Homo sapiens</b>    | <b>384</b> | <b>6907</b>  | <b>74380.0</b>  | <b>384</b> |
| msa-20140522-000003 | GN=LMNA PE=1 SV=1                      | 384        | 6907         | 74380.0         | 384        |
| LMNA_HUMAN          | msa-20140521-000021                    | 3813       | 74380.0      | 269             | 404        |
| LMNA_HUMAN          | msa-20140521-000024                    | 5730       | 74380.0      | 271             | 404        |
| LMNA_HUMAN          | msa-20140521-000003                    | 5527       | 74380.0      | 158             | 404        |
| LMNA_HUMAN          | msa-20140521-000011                    | 4966       | 74380.0      | 154             | 404        |
| LMNA_HUMAN          | msa-20140520-000057                    | 3678       | 74380.0      | 147             | 404        |
| LMNA_HUMAN          | msa-20140522-000003                    | 6907       | 74380.0      | 284             | 404        |
| LMNA_HUMAN          | msa-20140521-000021                    | 5813       | 74380.0      | 269             | 404        |
| LMNA_HUMAN          | msa-20140521-000024                    | 5730       | 74380.0      | 271             | 404        |
| <b>K2CB_HUMAN</b>   | <b>Keratin, type II cytoskeletal 8</b> | <b>266</b> | <b>5647</b>  | <b>53671.0</b>  | <b>266</b> |
| msa-20140521-000003 | OS=Homo sapiens GN=KRT8 PE=1 SV=2      | 266        | 5647         | 53671.0         | 266        |
| K2CB_HUMAN          | msa-20140521-000024                    | 5114       | 53671.0      | 280             | 403        |
| K2CB_HUMAN          | msa-20140522-000003                    | 5080       | 53671.0      | 253             | 403        |
| K2CB_HUMAN          | msa-20140521-000003                    | 4492       | 53671.0      | 150             | 403        |
| K2CB_HUMAN          | msa-20140521-000011                    | 4247       | 53671.0      | 155             | 403        |
| K2CB_HUMAN          | msa-20140520-000057                    | 2679       | 53671.0      | 134             | 403        |
| K2CB_HUMAN          | msa-20140522-000003                    | 5080       | 53671.0      | 253             | 403        |

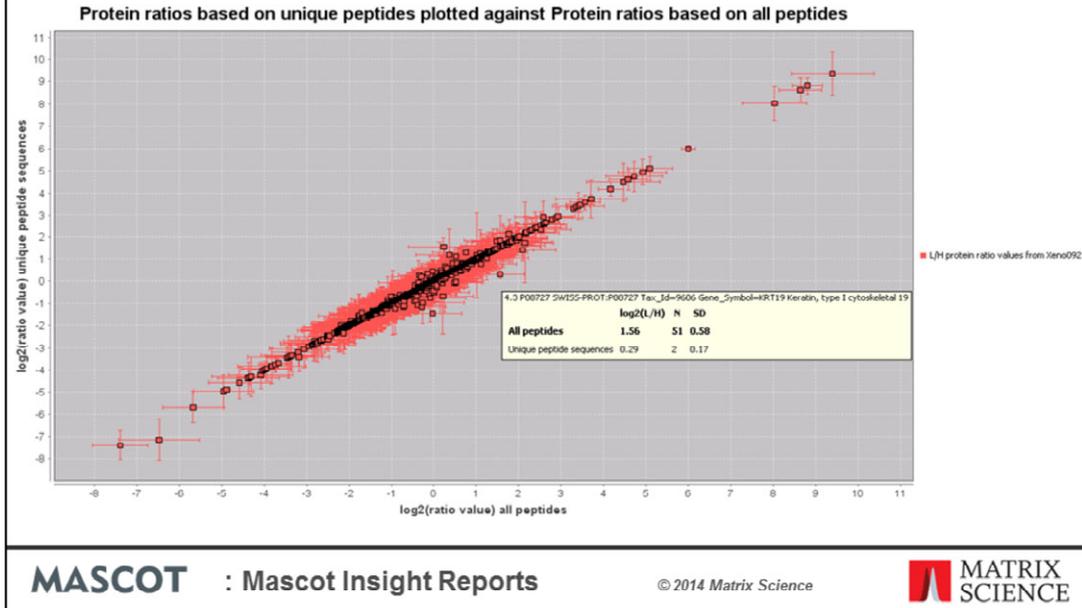
**MASCOT** : Mascot Insight Reports

© 2014 Matrix Science

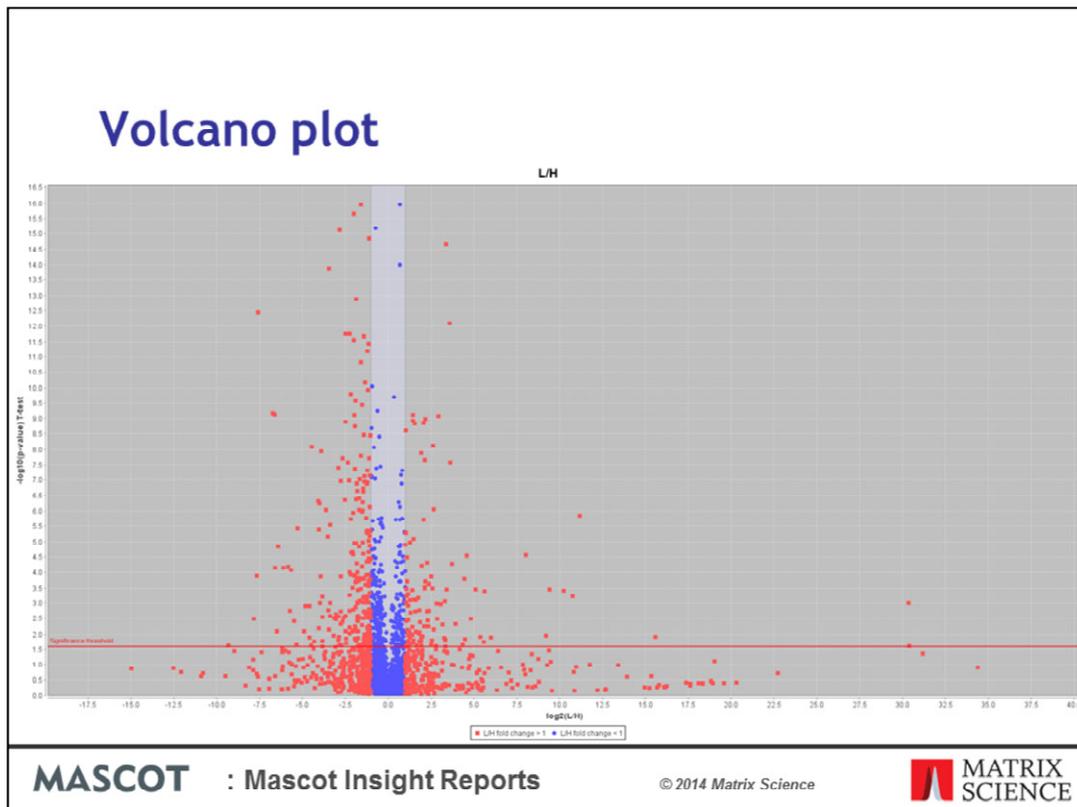


The peptide sharing between Keratins become obvious once you use the BLAST cluster algorithm to group them. There were 77 keratin sequences in the first group.

