# In-Depth Look at Target-Decoy Validation

**MASCOT**

## Guidelines

### Molecular & Cellular Proteomics
- http://www.mcponline.org/misc/ParisReport_Final.shtml

### Proteomic Standards Initiative / MIAPE
- http://www.psidev.info/index.php?q=node/91

### Proteomics
- http://dx.doi.org/10.1002/pmic.200500856

---

**MASCOT** : *Target-Decoy Validation*    © *2008 Matrix Science*    *MATRIX SCIENCE*
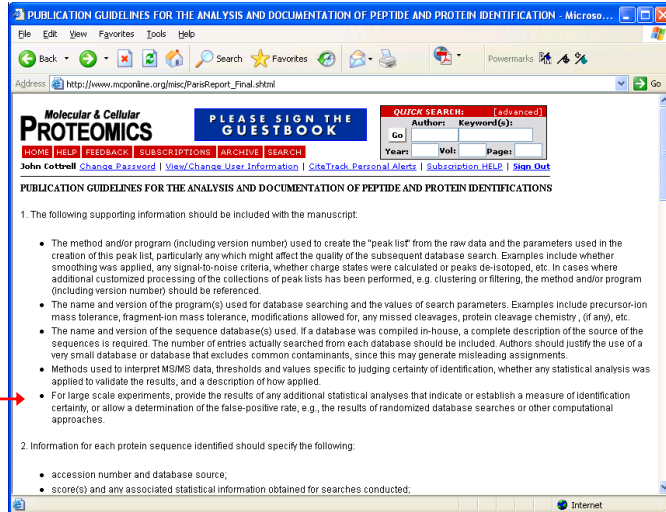
Whatever search algorithm you use, it is simply good scientific practice to validate the results. In particular, verify that the false discovery rate is not greater than claimed.

There is growing concern that some of the results in the literature, particularly from large scale searches, are not as reliable as they could be.

This has led to calls for greater stringency in the reporting of proteomics results. Most notably, the initiative taken by the Editors of Molecular and Cellular Proteomics, who held a workshop in 2005 to define a set of guidelines. The MCP guidelines have also been adopted by the Journal of Proteome Research

The PSI guidelines, which are part of the MIAPE standard, are very similar to those from MCP. The editorial board of Proteomics has also drafted a similar document. Hopefully, at some point, these parties and others will agree to a common standard.

# Guidelines

**MASCOT** : *Target-Decoy Validation*    *© 2008 Matrix Science*    *MATRIX SCIENCE*

One of the specific recommendations in the MCP guidelines is "For large scale experiments, provide the results of any additional statistical analyses that indicate or establish a measure of identification certainty, or allow a determination of the false-positive rate, e.g., the results of reversed or randomized database searches or other computational approaches."

## The Database

### "Decoy" database

- Direct estimate of false discovery rate (FDR)
- Requires large dataset to get accurate estimate of FDR
- Not a substitute for a reliable scoring scheme
- What makes a good decoy database?

**MASCOT** : *Target-Decoy Validation* © *2008 Matrix Science* *MATRIX SCIENCE*

I think it was the Gygi group first coined the term "decoy database search" for this approach. The idea is to repeat the search, using identical search parameters, against a database in which the sequences have been reversed or scrambled.

You do not expect to get any significant matches from the decoy database. So, the number of matches that are found is a good estimate of the number of false positives in the results from the target database.

This is an excellent validation method for MS/MS searches of large data sets. It is not as useful for a search of a small number of spectra, because the numbers are too small to give an accurate estimate of the false discovery rate. Hence, it is not a substitute for a reliable scoring algorithm.

What are the requirements for a decoy database?

# The Database

## We want database entries that

- Look like "real" proteins to the search algorithm
- Do not contain any genuine matches.

This is actually quite a difficult question.

We want database entries that look like "real" proteins to the search algorithm

However, we want database entries that don't contain genuine peptide sequences

## The Database

### Reversed entries

- Common approach for MS/MS with enzyme
- Not suitable for MS/MS without enzyme
  - Can get C-term series swap with N-term
- True palindromic peptides are rare
  - H12_MOUSE: K.AVKPKAAKPKVA.K
- Preserves entries with unusual compositions
  - For example, collagens have high abundances of G & P
- Preserves redundancy
  - Get approx same number of unique precursors

**MASCOT** : *Target-Decoy Validation*   © 2008 Matrix Science   *MATRIX SCIENCE*

The Gygi group advocate simply reversing the entries in the original database. This is a reasonable approach for an MS/MS search where a specific enzyme has been used.

