# Integration of spectral library searching into Mascot Server

MATRIX SCIENCE

In Mascot Server 2.6, we added the capability to search spectral libraries using MSPepSearch from Steven Stein's group at NIST. When submitting a search, any combination of amino acid Fasta files, nucleic acid Fasta files, and spectral library files can be selected. Here, we perform a simple search of some data from CPTAC study 6 against a NIST yeast library

Most search parameters – modifications, enzyme, missed cleavages, taxonomy, and instrument – simply don't apply to a library search. All that matters is how well the experimental spectrum matches the one in the library. The main exceptions are the precursor and fragment mass tolerances.

On completion of the search, the matches are reported in a protein family summary. In order to generate such a report, we need reliable and accurate protein inference

# Protein inference for library matches

- Library entries are peptides, not proteins, which means that protein information is only present as annotations
- Such annotations are optional, and may be missing
- Even when accessions are present
  - Reliability is unknown
  - Accession may not have any external meaning
  - Will rarely extend to more than a single accession per library entry

**MASCOT** : *Integration of spectral library …* © *2017 Matrix Science*   MATRIX SCIENCE

The are some difficulties associated with Protein inference for library matches. First of all, library entries are peptides, not proteins, which means that protein information is only ever present as annotations. Such annotations are optional, and may be missing, as in the case of most PRIDE libraries.

Even when present, the reliability is unknown. The accession could be a meaningless number or string. And, I've never seen a library with more than a single accession per library entry, so protein inference will be inaccurate for shared peptides.

Our solution is to require a reference Fasta database to be assigned to each library file when it is added to the system. The default is SwissProt, with an appropriate taxonomy filter, but any online Fasta can be chosen. This allows Mascot to map most of the library peptides to accessions in the reference database. This mapping is done at the sequence level, with no constraints from enzyme specificity. If a library entry has a novel sequence, not found in the reference database, the accession in the library annotations is used. If there is no accession, the peptide sequence is treated as the accession, so that duplicate matches to the same peptide can be grouped, if nothing else.

Here's the library search report again. Protein inference allows us to create a report for library matches that is near identical to a report for Fasta matches.

If only libraries are searched, MSPepSearch scores are converted to arbitrary expect values. A score of 300 becomes an expect value of 0.05 and the maximum score of 1000 becomes an expect of 5E-9

Coming back to the search form, you can search any mixture of amino acid and nucleic acid databases and spectral libraries.

This is a report for a search of both library and Fasta files. We can see all three databases listed at the top of the result report, and each is assigned an index so that we know where each accession comes from. The top hit has an index 3 which corresponds to the UniProt proteome. There are two important differences between this 'integrated' report and a library-only report.

## Integrated searches (Fasta + library)

- **Protein inference**
  - Library matches are mapped to accessions from the Fasta file
  - Reference database accessions or original library annotations only used where this fails

- **Library match scores**
  - Take the set of queries where the library and Fasta matches agree and the Mascot score is significant
  - Find scaling factors for library scores in this set such that their mean and standard deviation are the same as Mascot scores
  - Assign expect values based on the scaled scores, using the Mascot expect value formula

**MASCOT** : *Integration of spectral library …* © 2017 Matrix Science  MATRIX SCIENCE

For protein inference, if the peptide sequence can be mapped to one of the Fasta files being searched, this becomes the preferred accession. The accession from the reference database is only used when this fails.


In an integrated search, we can use the Fasta matches to create a simple empirical estimate of library score significance. This is achieved by calibrating library scores based on the set of queries where the library and Fasta searches return the same match and the Mascot score is significant. The shapes of the library and Mascot score distributions in this set are similar and they often have a fairly high correlation. Next, scale these library scores so that they have the same mean and standard deviation as Mascot scores. This produces values on the same scale as Mascot scores. We can now assign expect values to library matches using the same expression as for Mascot matches

Here, the top hit has been expanded. You can see that the top ranking PSMs come from both library and Fasta. In most cases, the same match is found in two or all three databases, and the listed match is the one with the lowest expect value. An exception can be seen here for query 8233. This peptide is non-specific at the amino terminus, and is only found in the library. It will not be matched in the Fasta files because the enzyme for the search was strict trypsin.

Let's turn our attention to administration aspects. Library files in MSP format are handled in Database Manager much the same as Sequence databases in Fasta format. This slide shows the top level screen of Database Manager, with a mixture of Fasta and library databases configured for searching. The 'Type' column shows which are AA or NA Fasta and which are spectral library. Most have 'predefined' configuration settings – that is, Matrix Science maintains a master file of configuration settings that is downloaded by Database Manager.

