# Five Common Causes
# of
# Inaccurate Quantitation

**MASCOT**

*MATRIX SCIENCE*

I must make it clear that I am only considering data processing problems

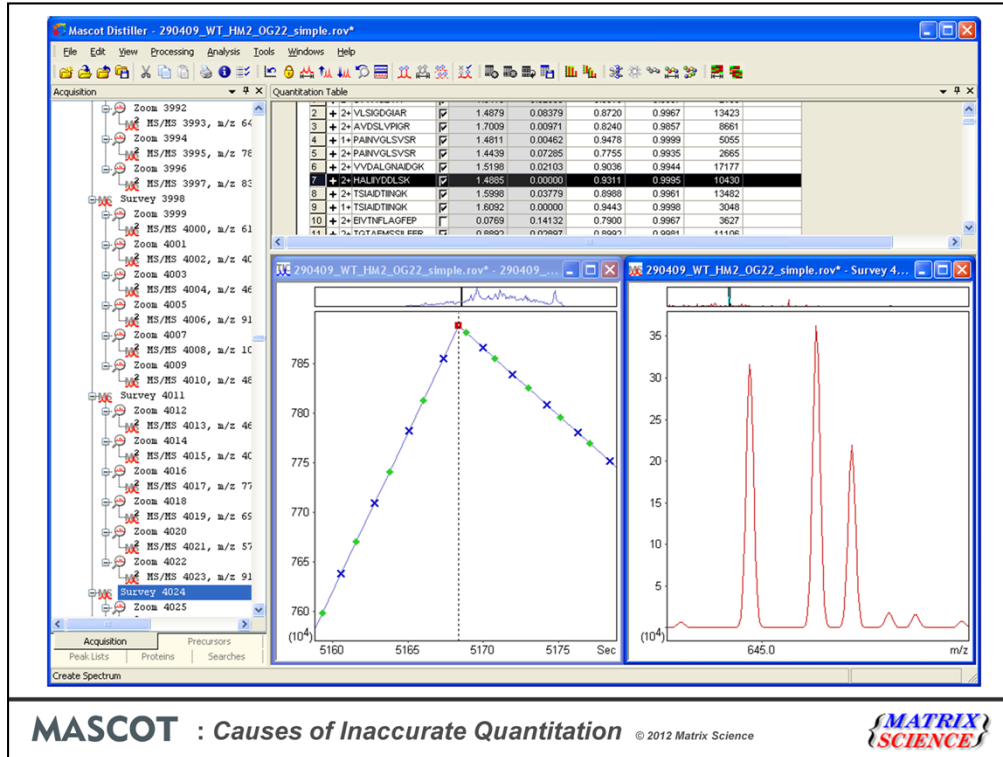MASCOT : *Causes of Inaccurate Quantitation* © 2012 Matrix Science

I hardly dare think about all the things that can go wrong at the bench, or when acquiring the mass spec. data

# 1. Low mass resolution
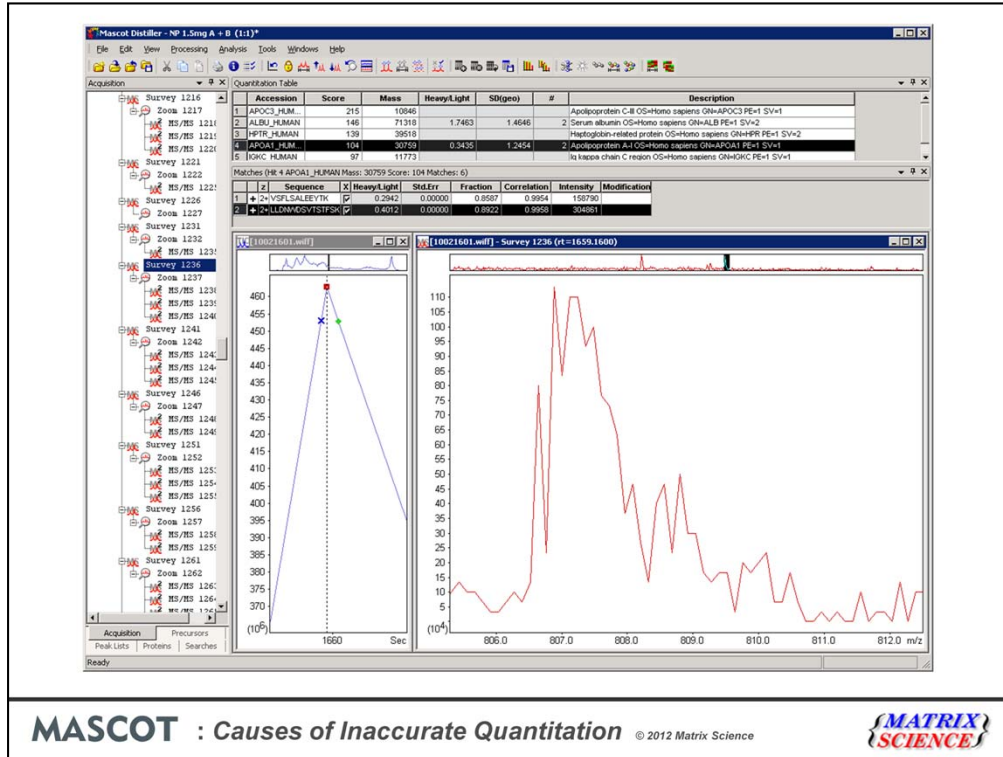
**For reliable quantitation, you need at least partial isotopic resolution**