It is not suitable for no-enzyme MS/MS searches, especially when there are several variable mods, because it is possible to get mass shift at each end of a reversed peptide sequence that just happens to transform a genuine y series match into a false b series match or vice versa. (True palindromic peptides also exist, but are rare).

One advantage of using reversed entries is that entries with unusual amino acid compositions are preserved. For example, the very high levels of G and P in collagens or the runs of near poly-G in many keratins

Reversing the entries also preserves the degree of redundancy present in the target database, which means that the number of unique sequences will be similar. This will not be the case if the sequences are randomised.

# The Database

## Randomised database

- Option 1
    Scramble individual protein entries
- Option 2
    Replace each entry with a random sequence of the same length but with the average AA composition of the database as a whole
- Option 3
    Replace each entry with a randomised sequence that preserves the statistics of the original, e.g. same di-peptide and tri-peptide frequencies, etc.

**MASCOT** : *Target-Decoy Validation*   © 2008 Matrix Science   **MATRIX SCIENCE**

The other approach is to randomise the sequences, rather than simply reverse them. This still leaves us with a number of possibilities.:

With option 1, an entry that is Cys rich or Gly rich will remain so. This can create problems with entries which have long runs of poly R or K. When scrambled, these go from producing fewer tryptic peptides than you would expect for a protein of that size to producing far too many

Option 2 is our preferred approach and has been implemented as an automatic part of a Mascot search, as I will describe shortly.

Option 3 is unnecessary unless the search algorithm is known to take di-peptide and tri-peptide frequencies into account. This is not the case with Mascot and I suspect not the case with any other algorithm. The only time this might have a noticeable effect is with enzyme cleavage specificity. To get the same number of peptides with strict trypsin, you need to preserve the frequencies of KP and RP when randomising.

**The Database**

**Separate or concatenated?**
- Threshold score 30
- Match in target database 50
- Match in decoy database 40

**False positive?**
- Concatenated: No
- Separate: Yes

**MASCOT** : *Target-Decoy Validation*    © 2008 Matrix Science    *MATRIX SCIENCE*

The Gygi group advocate searching a database in which the real and decoy sequences have been concatenated. A more conservative approach is to search the two databases independently. If the Mascot score threshold for a given spectrum is (say) 30, and we get a match of 50 from the real database and 40 from the decoy database, this would not count as a false positive from a concatenated database, but it would count as a false positive if the two had been searched independently.

Note that, when you search a concatenated database, you must double the number of matches found in the decoy because a random match is equally likely to occur in the target half.

On our public web site there is a help page devoted to decoy database searches. You can download a utility program from this page that allows you to create a randomised or reversed database. You can use this if you have an old version of Mascot, or if you want to use a decoy with reversed entries.

MASCOT : *Target-Decoy Validation*  © 2008 Matrix Science

Because more and more people wish to perform decoy searches routinely, we added this into Mascot 2.2 as a built-in part of the search. If you choose the Decoy checkbox on the search form, every time a protein sequence from the target database is tested, a random sequence of the same length is automatically generated and tested. The average amino acid composition of each random sequence is the same as the average composition of the target database. The matches and scores for the random sequences are recorded separately in the result file. The net result is identical to searching a separate, randomised database.

When the search is complete, the statistics for matches to the random sequences, which are effectively sequences from a decoy database, are reported in the result header. If you change the significance threshold, the numbers are recalculated. For example, if we change the threshold from 5% to 0.5% …

The false discovery rate drops accordingly. Of course, so does the number of true positives
If you click the hyperlink, you display the results from searching the randomised database.

## Terminology

**BOX 1 TERMINOLOGY AND GENERAL STATISTICAL METHODS FOR CONTROLLING FALSE DISCOVERY RATE**

One can view spectral identification as a process of hypothesis testing in which the hypothesis $H_0$ 'random chance identification' is tested for each spectrum against the alternative hypothesis $H_a$ 'correct identification'.