To enable a predefined library is a matter of a few mouse clicks

If the library you want to search is not on the predefined list, you use the 'Create New' Wizard to configure it as a custom database. A particularly interesting case is if you want to create your own library from Mascot search results. This is easily accomplished, as illustrated in the next few slides. Suppose that we are working on aardvark, and want to make a custom library for the aardvark proteome. We choose a name and select 'Create from search results'

The next screen just gives an opportunity to change the default location for the files

The reference database is used to assign protein accessions to the library entries. Normally, you wouldn't choose NCBIprot because it is such a large and redundant database. But, since SwissProt only contains 10 aardvark entries, we don't have much choice. We must also provide an estimate of suitable MS/MS tolerances for the library contents. If the search results come from multiple instruments, you need to base this on the least accurate of them.

Peptide match filters are used to select matches for inclusion in the library. We choose 'Edit filters'

**Peptide match filters for aardvark**

The library must have at least one score or expect value filter, typically expect < 0.01.

Each individual filter is in a filter group. To add more filters to the group, use the OR button. To add more groups, use the AND button. The peptide match must pass all filter groups to be accepted, but within each group, only one filter needs to succeed.

To remove a filter, leave its value field empty. To remove a filter group, remove all its filters.

Filters are used in two complementary ways:

1. When Database Manager chooses results files to process, only files that might contain suitable peptide matches are included.
2. When Database Manager loops over peptide matches in a results file, only matches that pass the filter are imported to the library.

For example, if you have a filter DB = SwissProt and no other DB filters, then only results files that were searched against SwissProt are processed. (Or in a multi-database search, had SwissProt as one of the databases.) When Database Manager loops over its peptide matches, only those that actually come from SwissProt are imported.

MASCOT : *Integration of spectral library …* © 2017 Matrix Science

There is a lot of flexibility here, and we don't have time to go into all the possibilities. This would be a simple filter for PSMs that can be assigned to a specific organism. We only want strong, confident matches in our library, so we require the match to have an expect value less than 0.01 and a score greater than 50. If the set of search results includes duplicate PSMs, only the one with the highest score goes into the library. We choose Save …

Which takes us back to the previous page, and we are ready to import search results

The only other thing we need to decide is which search result files to crawl. This can be specified as a date range or a wild card file path or some combination of the two. Finally, we add the import task to the queue and the selected files will be crawled as a background task.

You can even schedule automatic updates for such a database, which means that matches can be imported from new result files, created since the last import

## Cleaning out the Litterbox of Proteomic Scientists' Favorite Pet: Optimized Data Analysis Avoiding Trypsin Artifacts

Matthias Schittmayer,[†,‡,∥] Katarina Fritz,[†,‡,∥] Laura Liesinger,[†,‡] Johannes Griss,[§] and Ruth Birner-Gruenberger[*,†,‡]

[†]Research Unit Functional Proteomics and Metabolic Pathways, Institute of Pathology, Medical University of Graz, 8010 Graz, Austria

[‡]Omics Center Graz, BioTechMed-Graz, 8010 Graz, Austria

[§]Department of Dermatology, Medical University of Vienna, 1090 Vienna, Austria

**⑤** *Supporting Information*

**ABSTRACT:** Chemically modified trypsin is a standard reagent in proteomics experiments but is usually not considered in database searches. Modification of trypsin is supposed to protect the protease against autolysis and the resulting loss of activity. Here, we show that modified trypsin is still subject to self-digestion, and, as a result, modified trypsin-derived peptides are present in standard digests. We depict that these peptides commonly lead to false-positive assignments even if native trypsin is considered in the database. Moreover, we present an easily implementable method to include modified trypsin in the database search with a minimal increase in search time and search space while efficiently avoiding these false-positive hits.

**KEYWORDS:** *proteomics, autolysis protected trypsin, database search, search space restriction, misassigned spectra, false positives*

**MASCOT** : *Integration of spectral library …* © 2017 Matrix Science **MATRIX SCIENCE**

Let's look at a practical example of how these new features might be used. This recent paper JPR reminded us that sequencing grade trypsin is modified by methylation or acetylation of the lysines. Unless these variable modifications are selected in a search, simply including a contaminants database will not be sufficient to catch all trypsin autolysis peptides. The authors suggested a solution based on editing the sequence of trypsin in the Fasta, replacing K with J, and defining J as the mass of dimethylated lysine. This is fine, as far as it goes, but it misses many of the other modifications that are present, not to mention extensive non-specific cleavage.