- In the scans used for quantitation
- Especially for $^{18}O$

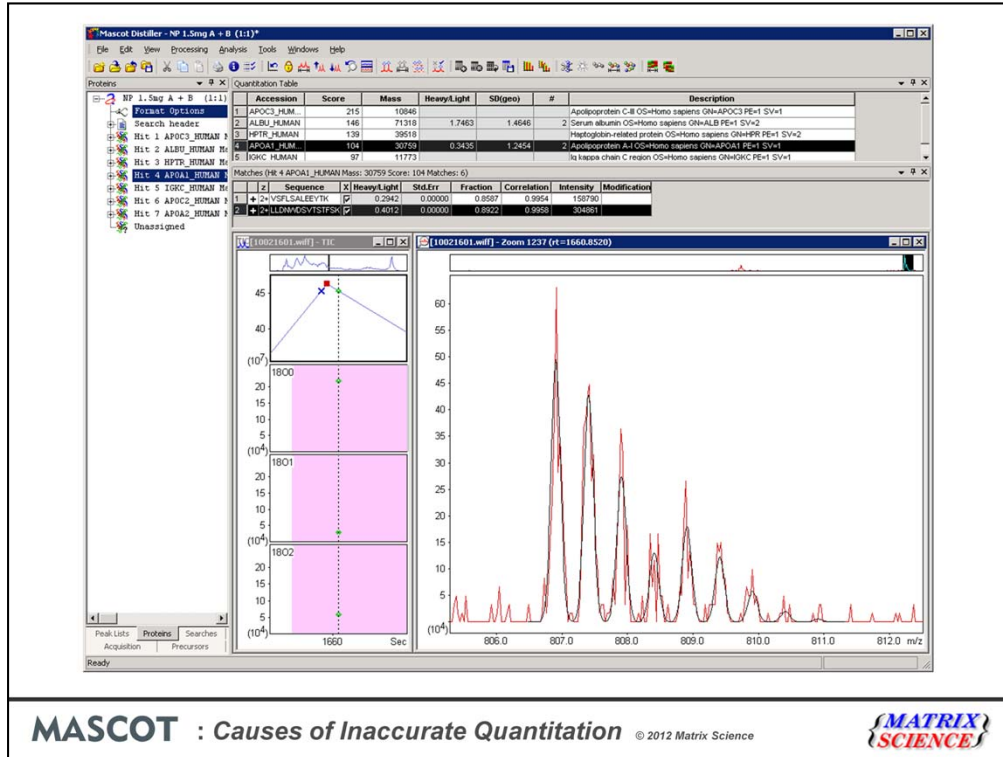**MASCOT** : *Causes of Inaccurate Quantitation* © 2012 Matrix Science

I think most people would agree that you need some degree of isotopic resolution to get reliable quantitation. If you are using a 'classic' ion trap, 3+ and even 2+ peaks may not show any isotopic resolution

MASCOT : *Causes of Inaccurate Quantitation* © 2012 Matrix Science

This is some 18O data from a 'classic' trap. This particular peptide gets a strong match. At first glance, the resolution in the precursor region of a survey scan looks pretty good. But, if you look more carefully, these peaks are not from resolved isotope distributions. For a start, there aren't enough peaks. The data have been saved to the raw file as centroids, not profile, as is common practice with traps.

MASCOT : *Causes of Inaccurate Quantitation* © 2012 Matrix Science

This is a different file where the survey scans have been saved as profile data. Now we can see the true picture. When an unresolved distribution such as this is centroided, it gets broken up into peaks in an arbitrary way. Trying to use such survey scans for any type of quantitation would be difficult, whether saved as profile data or centroids. For 18O labelling, the situation is hopeless because the separation between heavy and light is only 4 Da, and it is essential to deconvolute the distributions.
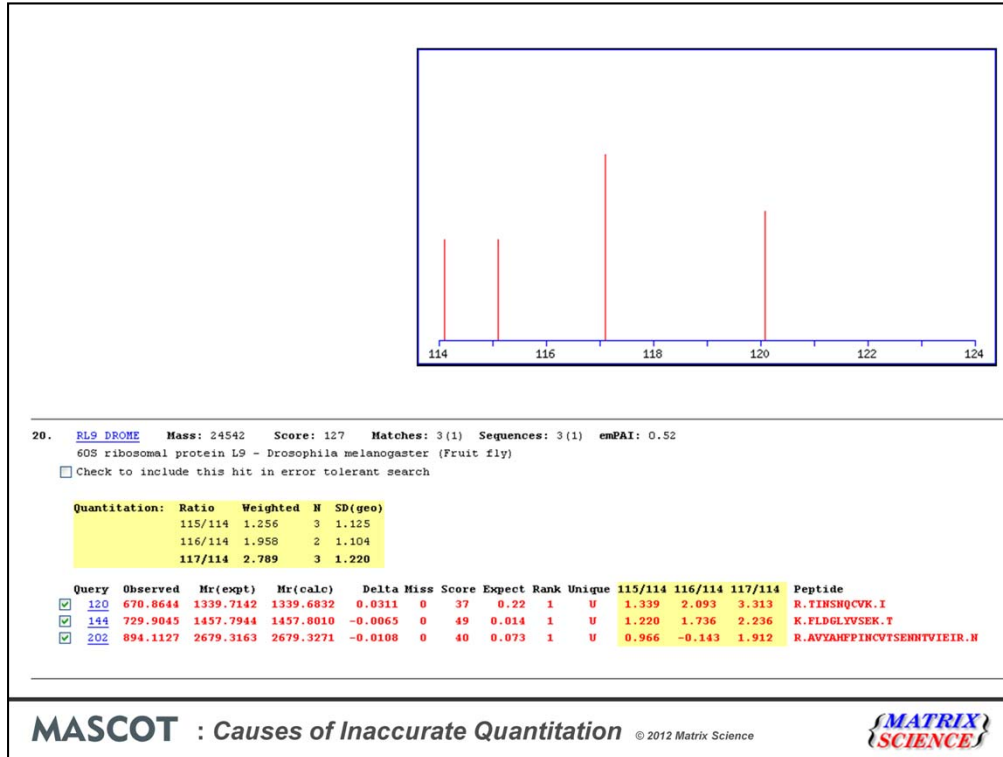
MASCOT : *Causes of Inaccurate Quantitation* © 2012 Matrix Science