**Table 2** summarizes the outcome of identification of $m$ MS/MS spectra. Counts $U$, $V$, $T$ and $S$ are unknown and random due to the stochastic nature of mass spectra. The total number of incorrectly identified spectra $m_0$ is unknown but fixed. Although $V$, the number of false positive identifications, is unknown, it is possible to estimate or bound various error rates that involve the expected value of $V$ (that is, the average value that one would obtain after an infinite repetition of the experiment):

• False positive rate (FPR), or type I error, is a property of a single MS/MS spectrum, and is defined as the probability that a randomly matched spectrum is judged correct: $FPR = E(V)/m_0$.

• Family-wise error rate (FWER) is a property of $m$ MS/MS spectra, and is defined as the probability of making at least one incorrect identification among all identifications judged correct: that is, $FWER = p(V \geq 1)$. Example of method controlling FWER: Bonferroni[102].

• False discovery rate (FDR) is a property of $m$ MS/MS spectra, and is defined as the expected proportion of incorrect identification among all identifications judged correct: that is, $FDR = E(V)/R$. Examples of methods controlling FDR: step-up[19], permutation-based[103], empirical Bayes[22,104].

**Table 2** | Outcomes of applying a classification rule

| | No. of matches judged incorrect | No. of matches judged correct | Total |
|---|---|---|---|
| Number of truly incorrect matches | $U$ | $V$ | $m_0$ |
| Number of truly correct matches | $T$ | $S$ | $m - m_0$ |
| Total | $m - R$ | $R$ | $m$ |

Nesvizhskii, A. I., *et al*., Nature Methods 4 787-797 (2007)

**MASCOT** : *Target-Decoy Validation*  © *2008 Matrix Science*  **MATRIX SCIENCE**

There is some confusion in the literature over terminology. Should we talk about false discovery rate or false positive rate? Some explanations are quite difficult to understand. This recent review gives a rather mathematical definition.

## Terminology

**The MS/MS spectrum comes from a peptide sequence in the database**

|  | | True | False |
|---|---|---|---|
| **Search reports a match to the correct sequence** | **True** | True positive | False positive |
| | **False** | False negative | True negative |

False Discovery Rate
= FP / (FP + TP)

True Positive Rate
= TP / (TP + FN)

False Positive Rate
= FP / (FP + TN)

**MASCOT** : *Target-Decoy Validation*     © 2008 Matrix Science     *MATRIX SCIENCE*

If TP is true positive matches and FP is false positive matches, the number of matches in the target database is TP + FP and the number of matches in the decoy database is FP. The quantity that is reported is the False Discovery Rate = FP / (FP + TP)

True Positive Rate and False Positive Rate are different quantities, and I'll return to this topic later.

14

## Terminology

**The FDR for proteins can be higher or lower than that for peptides, depending on the classification rule(s).**

**For example: You have a set of peptide matches with 5% FDR**

- If you report all proteins in which any of these peptides are found, protein FDR is likely to be higher than peptide FDR
- If you require a protein to have (say) two unique peptide matches, and group together proteins that contain the same set or a sub-set of peptide matches, protein FDR is likely to be lower than peptide FDR

**MASCOT** : *Target-Decoy Validation*       *© 2008 Matrix Science*       *MATRIX SCIENCE*

It is very important to distinguish between peptide and protein FDR. They are usually very different.

## Validation

| | | IPI human 2.18 | percent matched | decoy database | False discovery rate |
|---|---|---|---|---|---|
| **A8 Dataset** | no match | 722 | | 839 | |
| | score below threshold | 75937 | | 82125 | |
| | score above 5% identity threshold | 4307 | 5.2% | 23 | 0.5% |
| | score above 5% homology threshold | 6657 | 8.0% | 352 | 5.0% |

**MASCOT** : *Target-Decoy Validation*   *© 2008 Matrix Science*   *MATRIX SCIENCE*

Lets look at some typical numbers for a Mascot search of MudPIT data from a standard ion trap. The significance threshold was the default setting of 5%.

The first column of figures is for a search of IPI human. The third column of figures is for a search of the same data against the reversed database using identical search parameters.

You can see that the false discovery rate for the identity threshold is very conservative. A factor of 10 below the predicted rate. This is often the case for ion trap data, because the mass accuracy and signal to noise are limited.

The false discovery rate for the homology threshold is spot on. By using the homology threshold, we get a large number of additional true positives without exceeding our 5% false discovery rate.

Notice that only some 8% of the spectra give significant matches. This is not unusual. In fact, it is quite good. I would say 5% is average.