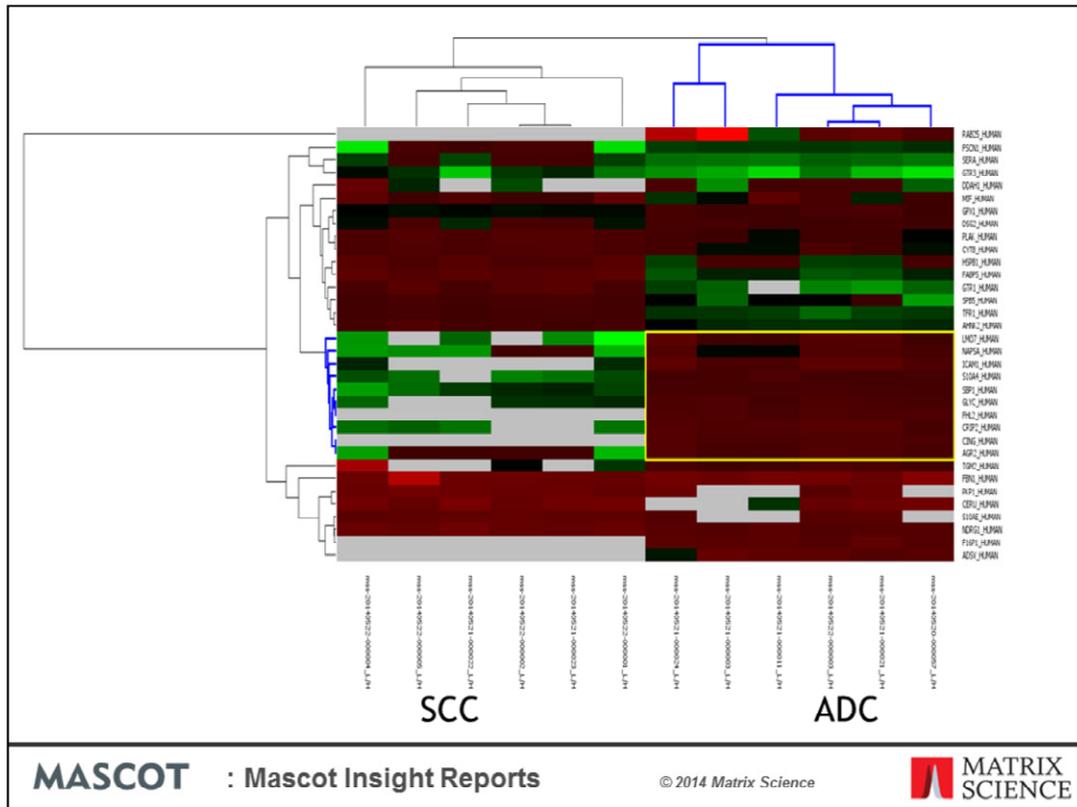
## Protein ratios based on unique peptides plotted against Protein ratios based on all peptides



One final quality plot that illustrates the problem with shared peptides. This is a plot of the Log<sub>2</sub> transformed Protein ratios based on unique peptides plotted against Log<sub>2</sub> transformed Protein ratios based on all peptides. Proteins with no shared peptides or no difference fall on the diagonal line. However some of the proteins are quite a way of the line with a wide standard deviation, as shown by the error bars. Mousing over the data point or clicking on them give the accession number and information on the protein in question. In this case Keratin type 1 C19 whose unique peptides barely changed with a non significant Log<sub>2</sub> transformed ratio of 0.29 (close to 1.2 ), compared to its shared peptides with a significant up regulation of Log<sub>2</sub> transformed ratio of 1.56 (close to 2.9).



Now that I am satisfied that the overall quality of the data set is good I want to know which are the differently regulated proteins. This is a Volcano plot, generated for one of the ADC fractions, Xeno092-1. A Volcano plot combines the t-statistic for the protein quantitation ratio with the protein quantitation ratio fold change. Using the Volcano plot we can identify the proteins whose protein quantitation ratio value are above the significance threshold and their Protein ratio values after log<sub>2</sub> transformation are greater than 1 or -1. These proteins are the significantly up or down regulated proteins of interest. The original paper used protein ratio values after log<sub>2</sub> transformation greater than 2 or -2.

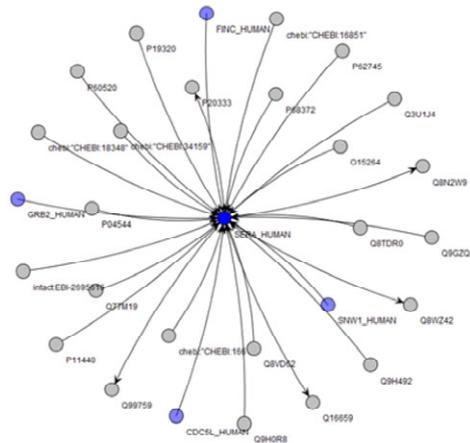


The top 40 up and down regulated proteins as identified from a volcano plot, minus the 4 Keratin proteins that were shown to have not normal peptide probability plots, were used for the Hierarchical clustering report. The Hierarchical clustering report uses unsupervised hierarchical clustering to group both the proteins and the quantitation ratios into clusters and then generates a heatmap of the results. Here we are looking at the heatmap generated from. Red is up and green is down regulated – the hierarchical clustering algorithm has clustered the matching ratios from each sample, and we can see from the heatmap that we can identify candidate proteins to differentiate the two types of tumour based on up and down regulated proteins.

In the original publication the authors go on to use the candidate signature proteins to differentiate an additional 12 tumour samples that were quantified by label-free spectral counting.