If we look at an adjacent zoom scan, we can see what the isotope pattern should look like. Signal to noise is still not great, but deconvolution becomes possible when the peaks are fully resolved. So, with zoom scans, even though you might only have a single scan for each precursor, you can get reasonable results from a standard trap.
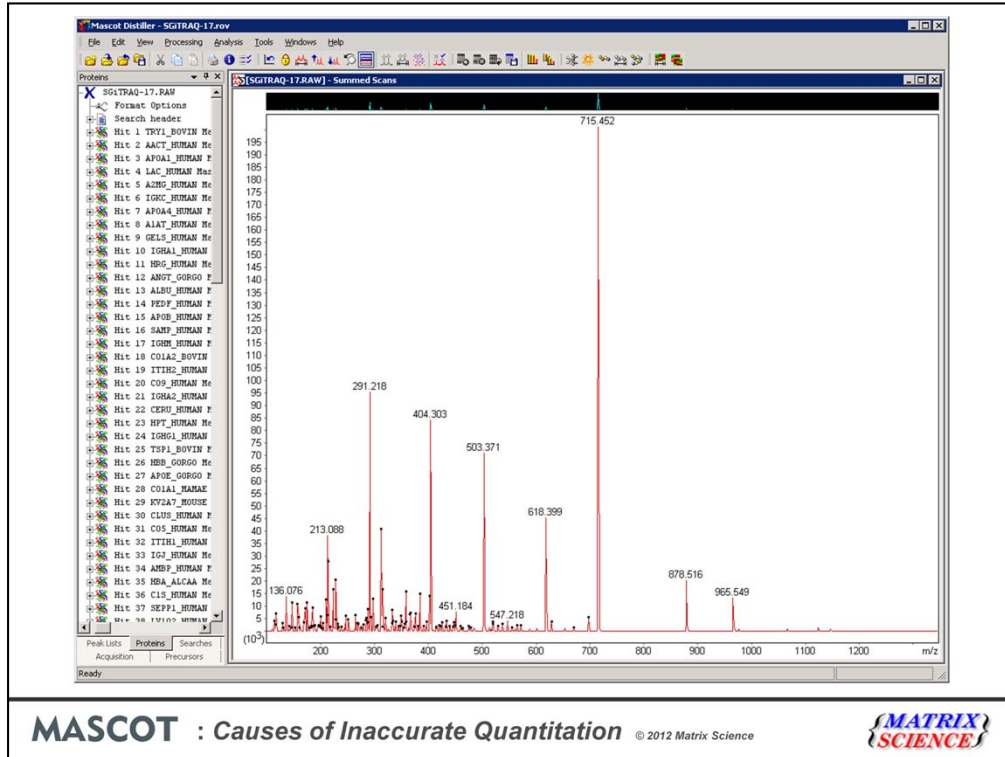
## 2. Reporter ion peak picking

**Reporter ions are not peptide fragments**

If you are using iTRAQ or TMT, it is very important to understand that the reporter ions are not peptide fragments. Make sure your peak picking software doesn't try to apply some standard de-isotoping algorithm, designed for peptides. This can only distort the relative intensities of reporter ion peaks.

**20.** RL9_DROME  Mass: 24542  Score: 127  Matches: 3(1)  Sequences: 3(1)  emPAI: 0.52
  60S ribosomal protein L9 - Drosophila melanogaster (Fruit fly)
  ☐ Check to include this hit in error tolerant search

| Quantitation: | Ratio | Weighted | N | SD(geo) |
|---|---|---|---|---|
| | 115/114 | 1.256 | 3 | 1.125 |
| | 116/114 | 1.958 | 2 | 1.104 |
| | 117/114 | 2.789 | 3 | 1.220 |

| Query | Observed | Mr(expt) | Mr(calc) | Delta | Miss | Score | Expect | Rank | Unique | 115/114 | 116/114 | 117/114 | Peptide |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| ☑ 120 | 670.8644 | 1339.7142 | 1339.6832 | 0.0311 | 0 | 37 | 0.22 | 1 | U | 1.339 | 2.093 | 3.313 | R.TINSNQCVK.I |
| ☑ 144 | 729.9045 | 1457.7944 | 1457.8010 | -0.0065 | 0 | 49 | 0.014 | 1 | U | 1.220 | 1.736 | 2.236 | K.FLDGLYVSEK.T |
| ☑ 202 | 894.1127 | 2679.3163 | 2679.3271 | -0.0108 | 0 | 40 | 0.073 | 1 | U | 0.966 | -0.143 | 1.912 | R.AVYAHFPINCVTSENNTVIEIR.N |