- **Download data set from PRIDE**
- **Find "optimal" set of mods with error tolerant searches**
- **Search with these mods against SwissProt**

| | |
|---|---|
| Type of search | : MS/MS Ion Search |
| Enzyme | : semiTrypsin |
| Fixed modifications | : Carbamidomethyl (C) |
| Variable modifications | : Carbamidomethyl (N-term), Methyl (K), Methyl (N-term), Dimethyl (K), Dimethyl (N-term), Dehydro (C), Deamidated (NQ) |
| Mass values | : Monoisotopic |
| Protein mass | : Unrestricted |
| Peptide mass tolerance | : ± 10 ppm |
| Fragment mass tolerance | : ± 0.5 Da |
| Max missed cleavages | : 2 |
| Instrument type | : ESI-TRAP |
| Number of queries | : 26,505 |

- **Large search space, low sensitivity, but many matches to Trypsin**

| | | Score | Mass | Matches | Sequences |
|---|---|---|---|---|---|
| 1.1 | TRYP_PIG | 9334 | 25078 | 714 (714) | 59 (59) |
| 1.2 | TRY1_HUMAN | 247 | 27111 | 15 (15) | 4 (4) |
| 1.3 | TRY1_BOVIN | 100 | 26453 | 10 (10) | 3 (3) |
| 1.4 | TRY3_RAT | 75 | 26936 | 6 (6) | 3 (3) |

1 TRYP_PIG
4 TRY3_RAT
2 TRY1_HUMAN
3 TRY1_BOVIN

| Modification | Site | Above thr. |
|---|---|---|
| Carbamidomethyl | N-term | 539 |
| Deamidated | N | 403 |
| Carbamidomethyl | C | 357 |
| Dimethyl | K | 167 |
| Deamidated | Q | 142 |
| Methyl | N-term | 100 |
| Methyl | K | 71 |
| Dehydro | C | 70 |
| Dimethyl | N-term | 62 |

- **Import TRYP_PIG matches as new spectral library "Trypsin"**

**MASCOT** : *Integration of spectral library …* © 2017 Matrix Science

**MATRIX SCIENCE**

We downloaded the raw files for one of the data sets in this study from PRIDE and tried a variety of error tolerant searches to discover exactly what was present. Based on these results, we chose these search settings. The enzyme specificity was semiTrypsin because peptides show very extensive C-terminal 'ragged ends'

This makes the search space very large. The search takes a long time and overall sensitivity is not as good as it would be for a simple search with strict trypsin and only one or two variable modifications. The answer, of course, is to make a library of the trypsin matches and include this in the vanilla search. This is a very powerful option, since it allows any number of modified and non-specific peptides from any number of contaminants to be intercepted with no increase in the search space.

MASCOT Search Results

User        :
E-mail      :
Search title : Yeast_In-gel_digest
MS data file : C:\ProgramData\Matrix Science\Mascot Daemon\MGF\812 trypsin\mascot_daemon_merge.mgf
Databases  : 1: Trypsin (124 library entries)
             2: Yeast 20160706 (587,876 sequences; 285,263,038 residues)
Timestamp  : 24 Mar 2017 at 16:38:10 GMT

Search parameters
   Type of search        : MS/MS Ion Search
   Enzyme                : Trypsin/P
   Fixed modifications   : Carbamidomethyl (C)
   Variable modifications : Carbamidomethyl (N-term)
   Mass values           : Monoisotopic
   Protein mass          : Unrestricted
   Peptide mass tolerance : ± 10 ppm
   Fragment mass tolerance : ± 0.5 Da
   Max missed cleavages  : 2
   Instrument type       : ESI-TRAP
   Number of queries     : 26,505
Score distribution
Modification statistics

| Modification      | Site   | Above thr. |
| ----------------- | ------ | ---------- |
| Carbamidomethyl   | N-term | 420        |
| Carbamidomethyl   | C      | 312        |
| Dimethyl          | K      | 297        |
| Deamidated        | N      | 235        |
| Methyl            | K      | 64         |
| Deamidated        | Q      | 39         |
| Carbamidomethyl   | S      | 36         |
| Dimethyl          | N-term | 24         |
| Methyl            | S      | 4          |
| Methyl            | L      | 2          |
| Carbamidomethyl   | E      | 1          |

Here, we search a yeast database plus the tryptic autolysis library with strict trypsin and a single variable mod - yet still obtain matches to all the modified and non-specific trypsin autolysis peptides

This removes 1190 spectra which otherwise might have given rise to false positives.

If you're wondering about the ridiculous emPAI value, it's because the assumption behind emPAI is strict tryptic cleavage. However, the library search is giving all kinds of semitryptic matches, so the model assumptions are not satisfied.

# Summary

- Mascot Server 2.6 uses NIST MSPepSearch for spectral library searches
- You can search any combination of Fasta and spectral library files
- Results are presented using the protein family summary report
- A reference Fasta is assigned to each library file to ensure accurate protein inference
- For an integrated search, library match expect values are determined from the set of matches that have significant Mascot score and where the library and Fasta searches agree
- MSP files are configured and updated just like Fasta files
- Libraries can be created by importing Fasta search results

**MASCOT** : *Integration of spectral library ...* © 2017 Matrix Science    MATRIX SCIENCE

To summarise