Other quantitation specific plots include a Bland-Altman plot, a QQ plot for peptide ratios matching a protein hit

## Phosphoglycerate dehydrogenase (SERA\_HUMAN) interaction map

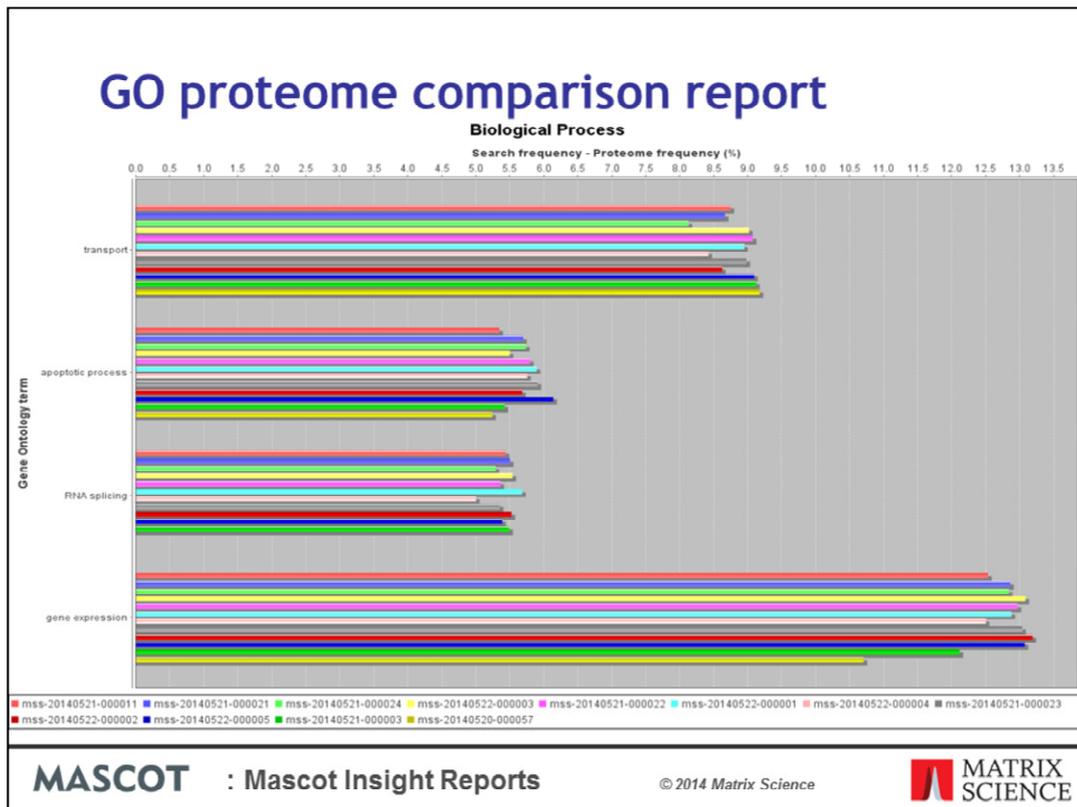


**MASCOT** : Mascot Insight Reports

© 2014 Matrix Science

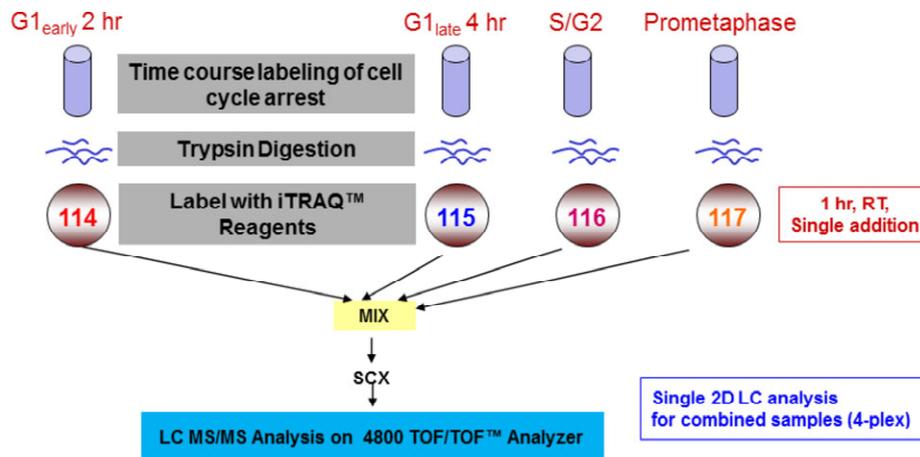


We can explore the candidate “signature” proteins identified in the clustering by running an interaction report on the accession numbers. Here we look at Phosphoglycerate dehydrogenase (SERA\_HUMAN) which is known as a key enzyme involved in the re-routing of glycolytic carbon into serine biosynthesis. It was found more highly expressed in SCC, and is known to be upregulated in 70% of estrogen receptor negative breast cancers, and is required for tumorigenesis in in vivo breast cancer models. Proteins in blue have been identified in the data set.



As a final analysis I have plotted a Gene Ontology proteome comparison report. Proteome frequencies are calculated using UniProt species proteomes. The plot is search frequency - proteome frequency - so over represented GO terms are positive values, under represented would be negative. The report here has been filtered to show on differences  $\geq 5\%$ ; so nothing highly under represented but we have a large over representation of proteins involved in gene expression, RNA splicing, apoptosis etc. not unexpected for a rapidly dividing cell line.