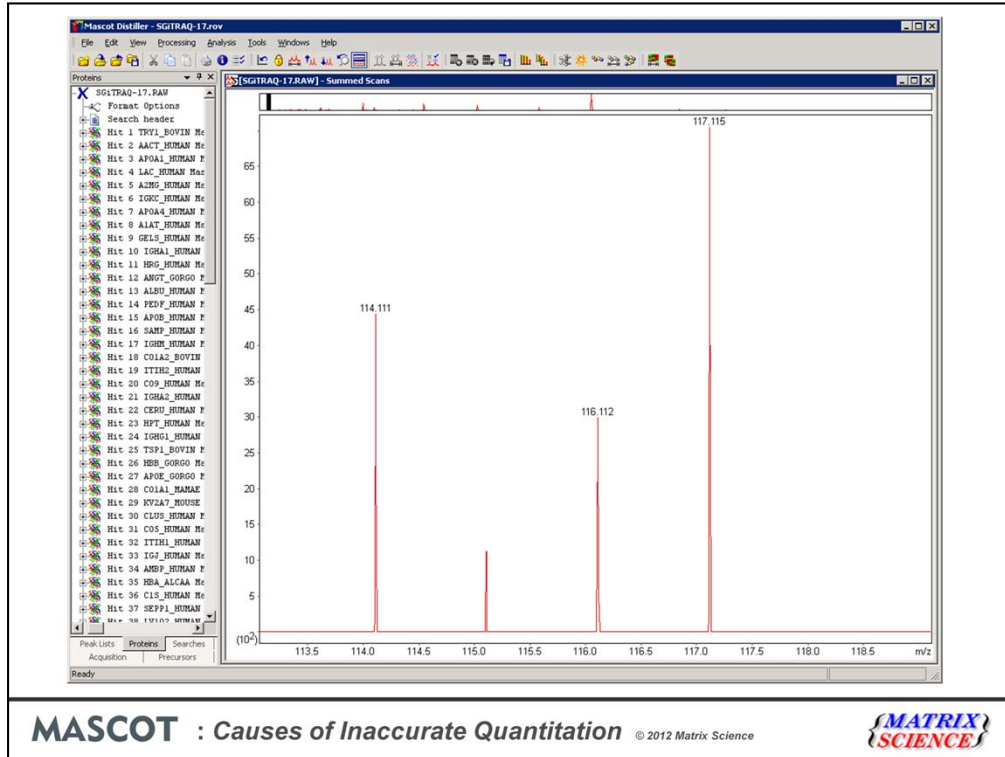**MASCOT** : *Causes of Inaccurate Quantitation* © 2012 Matrix Science  MATRIX SCIENCE

A more serious problem is unreliable peak picking. If you look at a Mascot quantitation report for a reporter ions experiment, you may see large numbers of negative ratios. These are where the peak picking has missed a peak completely, giving a raw intensity of zero. The isotope correction then removes a little more intensity, and donates it to the adjacent peaks, so that the missing peak goes negative. Here, for example, the 116 has been missed.

We decided not to suppress these negative ratios because they are a strong indicator that something is wrong with the peak picking. Usually, the problem is a setting that would be fine for sequence ions, such as 'ignore peaks less than 1% of the base peak intensity'

Here, for example, the base peak has an intensity of approx 200,000 widgets. The reporter ions are relatively weak peaks, down at the bottom left

MASCOT : *Causes of Inaccurate Quantitation* © 2012 Matrix Science

If we zoom in, we can see that the 115 peak has been missed. 1% of the base peak is an intensity of 2000 and the 115 is below this. Not a problem for protein identification but a huge problem for iTRAQ quantitation. So, very important to ensure your peak picking settings are correct for these peaks

# 3. SILAC: Arg-Pro Conversion

**MASCOT** : *Causes of Inaccurate Quantitation*  © 2012 Matrix Science   *MATRIX SCIENCE*

SILAC is extremely popular. Not everyone is aware of Arg-Pro conversion. Ong and colleagues reported how cells grown in media containing labelled arginine could yield peptides containing labelled proline. To obtain an accurate ratio, it becomes necessary to account for the label distributed across these additional peaks.

3. SILAC: Arg-Pro Conversion
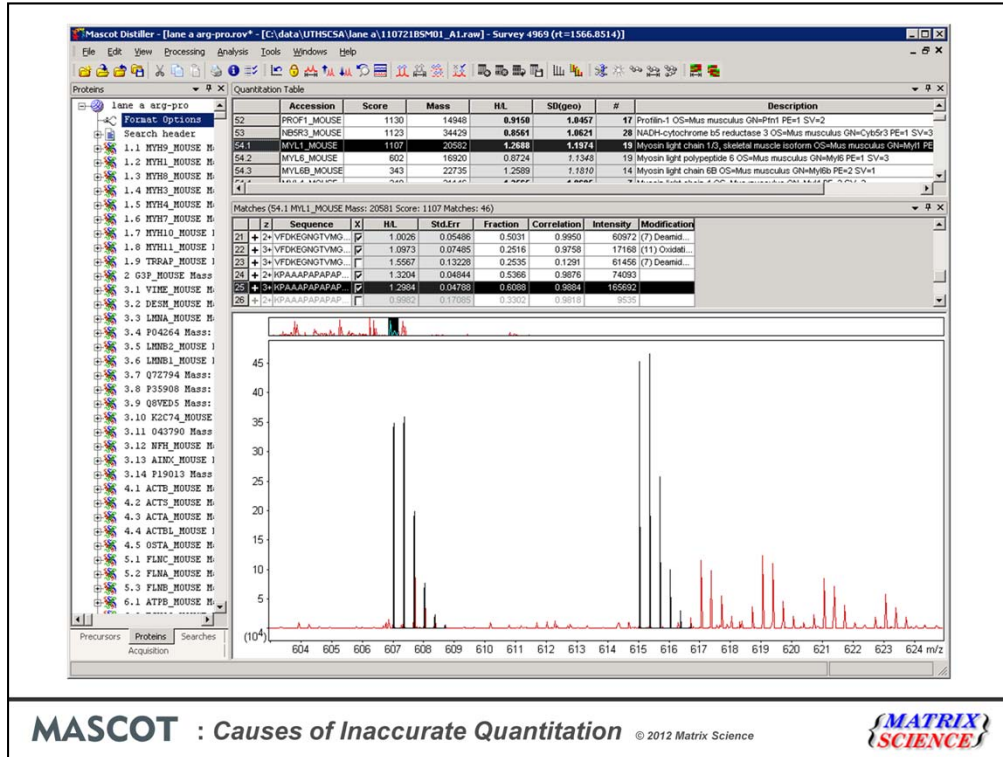
MASCOT : *Causes of Inaccurate Quantitation* © 2012 Matrix Science