What about those false positives? Let's have a closer look.

This is the select report for the a8 search of the reversed database. Our highest scoring false positive has a score of 61

Monoisotopic mass of neutral peptide Mr(calc): 891.45
Fixed modifications: Carbamidomethyl (C)
Variable modifications:
M2       : Oxidation (M)
Ions Score: 61  Expect: 0.00084
Matches (**Bold Red**): 10/60 fragment ions using 11 most intense peaks

| # | b | b⁺⁺ | b* | b*⁺⁺ | b⁰ | b⁰⁺⁺ | Seq. | y | y⁺⁺ | y* | y*⁺⁺ | y⁰ | y⁰⁺⁺ | # |
|---|---|-----|----|------|----|------|------|---|-----|----|------|----|------|---|
| 1 | 115.05 | 58.03 | 98.02 | 49.52 | | | N | | | | | | | 7 |
| 2 | 262.09 | 131.55 | 245.06 | 123.03 | | | M | **778.41** | **389.71** | 761.39 | 381.20 | 760.40 | 380.70 | 6 |
| 3 | 391.13 | 196.07 | 374.10 | 187.55 | 373.12 | 187.06 | E | **631.38** | 316.19 | 614.35 | 307.68 | **613.37** | 307.19 | 5 |
| 4 | 504.21 | 252.61 | 487.19 | 244.10 | 486.20 | **243.60** | L | **502.33** | 251.67 | 485.31 | **243.16** | | | 4 |
| 5 | 632.31 | 316.66 | 615.28 | 308.14 | 614.30 | 307.65 | K | **389.25** | 195.13 | 372.22 | 186.62 | | | 3 |
| 6 | 746.35 | 373.68 | 729.32 | 365.17 | 728.34 | 364.67 | N | **261.16** | 131.08 | 244.13 | 122.57 | | | 2 |
| 7 | | | | | | | K | **147.11** | 74.06 | 130.09 | 65.55 | | | 1 |

And this is what it looks like. A near perfect match from the reversed database. Tryptic peptide, complete run of y ions, only one large peak left unmatched.

Asking whether it is correct or wrong becomes almost a philosophical question.

The fact is, when we search large numbers of spectra against large sequence databases, we can get such matches by chance. No amount of expert manual inspection will prevent this. Database matching is a statistical process and, for this search, the number and magnitude of the false positives is within the predicted range.

ROC plot – trypsin (IPI db)

**Tryptic search**

Sensitivity / 1-Specificity

- ■ Mascot Ion score (AUC=0.98)
- ■ PeptideProphet (AUC=0.96)
- ■ Sonar (AUC=0.94)
- ■ Tandem (AUC=0.93)
- ■ Spectrummill (tag) (AUC=0.91)
- ■ Sequest XCorr (AUC=0.91)
- ■ Spectrummill (AUC=0.86)

A

Kapp E. A., *et al.*, Proteomics (HUPO-PPP special issue), August 2005

**MASCOT** : *Target-Decoy Validation*    © 2008 Matrix Science    *MATRIX SCIENCE*

The performance of a scoring scheme is sometimes illustrated as a Receiver-Operating Characteristic or ROC Curve. Here is an example from the publication summarising the statistics for the data collected for the HUPO plasma proteome project.

**Receiver-Operating Characteristic (ROC Curve)**

MASCOT **: Target-Decoy Validation** © 2008 Matrix Science — MATRIX SCIENCE

A ROC Curve plots true positive rate and false positive rate as a function of a discriminator, such as a score threshold. A good scoring scheme will try to follow the axes, as illustrated by the red curve, pushing its way up into the top left corner. A useless scoring algorithm, that cannot distinguish correct and incorrect matches, would follow the yellow dashed diagonal line.

The origin of the ROC curve has unit specificity, i.e. zero false positives, but also zero true positives. Not a useful place to be. The top right of the ROC curve has unit sensitivity, i.e. 100% true positives, but also 100% false positives, which is equally useless. By setting a significance threshold in Mascot, you effectively choose where you want to be on the curve.