## Cell Cycle Arrest and Labeling Workflow (HeLa S3)



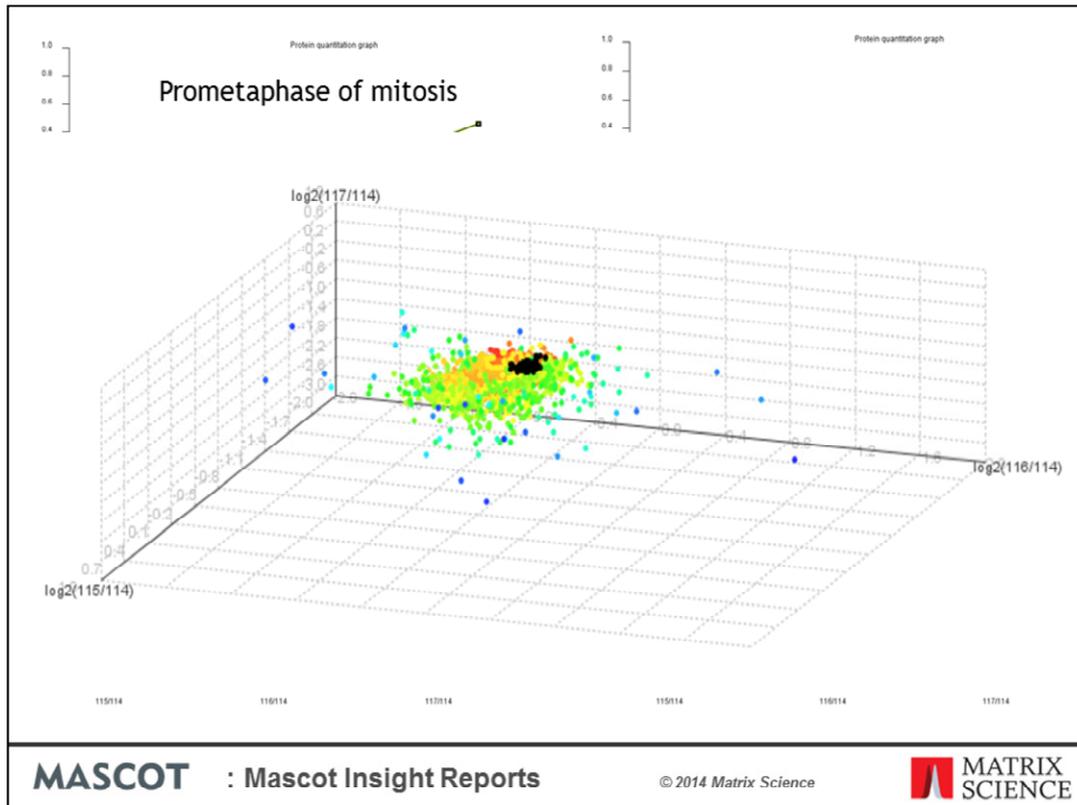
Computational protein profile similarity screening for quantitative mass spectrometry experiments.  
Kirchner M., et al. Bioinformatics. 2010 Jan 1;26(1):77-83.

**MASCOT** : Mascot Insight Reports

© 2014 Matrix Science

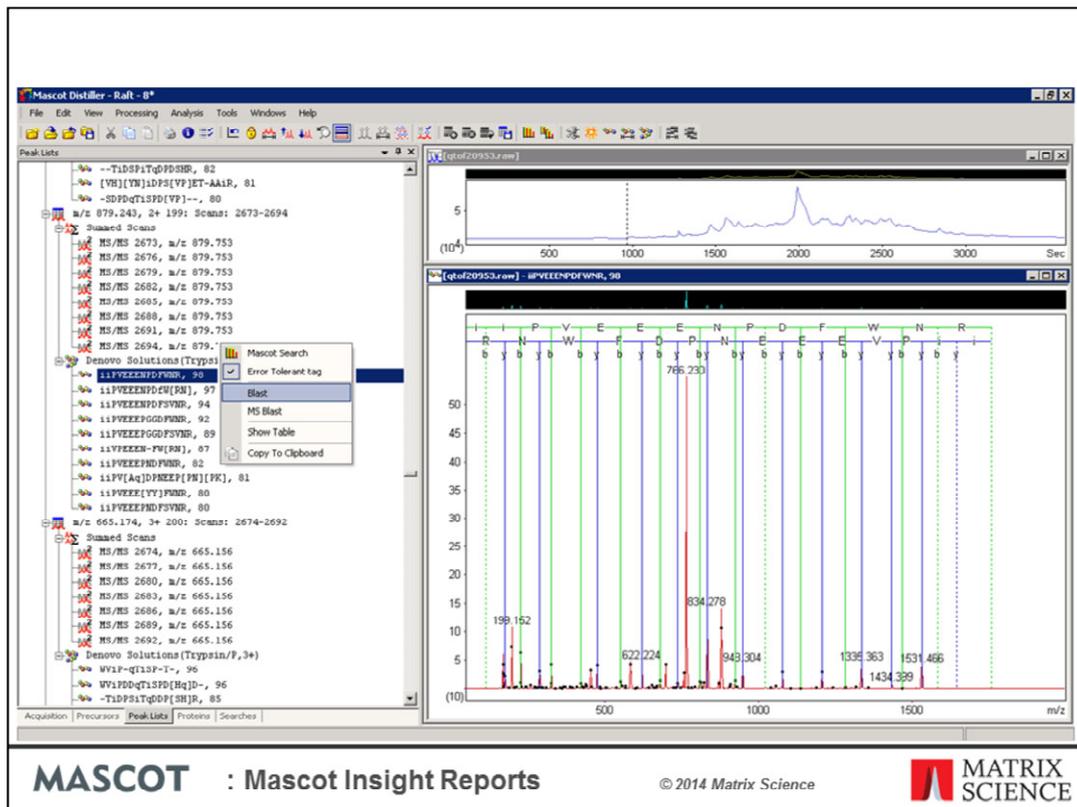
**MATRIX SCIENCE**

I want to show one other clustering plot, the K-means clustering report. The nature of the Super-SILAC data set, two tumors compared to a labelled standard was not a suitable data set for this type of analysis. Instead I have used a 4plex iTRAQ data set based on samples taken from 4 different points in the cell cycle. The data is from a collaboration between the labs of Darryl Pappin at Cold Spring Harbour Laboratory and Hanno Steen at the Children's Hospital, Boston.



The main difference between K-means clustering and hierarchical clustering is that K-means is a supervised clustering algorithm, which means that you have to specify the number of clusters the algorithm should go on and create. Here we have four different clusters from an iTRAQ dataset, where each ratio represents a different point in the cell cycle, G1 early, G1 late, S/G2 and Prometaphase. The four different clusters are showing different behaviour across those points. For example the proteins clustered together in the top left plot increase in concentration during the Prometaphase of mitosis period of the cell cycle while those below it are at their highest concentration during the Synthesis/Gap 2 period.

We can also plot a scatter graph of the clusters – here we have the data for that iTRAQ dataset plotted as a 3D scatter plot. The clusters generated are colour coded. You can rotate 3D scatter plots to get a better view of a particular cluster, hover over a data point and you get a tooltip giving you the protein name. You can also select individual clusters in order to highlight them. Here, I've selected the cluster which we're showing in the top left had graph.



This is a very different type of report from the previous ones in that it has nothing to do with quantitation. The Mascot Distiller search toolbox includes a powerful de novo sequencing algorithm. One of the most common reasons for carrying out a de novo sequence search is to try to find matches to MS/MS spectra which failed to get a match from a standard database search, which is what I've done here.