Here is an example for Arginine labelled with $13C(6)15N(4)$, $+10$. Some of the label has been incorporated as Proline The proline label is not identical to the arginine label. In this case, it is $13C(5)15N(1)$, $+6$. To get an accurate ratio, you need to sum the area of the two heavy distributions.

MASCOT : *Causes of Inaccurate Quantitation* © 2012 Matrix Science

Not everyone sees this problem, and there are ways to minimise it. But, take a close look at your data from time to time. Here is a case where it is very strong and the ratios are seriously distorted. Without a correction, we only integrate the first heavy distribution, overlayed in black, and the ratio is 0.2 rather than 1

MASCOT : *Causes of Inaccurate Quantitation* © 2012 Matrix Science

With a correction, we sum all of the distributions and the ratio is closer to those of the non-proline containing peptides

**4. Modified peptides**

**Two samples**

**Abundance of protein X is same in both**
- In sample 1, peptide Y is 1% phosphorylated
- In sample 2, peptide Y is 3% phosphorylated

**Do we want to include phospho-peptide Y in the quantitation of protein X?**

**What about unmodified peptide Y?**

MASCOT : *Causes of Inaccurate Quantitation* © 2012 Matrix Science

How should modified peptides be handled when we are interested in relative quantitation of proteins? I get the distinct impression that many people don't give this a great deal of thought.

Consider this case. The abundance of the protein is the same in both samples but one of the peptides carries a low level of phosphorylation: 1% in one sample and 3% in the other. Clearly, we want to exclude this peptide because it will give us a ratio of 1:3 rather than 1:1.

Using the unmodified peptide is fine, because we'll get a ratio of 99:97, which in most cases will be indistinguishable from 1:1

# 4. Modified peptides

**Two samples**

**Abundance of protein X is same in both**

- In sample 1, peptide Y is 30% deamidated
- In sample 2, peptide Y is 70% deamidated

**Do we want to include deamidated peptide Y in the quantitation of protein X?**

**What about unmodified peptide Y?**

**MASCOT** : *Causes of Inaccurate Quantitation* © 2012 Matrix Science  {MATRIX SCIENCE}

What about a peptide that is more extensively modified? Again, the abundance of the protein is the same in both samples. In one sample, a hypothetical peptide is 30% deamidated, in the other 70% deamidated. We want to exclude this peptide because it will give us a ratio of 3:7 rather than 1:1. Unlike the previous case, the unmodified peptide is no better, giving a ratio of 7:3

16

# 4. Modified peptides

## If you are trying to quantify proteins:

- Perform a preliminary search to identify abundant, non-quantitative modifications
- Include such modifications in the search but exclude both modified and unmodified peptides from quantitation
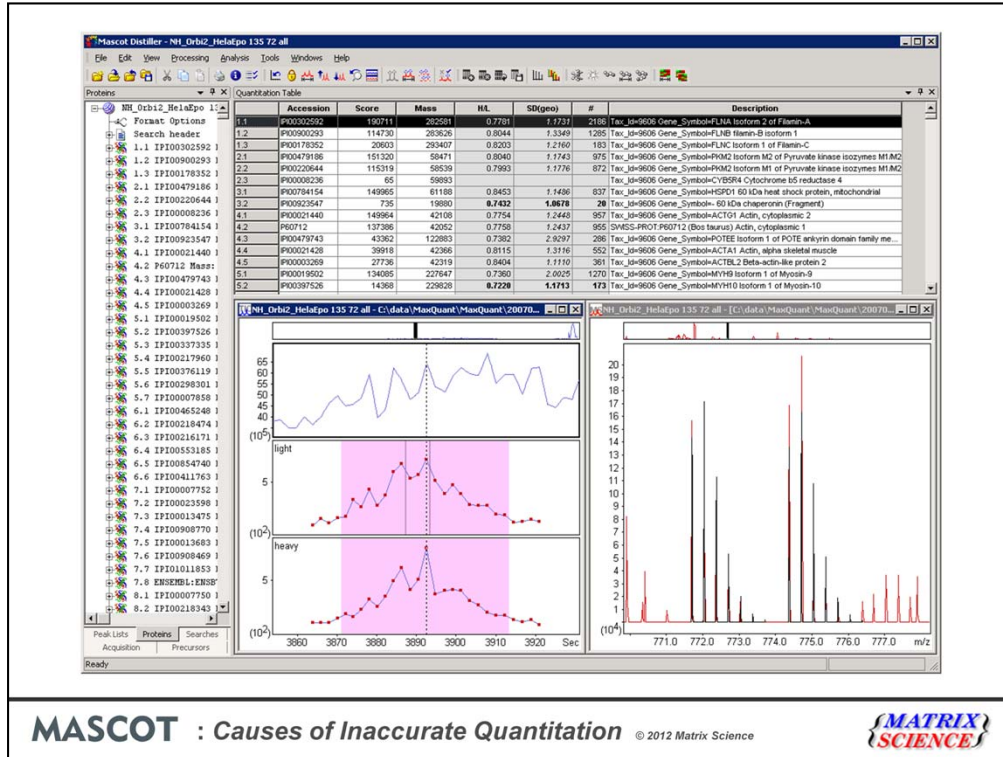- For low abundance modifications, don't include the modification in the search