**Terminology**

The MS/MS spectrum comes from a peptide sequence in the database

|  | | True | False |
|---|---|---|---|
| Search reports a match to the correct sequence | **True** | True positive | False positive |
| | **False** | False negative | True negative |

False Discovery Rate
= FP / (FP + TP)

True Positive Rate
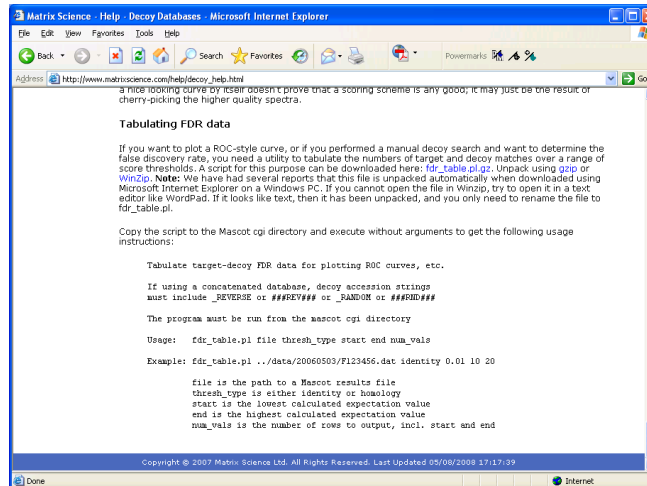= TP / (TP + FN)

False Positive Rate
= FP / (FP + TN)

**MASCOT** : *Target-Decoy Validation* © 2008 Matrix Science

You will remember seeing this slide earlier. To plot an authentic ROC curve, we need estimates of the numbers of true negatives (TN) and false negatives (FN), because true positive rate = TP / (TP + FN) and false positive rate = FP / (FP + TN). However, for real-life datasets, where we are dealing with unknown samples, we do not know TN and FN. So, what is presented as a ROC curve is often just a plot of the fraction of spectra matched in the target database versus the fraction matched in the decoy, or something similar.
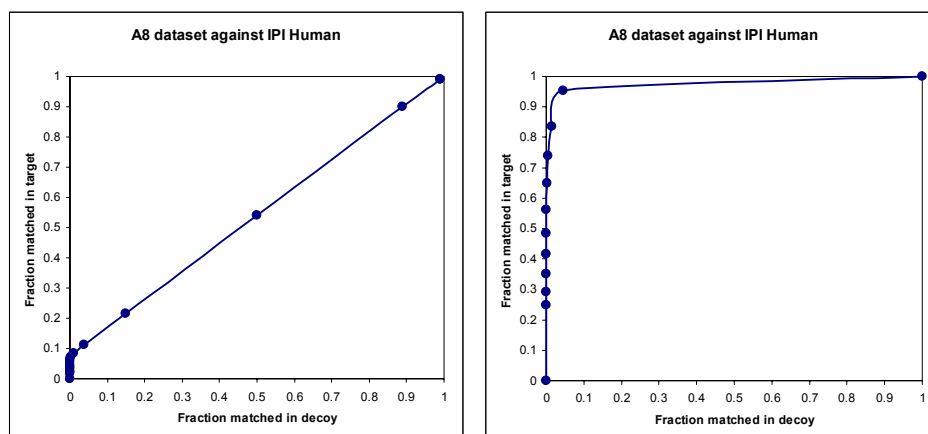
# Receiver-Operating Characteristic (ROC Curve)



If you want to plot a ROC-style curve, or if you performed a manual decoy search and want to determine the false discovery rate, you need a utility to tabulate the numbers of target and decoy matches over a range of score thresholds. A script for this purpose can be downloaded from the decoy help page on our web site

Receiver-Operating Characteristic (ROC Curve) — A8 dataset against IPI Human. MASCOT : Target-Decoy Validation © 2008 Matrix Science. MATRIX SCIENCE

If we take the MudPIT search I showed you a few slides back, and plot a ROC-style curve for the entire data set, you will get a very poor looking curve, like this one, because no score can discriminate the unmatchable spectra. In other words, as the score threshold is reduced towards zero, additional matches are equally likely to come from the decoy as from the target, and the ROC curve tends towards a diagonal line.

If you exclude the unmatchable spectra, then you can get a nice looking curve from exactly the same set of search results, like the one on the right. However, deciding which spectra to include is arbitrary. I could choose just the 10 strongest spectra, in which case any scoring scheme will give a beautiful ROC curve.

As long as all curves use the same dataset and assumptions, ROC curves are fine for comparisons. But, a nice looking ROC curve by itself doesn't necessarily prove that a scoring scheme is any good.