The data set is from a MHC (major histocompatibility complex) analysis using no enzyme and no modifications in the search settings. After the initial search of 1064 spectra there were a total of 290 matches and 774 unassigned spectra.

After the de novo search has completed we have a series of de novo sequence matches of various quality. One of the problems de novo sequencing is trying to verify the results – In this case, have we identified any additional matches to proteins we identified in the initial Mascot database search? The standard way try to answer this question would be to submit the high scoring de novo solutions for a sequence homology search, such as a BLAST or MS-BLAST search, or for an error tolerant sequence tag search. These approaches can be very time consuming as the process of submitting the searches is largely manual.

**Mascot Insight**

User: Patrick Emery | 10 unread notifications | Preferences | Administration | Help | Logout

Explore | Protein homology search for Mascot Dsiler de novo results from mhc.csi

Alignment summary | Alignment detail by protein | Alignment detail by query | Search details

| Accession | Description  | No. align... |
|-----------|--|--------------|
| P01964    | P01964 (Bos taurus) Hemoglobin subunit alpha                           | 27           |
| P02769    | P02769 (Bos taurus) Bovine serum albumin pr...                         | 16           |
| P00761    | P00761 (Sus scrofa) Sus scrofa (Pig)...                                | 15           |
| gi4502143 | cathepsin D preproprotein [Homo sapiens]                               | 14           |
| gi183766  | cathepsin D preproprotein [Homo sapiens]                               | 11           |
| gi406387  | L1F excision repair protein (RAD23) homolog B isoform 1 [Homo sapiens] | 9            |
| CCSD009   | T5246L_Q250102 (Bos taurus) similar to HSG protein                     | 8            |
| gi7661822 | dynem light chain roadblock-type 1 [Homo sapiens]                      | 6            |
| gi7657162 | pro-fascin subunit 6 [Homo sapiens]                                    | 6            |
| P81644    | P81644 (Bos taurus) Apolipoprotein A-II prec...                        | 6            |
| gi1912494 | HYNCA ribonucleoprotein complex subunit 2 isoform a [Homo sapiens]     | 3            |
| gi7657369 | NADH dehydrogenase [ubiquinone] 1 alpha subcomplex sub...              | 5            |
| gi406701  | H65 ribosomal protein S23 [Homo sapiens]                               | 5            |
| gi406695  | H65 ribosomal protein S19 [Homo sapiens]                               | 4            |
| gi4760774 | NADH dehydrogenase [ubiquinone] 1 beta subcomplex sub...               | 3            |
| gi4616928 | mitochondrial import receptor subunit TOM5 homolog isofo...            | 3            |
| gi8912714 | mitochondrial import inner membrane translocase subunit TL...          | 2            |
| gi6912634 | H65 ribosomal protein L13a isoform 1 [Homo sapiens]                    | 2            |

Match to Query 2159 rank 2 de novo score 34  
 de novo solution: DqEiDPiqqi  
 Expanded peptide sequence: DqEiDPiqqi length 10  
 Alignment score: 31  
 33 NKELDPIQKL 42  
 +KELDPIQKL  
 1 DqEiDPiqqi 10

**MASCOT** : Mascot Insight Reports © 2014 Matrix Science **MATRIX SCIENCE**

In Mascot Insight, we have implemented a report which allows you to use the protein hit sequences from a selected search result to carry out a BLAST like sequence homology search using de novo solutions, in order to try to find additional possible matches to spectra from your dataset, drilling down into the unassigned MS/MS spectra. The sequence homology search has been tailored to de novo solution data, and allows for Q->K and F->M\* in the alignment without penalty for example.

Out of the de novo results we took the 156 spectra with 1 or more solutions and a *De novo* score  $\geq 40$ , that is the good quality de novo matches. These 156 spectra expanded into 1082 *de novo* solutions which were tested against the 106 proteins identified from the original database search (does many more alignments than this because ambiguity in sequences is 'exploded'). The search results were 79 spectra with good matches to a previously identified protein.

(e.g. you could potentially ignore half of the of the spectra (79/156) with the best de-novo solutions when MHC hunting)

This allows us to see if we have any convincing de novo sequence matches to peptides from the protein hits from the initial Mascot database search – for example, this looks like a good match to a deamidated sequence also identified in the original search where the differences at the n-terminus prevented us from getting a sequence database match.

## Types of reports

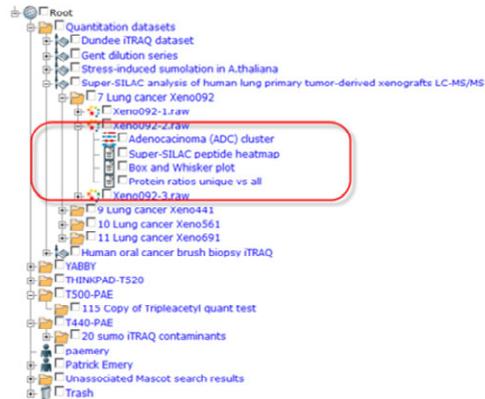
- **BLAST Cluster**
  - Use BLASTClust to make non-redundant sequence sets from one or more searches
- **Search level reports**
  - Many of the reports compare search results, some combine search results
- **Protein level reports**
  - Protein interaction reports, K-means clustering and statistical reports
- **Peptide level reports**
  - Ratio distributions, K-means clustering and statistical reports
- **Some of the reports only work with quantitative data** 
- ***De novo* report**

That's the end of the data analysis examples and I would like to quickly review the types of reports that Mascot Insight ships with.

And there is one special case of report a *De novo* report, which I have just shown, that works with the *De novo* results generated by Mascot Distiller.

## Sharing reports

- **Export reports**
  - Variety of formats, including CSV, SVG, PNG and HTML.
  - Publication export (MCP guidelines)
  - Or saved to the dataset explorer tree
- **Direct access to Mascot Insight**
  - For interactive reports collaborators can log on to Insight and evaluate the results directly
  - No limits to the number of concurrent users
  - Control access to the reports and Insight features via built in security



**MASCOT** : Mascot Insight Reports

© 2014 Matrix Science



I have shown examples of many different kinds of reports and I want to make it clear that reports are not locked into Insight. Once a report is generated, they can be exported in a variety of formats, including CSV, SVG, PNG and HTML, or saved to the dataset explorer tree to be easily viewed at a later date.

Members of the lab or collaborators can be given access to the data to make their own analysis or view any saved reports. As mentioned in the previous talk there is no limit on the number of users for the software and that goes for **concurrent users** too.