**MASCOT** : *Causes of Inaccurate Quantitation* © 2012 Matrix Science  {MATRIX SCIENCE}

So, for relative quantitation of proteins…

**5. Too few peptides**

**"There is safety in numbers"**

MASCOT : *Causes of Inaccurate Quantitation* © 2012 Matrix Science
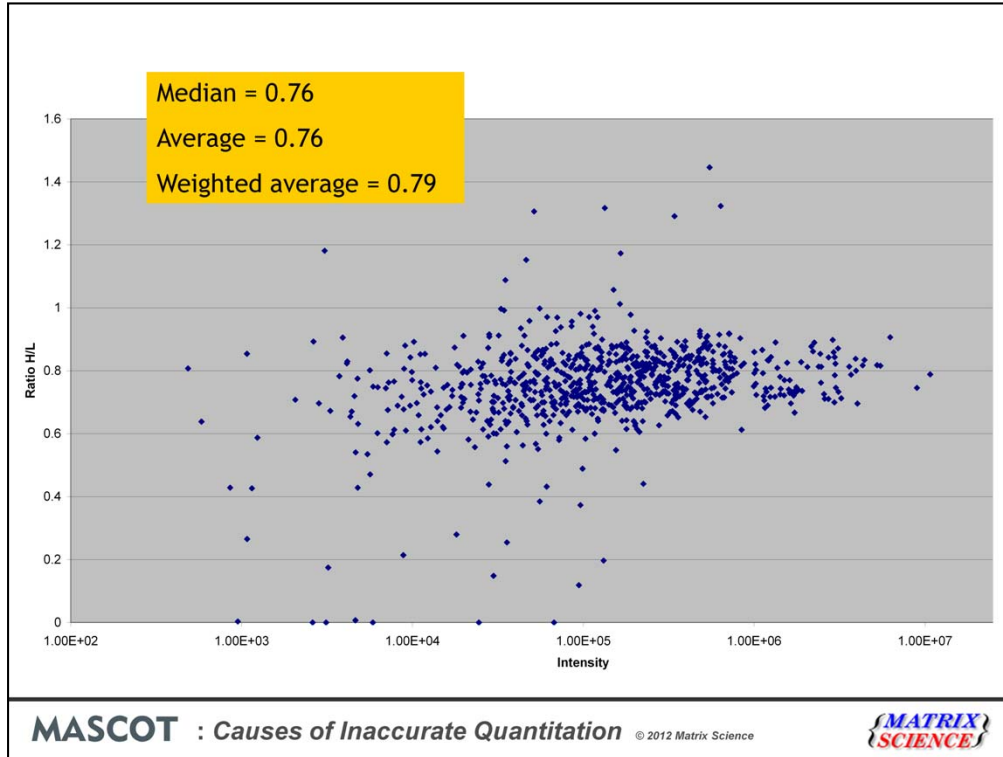
Finally, I suggest the main cause of inaccurate quantitation in a discovery experiment is having insufficient data. Peptide abundance is a surrogate for protein abundance. We assume, or rather hope, that the two are closely coupled. This only becomes a safe assumption when you look at a good sized population of peptides, and eliminate the outliers.
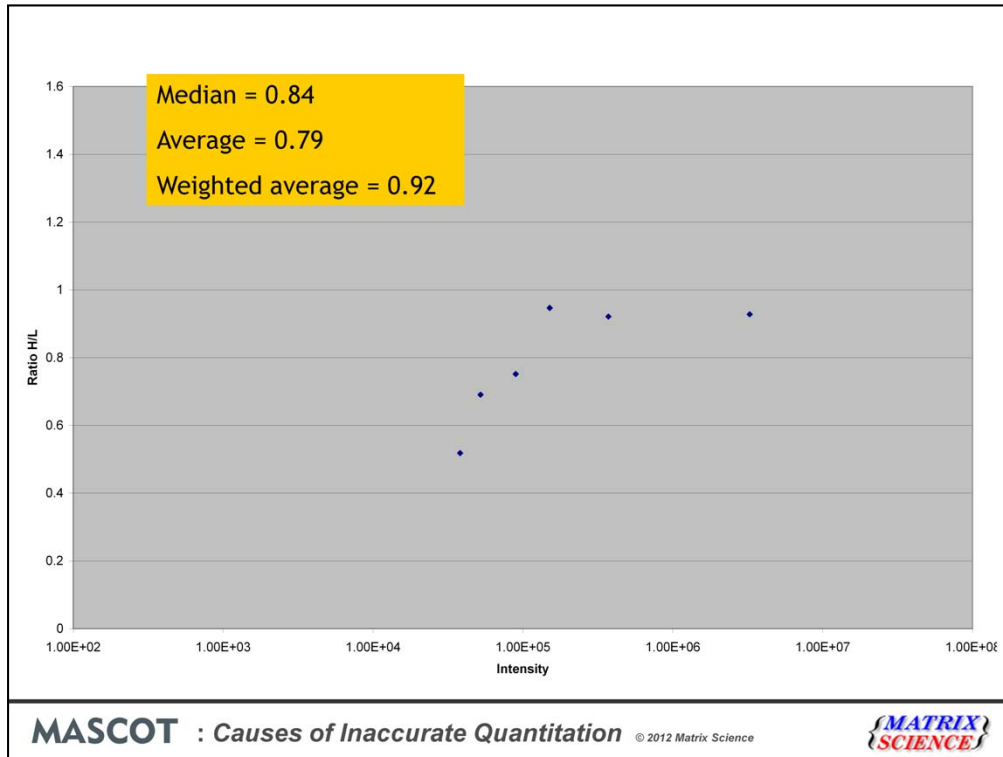
One of the main reasons for peptide abundance being different from protein abundance was just discussed: modified peptides. For post-digest labelling, another factor might be the enzyme digest conditions.

Here's a very nice SILAC data set, containing over 4000 proteins in the minimal list and some of these have over a thousand peptide matches. High mass resolution and accuracy
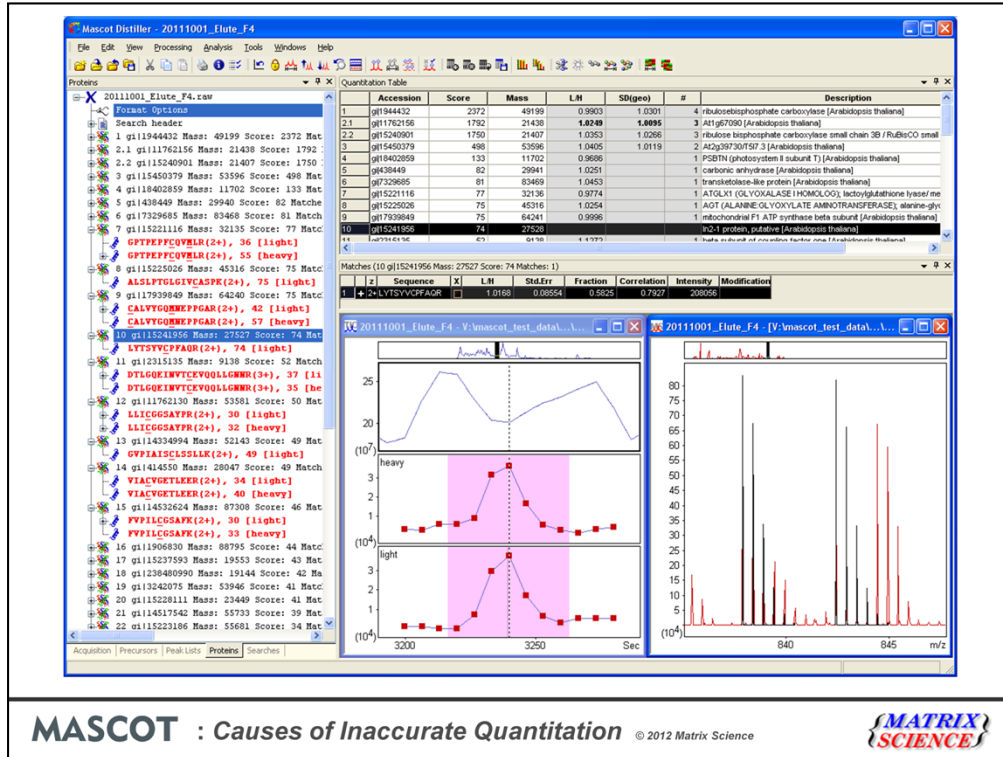
Here is a plot of 838 peptide ratios for a protein near the top of the list. The y axis should really be logarithmic, but a linear scale makes it easier to visualise the data. Possibly the extreme measurements are outliers caused by some failure in peak picking or chromatogram integration. Possibly they are peptides that have been misassigned. Possibly they are modified or processed in some way that makes them not representative of the abundance of the protein. In doesn't really matter because, when you have this many measurements, you can see where the centre of gravity is. Somewhere just under 0.8, yes? It doesn't matter whether you take the mean or the weighted mean or the median … we still get the same ratio.

MASCOT : *Causes of Inaccurate Quantitation* © 2012 Matrix Science

Go further down the list, to the low abundance proteins, where you have handful of measurements. You don't need statistics to tell you that this is a less reliable measurement. For these 6 ratios, the average and weighted average are quite different. When people ask me which is the 'best' way to calculate the protein ratio from the peptide ratio, I'm tempted to reply that, if it makes much difference, you need more data.

The really dangerous situation is when these 6 peptide ratio measurements are all for the same peptide sequence, or maybe two sequences. Then it becomes a lottery.

MASCOT : *Causes of Inaccurate Quantitation* © 2012 Matrix Science

For this reason, I'm not a fan of methods that focus on a small number of peptides for each protein. Methods such as ICAT or COFRADIC. The idea is to simplify the problem. But I feel they throw out the baby with the bathwater.