Access to the results is controlled by Mascot Server security using the standard groups and uses configuration.

## Search level reports - export

- **Statistics export**
  - Dataset export in tab delineated format for use in statistical analysis software (e.g. R, Perseus)
- **JTreeView export**
  - Export in JTreeView compatible pcl file format
- **Standard report**
  - Standard report for selected search in Excel format
- **Molecular and Cellular Proteomics compliant report**
  - Based on the MCP guidelines for publication

You can also export data out of Insight for use with other software.

The different export reports work for one or more search results and format then export data for use by third party software, simple excel use or for publication compliant reports.

Files created by the Statistics export, which exports protein, peptide and quantitation data in a tab delineated format, can be used by third party software like the R statistics package or Perseus from the Matthias Mann group at MPI Munich.

The Java TreeView format generates input files from quantified data read for use with the Java TreeView. Java TreeView is a program that visualizes hierarchical clustering of gene expression data, or in our case protein expression data.

The Standard report is equivalent to the a standard csv export from a Mascot Server web page report. It is only available for search results that have been merged together in Mascot Insight and for MzIdentML and ProtXML results. If the search result is available on the Mascot Server then this option will be disabled.

The MCP report meets the reporting guidelines from MCP for publication of data. The end result is a zip file that contains a report with all the protein, peptide and spectral information required by MCP.



The screenshot displays a Microsoft Excel spreadsheet titled "Publication\_MSMS\_Report\_for\_result\_mss-20130610-000001\_1377697837140 (1).xls". The spreadsheet is divided into several sections:

- Search Results:** A list of MS/MS spectra for various queries, including "MS/MS spectrum for query 3418 rank 1 from mss-20130610-000001 (Merged searches (Gent dilution series L/H 100)) 2 Oxidation (H/W)".
- Mass Spectrum Plot:** A plot showing relative intensity versus m/z, with a base peak at m/z 1469.09.
- Match Summary:**
  - Monoisotopic mass of neutral peptide M(calc):** 1469.09
  - Fixed modifications:** Carbamidomethyl (C)
  - Variable modifications:** M(10): Oxidation (M), with neutral losses 0 (shown in table), 64
  - Ions Score:** 63 Expect: 2.549E-6
  - Matches (Bold Red):** 25/204 fragment ions using 36 most intense peaks
- Match Table:** A table with columns: #, b, b\*\*, b<sup>0</sup>, b<sup>0\*\*</sup>, Seq, y, y\*\*, y<sup>0</sup>, y<sup>0\*\*</sup>, #. The table lists 8 matches, with the most significant ones highlighted in bold red text.

| # | b             | b**    | b <sup>0</sup> | b <sup>0**</sup> | Seq | y       | y**    | y <sup>0</sup> | y <sup>0**</sup> | #       |        |
|---|---------------|--------|----------------|------------------|-----|---------|--------|----------------|------------------|---------|--------|
| 1 | 72.04         | 36.53  |                |                  | A   | 1399.66 | 700.34 | 1382.64        | 691.82           | 1381.65 | 691.33 |
| 2 | 129.07        | 65.04  |                |                  | G   | 1342.64 | 671.82 | <b>1325.62</b> | 663.31           | 1324.63 | 662.82 |
| 3 | 200.1         | 100.56 |                |                  | A   | 1271.6  | 636.31 | 1254.56        | 627.79           | 1253.59 | 627.3  |
| 4 | 257.12        | 129.07 |                |                  | G   | 1214.58 | 607.8  | 1197.56        | 599.28           | 1196.57 | 598.79 |
| 5 | 344.16        | 172.58 | 326.15         | 163.58           | S   | 1127.56 | 564.28 | 1110.52        | 555.77           | 1109.54 | 555.27 |
| 6 | <b>415.19</b> | 208.1  | <b>397.18</b>  | 199.1            | A   | 1066.61 | 528.76 | 1039.49        | 520.25           | 1038.5  | 519.70 |
| 7 | <b>516.24</b> | 258.02 | <b>498.23</b>  | 249.02           | T   | 955.47  | 478.24 | 938.44         | 469.72           | 937.46  | 469.23 |
| 8 | 629.33        | 315.17 | <b>611.31</b>  | 306.16           | L   |         |        |                |                  |         |        |

**MASCOT : Mascot Insight Reports** © 2014 Matrix Science

Here is an example of the files generated by the Publication export report, which exports data for selected search results in a report designed to meet a number of the criteria specified by the Molecular and Cellular Proteomics guidelines for publication. The Publication report generates a zip archive containing an index HTML page, with links to the main Excel report file and links to static HTML pages which contain the spectrum views for MS/MS peptide matches (or MS protein matches for PMF data). A separate archive is generated for each search result, or you can generate a single report for a merged dataset.

## Create your own report

- Well documented API to creating your own report using the Java programming language
- Source code and coding walkthroughs for two example reports included in the help
- Uses the JFreeChart library for the charting

That brings me to the end of my examples of the reports currently included with Mascot Insight. In addition to the shipped reports, you can write your own – to do this you will need to have a working knowledge of the Java programming language.

We include a well documented API to creating your own report along with two documented examples to get you started.

MI uses the popular JFreeChart library for the charting.