This is an ICAT example. The data quality is beautiful. High mass resolution. Clean and symmetrix XIC peaks. But, with only one peptide for most proteins, I simply don't feel confident that we are getting reliable protein quantitation.

| Accession | Description | All | | | | Unique | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | H/L | SD(geo) | # match | # seq | H/L | SD(geo) | # match | # seq |
| IPI00479186 | PKM2 Isoform M2 of Pyruvate kinase isozymes M1/M2 | 0.804 | 1.174 | 975 | 43 | 0.8385 | 1.142 | 107 | 4 |
| IPI00220644 | PKM2 Isoform M1 of Pyruvate kinase isozymes M1/M2 | 0.7993 | 1.178 | 872 | 41 | 0.8061 | 1.059 | 4 | 2 |
| IPI00007752 | TUBB2C Tubulin beta-2C chain | 0.7437 | 1.156 | 777 | 24 | 0.821 | 1.066 | 28 | 1 |
| IPI00023598 | TUBB4 Tubulin beta-4 chain | 0.7325 | 1.152 | 637 | 20 | 0.8768 | 1.062 | 23 | 1 |
| IPI00013475 | TUBB2A Tubulin beta-2A chain | 0.7292 | 1.157 | 629 | 22 | 0.7062 | 1.174 | 7 | 3 |
| IPI00021812 | AHNAK Neuroblast differentiation-associated protein | 0.8676 | 1.29 | 1154 | 152 | 0.8577 | 1.331 | 752 | 124 |
| IPI01012911 | clone CTONG2004264, moderately similar to AHNAK | 0.88 | 1.198 | 403 | 29 | 0.9462 | | 1 | 1 |
| IPI00024067 | CLTC Isoform 1 of Clathrin heavy chain 1 | 0.7672 | 1.263 | 725 | 64 | 0.75 | 1.159 | 493 | 49 |
| IPI00022881 | CLTCL1 Isoform 1 of Clathrin heavy chain 2 | 0.7996 | 3.112 | 219 | 15 | 0.000806 | 46.13 | 7 | 2 |
| IPI00385931 | PRO2051 | 0.8098 | 1.147 | 17 | 3 | 0.6116 | | 1 | 1 |
| IPI00643920 | cDNA FLJ54957, highly similar to Transketolase | 0.7896 | 1.143 | 859 | 29 | 0.8107 | 1.203 | 24 | 1 |
| IPI00793119 | cDNA FLJ56274, highly similar to Transketolase | 0.7893 | 1.14 | 835 | 28 | | | 0 | 0 |
| IPI00940673 | cDNA FLJ53217, highly similar to Transketolase | 0.7814 | 1.143 | 730 | 28 | 0.8422 | 1.029 | 4 | 2 |

**MASCOT** : *Causes of Inaccurate Quantitation* © 2012 Matrix Science    {MATRIX SCIENCE}

A related question is whether it's a good idea to restrict quantitation to "unique peptide". That is, peptides that are not shared with other proteins.  Here are a few examples from a large SILAC data set. There is solid evidence for the presence of all of these proteins from high scoring, unique peptides. But, when we look at the peptides that are quantified, a very high proportion are shared between isoforms. For example, these tubulins. Each have some 6 or 7 hundred matches to twenty odd distinct sequences. But, almost all of these are shared. When eliminated, we end up with just 1 distinct sequence each for two of the tubulins. Too few for any kind of reliable measurement.

Does removal of the shared matches reveal any up or down regulation? You may think you see one here, the H/L for this Clathrin goes from 0.8 to near zero. However, note that it is down to only 7 matches to two distinct sequences. If we look at what these are

| Sequence | Incl. | H/L | Std.Err. | Fraction | Correlation | Intensity | Modifications |
|---|---|---|---|---|---|---|---|
| AQILPVR | X | 0.000806 | 0.000062 | 0.3247 | 0.9816 | 5.87E+04 | Acetyl (Protein N-term) |
| AQILPVR | X | 0.000146 | 0.000239 | 0.4353 | 0.9878 | 3.08E+04 | Acetyl (Protein N-term) |
| AQILPVR | X | 0.000128 | 0.000002 | 0.3366 | 0.9888 | 7059 | Acetyl (Protein N-term) |
| AQILPVR | X | 0.000543 | 0.03647 | 0.3159 | 0.9813 | 8992 | Acetyl (Protein N-term) |
| AQILPVR | X | 0.002587 | 0.000374 | 0.3058 | 0.9835 | 1.91E+04 | Acetyl (Protein N-term) |
| QNLQLCVQVASK | X | 0.948 | 0.08599 | 0.6908 | 0.9823 | 5.93E+04 | |
| QNLQLCVQVASK | X | 0.8416 | 0.1362 | 0.3082 | 0.9818 | 4.76E+04 | |

**MASCOT** : *Causes of Inaccurate Quantitation* © 2012 Matrix Science    {MATRIX SCIENCE}

Matches to one of the sequences are both in the 0.8 ballpark. The other sequence is post-translationally modified, which makes it unreliable for quantitation of the protein.

It doesn't always make sense to limit discovery quantitation to unique peptides. Maybe better to study cases where the variance of the measurements is larger than expected and see whether there is evidence for the peptides belonging to two populations

**Five Common Causes of Inaccurate Quantitation**

1. Low mass resolution
2. Unreliable reporter ion peak picking
3. Arg-Pro conversion in SILAC
4. Modified peptides
5. Too few peptides

MASCOT : *Causes of Inaccurate Quantitation* © 2012 Matrix Science

To summarise