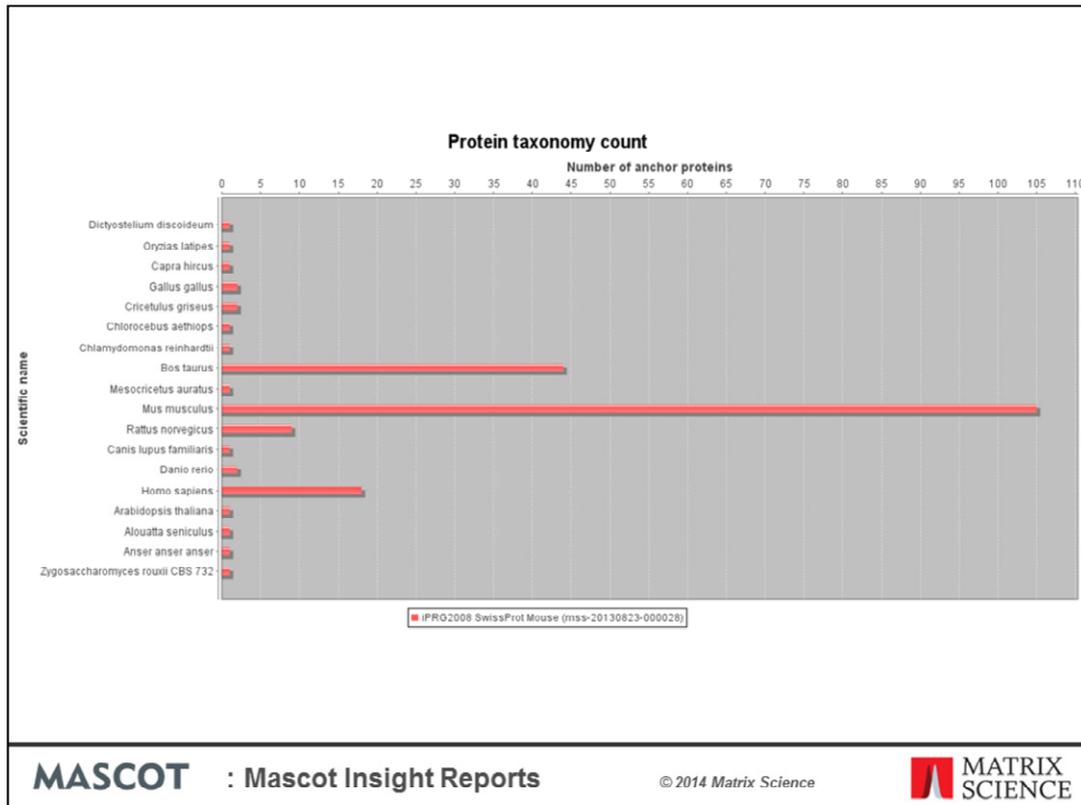
```

C:\backup\CVS_Home\Head\integral\src\java\com\matrixscience\tutorial\reports\ServerSideTaxonomyBarchart.java - Notepad++
File Edit Search View Encoding Language Settings Macro Run Plugins Window ?
ServerSideTaxonomyBarchart.java
187
188 /**
189  * This method is called on the Mascot Insight server side to carry out the
190  * work required to generate the data for the report
191  * @return true if the method runs, false otherwise
192  */
193 @Override
194 public boolean doMethodWork() {
195     try {
196         if(this.runReportType == RunReportType.SEARCH) {
197             // The JFreeChart barchart dataset object
198             DefaultCategoryDataset taxonomyDataset = new DefaultCategoryDataset();
199
200             // loop through each of the selected results
201             for(SearchEngineResult result : this.aResults) {
202                 HashMap<String,Integer> mTaxonomyCounts = new HashMap<String,Integer>();
203
204                 // check if we need to fetch the protein hit objects for the result
205                 boolean bFetchProteins = (result.getProteinHits() == null || result.getProteinHits().getNumberOfProteins() == 0);
206                 if(bFetchProteins) {
207                     // update the progress text to be displayed on the client side
208                     this.getProgressTracker().setProgressText("Fetching proteins for "+result.getSearchTitle());
209
210                     // use the setSearchAnchorProteinHits method from ReportPluginReport
211                     // to set the search protein hits. It is a helper method that
212                     // will populate the search anchor protein hits for you.
213                     this.setSearchAnchorProteinHits(result);
214
215                 }
216
217                 // use the getProteinHits().getVisibleProteins() method of Result
218                 // to fetch a Proteins array object containing only anchor protein hits that
219                 // have not been removed by any applied protein level filters
220                 Proteins proteins = result.getProteinHits().getVisibleProteins();
221
222                 // run through the protein objects in the Proteins container array
223                 for(int n = 1; n <= proteins.getNumberOfProteins(); n++) {

```

**MASCOT** : Mascot Insight Reports © 2014 Matrix Science 

You can write both chart type reports for use in MIRA, and export reports to generate export files in any format you choose. The example shown here will generate a barchart of taxonomy ids for a selected search result. The API we're using here is a high level API, so that the report works without any modification for Mascot results, and for searches imported from protXML and mzIdentML datasets.



And this is the output from the report run against the iPRG2008 dataset searched against the whole of SwissProt, with the largest number of protein hits being to mouse as expected since this is a mouse dataset.

**Mascot Insight**  
User: Patrick Emery | 0 unread notifications | Preferences | Administration | Back | Help | Logout

Explore  
C:\Program Files\Apach...

Searches  
VIT3\_DROME

| Hit | Accession   | Description   | Mass | Sca... | Peptides matched | sequence ... | % coverage | Reference | 115   | 116   |
|-----|-------------|---|------|--------|------------------|--------------|------------|-----------|-------|-------|
| 1   | VIT3_DROME  | Vitellogenin-3 OS=Drosophila melanogaster GN=Yp3 PE=1 SV=1                | 0    | 1.0    | 4                | D.1404762    | 14.05      | 0         | 0.868 | 1.86  |
| 2   | PDI_DROME   | Protein disulfide-isomerase OS=Drosophila melanogaster GN=Pd PE=1 SV=1    | 0    | 0.99   | 2                | D.050403236  | 5.04       | 0         | 0.91  | 2.027 |
| 3   | APU_DROME   | Adaptolipin OS=Drosophila melanogaster GN=Aplg PE=1 SV=2                  | 0    | 0.99   | 2                | D.005371531  | 0.54       | 0         | 0.259 | 1.947 |
| 4   | LGGG_DROME  | LGP-glucosylglycoprotein glucosyltransferase OS=Drosophila melano...      | 0    | 0.99   | 4                | D.031007752  | 3.1        | 0         | 0.391 | 1.348 |
| 5   | VDAC_DROME  | Voltage-dependent anion-selective channel OS=Drosophila melano...         | 0    | 0.99   | 2                | D.10283688   | 10.28      | 0         | 0.092 | 1.368 |
| 6   | SPTCA_DROME | Spectrin alpha chain OS=Drosophila melanogaster GN=alpha-Spec PE=1 SV=1   | 0    | 0.99   | 2                | D.007453416  | 0.75       | 0         | 0.74  | 1.214 |
| 7   | VL_DROME    | Putative vitellogenin receptor OS=Drosophila melanogaster GN=vl PE=1 SV=1 | 0    | 0.99   | 4                | D.022177419  | 2.22       | 0         | 0.071 | 0.396 |
| 8   | VIT1_DROME  | Vitellogenin-1 OS=Drosophila melanogaster GN=Yp1 PE=1 SV=1                | 0    | 0.99   | 3                | D.066059224  | 6.61       | 0         | 0.096 | 0.329 |
| 9   | GBLP_DROME  | Guanine nucleotide-binding protein subunit beta-like protein OS=Dros...   | 0    | 0.99   | 2                | D.056603774  | 5.66       | 0         | 0.003 | 1.224 |

No. proteins: 13

| Query No | Query Title | Charge | Observed | Inten... | M(Exp)  | M(Calc) | Delta | Start | End | Miss | Rank | ScaffoldPe... | Mascotscore                   | Mascotide... | Peptide          | Vars |
|----------|-------------|--------|----------|----------|---------|---------|-------|-------|-----|------|------|---------------|-------------------------------|--------------|------------------|------|
| 14       | Spec_551    | 2      | 1328.70  | 0        | 1328.76 | 1328.74 | 0.02  | 105   | 113 | 0    | 1    | 0.53          | 17.8                          | 6.41         | R LKVTETAK       | A    |
| 18       | Spec_366    | 2      | 735.41   | 0        | 1400.01 | 1400.70 | 0.03  | 200   | 207 | 0    | 1    | 0.94          | [ScaffoldPeptide Probability] |              | R ISDTLEYNAK     | S    |
| 26       | Spec_398    | 2      | 1119.11  | 0        | 2236.20 | 2236.26 | -0.05 | 190   | 208 | 0    | 1    | 0.94          | 60.73                         | 46.42        | K AASGLDLDL...   | R    |
| 27       | Spec_404    | 3      | 789.38   | 0        | 2365.12 | 2365.15 | -0.03 | 291   | 311 | 0    | 1    | 1             | 70.29                         | 44.47        | R ISGADPFVDAI... | C    |

No. peptides: 4

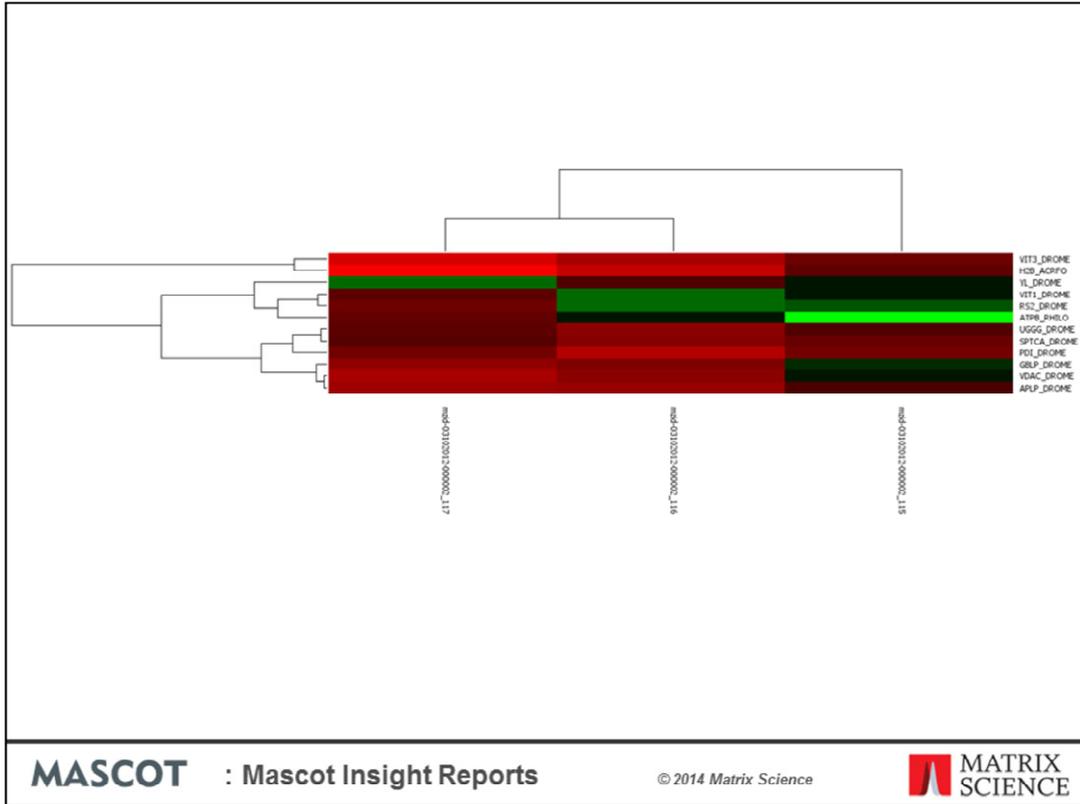
Protein View Taxonomy Family Annotation Spectrum Ions matched Error distribution Top spectrum matches Search Summary

Protein view for VIT3\_DROME Vitellogenin-3 OS=Drosophila melanogaster GN=Yp3 PE=1 SV=1

1 MHSLRICLLA TCLLVAAASG KQASNDRLKPK TSHLTATELE NVFSLNDITH  
51 ERLKNQPLEQ GARVIEKIYH VQIQKDLTP SFVSPSNVP VWIKNSIQK  
101 VECKLQNVY TARKQPFGE DEVTIVLIGL FKTSFAQQKA MRLLIQAVVQ  
151 KYHLQQLQNI AGEQQQLKLS SDYDTSSEE AADQKSAKA ASGDLTIIDL

**MASCOT** : Mascot Insight Reports © 2014 Matrix Science **MATRIX SCIENCE**

As I mentioned earlier in the presentation, you can import search result data from any source in either the protXML or mzIdentML standard formats. Here we are looking at an mzIdentML format result exported from Scaffold Q+, and we've also imported the Scaffold quantML quantitation result export file as well. We can use these results to generate reports, in the same way that we can do for Mascot search results.



As an example, here is the output from the hierarchical clustering report for this data set

## Summary

- **Generate reports covering**

- Quality control
- Quantitation
- Protein interactions

- **Without the use of additional software or Excel**

- **Facilitates in-depth analysis of data sets**
- **Reports can be saved for later review or exported**
- **Analyse Mascot as well as protXML and mzIdentML results.**
- **Write your own reports**
  - Java API

So, in summary, Mascot Insight ships with a large number of reports, which cover areas such as protein and peptide level comparisons between data sets and quantitation based reports including reports for comparison, ratio clustering and quality control. This is all done in Insight without the use of additional software or Excel.

You can also look at gene ontology and Protein interactions including mapping the quantitation results on the onto the interaction diagram.

As I have shown using Mascot Insight allows you to carry out in-depth analysis of a data set to determine both the quality of the data and to discover the significantly regulated proteins.

The results of these reports can be saved with the project or copied to Word or Powerpoint documents or exported into popular graphics formats.

You can analyse not only Mascot Server results but also results from any software that exports protXML or mzIdentML results

And, if you have a specific report you want to be able to generate, you can write your own reports in Java using a common API across different result formats