# Spectral Library
# Searching

## Spectral libraries

- Spectral library contains annotated MS/MS spectra
- Match observed spectra directly to library spectra
- Advantages
  - Faster and more specific than database search
  - Easily search non-tryptic peptides or uncommon modifications
- Disadvantages
  - Only identifies peptides that exist in the library
  - Requires good measurement reproducibility
  - Creating high-quality libraries is time consuming

A spectral library is a collection of annotated MS/MS spectra of peptides. Instead of searching observed spectra against a protein sequence database, you search the observed spectra directly against a spectral library. The observed peaks are compared to the annotated library peaks, then scored in some way based on similarity. Typically, the similarity score takes advantage of peak intensity patterns as well as peak masses. It may also utilise peak annotations, such as giving a higher score to b and y ion matches.

There are several advantages compared to protein sequence databases. A library search is often much faster than a database search, as spectral library typically has orders of magnitude fewer peptides than a tryptic digest of a sequence database. A library search can also be more specific than a database search, for example if it contains previously identified non-tryptic peptides or uncommon variable modifications. Selecting a semi-specific enzyme or including many uncommon variable modifications in a database search greatly increases the search duration. Searching a pre-prepared library of spectra of semi-specific peptides is much faster.

There are no free lunches, so of course spectral library searching has weaknesses. You can only identify peptides for which a spectrum exists in the library. If the library contains a peptide sequence with one missed cleavage, you will not get a match to a peptide with two missed cleavages. A library search also requires decent measurement reproducibility. If peak intensities vary wildly between repeat runs, it is harder to get a good library match. Finally, creating high-quality libraries is time consuming and

some care is needed. Fortunately, Mascot ships with several predefined spectral libraries to help you get started.

**MASCOT MS/MS Ions Search**

| | | | |
|---|---|---|---|
| Your name | ShirleyJackson ⌄ | Email | |
| Search title | | | |
| Database(s) | NIST_S.cerevesiae_IonTrap (SL) | | Test_DNA<br>ZEST_human<br>*Spectral library (SL)*<br>MassIVE_HumanHCD<br>NIST_HSA_IonTrap<br>NIST_Human_HCD<br>NIST_Human_IonTrap<br>NIST_Rat_QTof<br>PiM_PiZ_SL<br>PRIDE_Contaminants |
| Taxonomy | All entries ⌄ | | |

| | | | |
|---|---|---|---|
| Peptide tol. ± | 10   ppm ⌄  # ¹³C 0 ⌄ | MS/MS tol. ± | 0.6   Da ⌄ |
| Peptide charge | 2+ ⌄ | Monoisotopic ● Average ○ | |
| Data file | Browse… No file selected. | | |
| Data format | Mascot generic ⌄ | Precursor | ___ m/z |
| Instrument | Default ⌄ | Error tolerant ☐ | |
| Decoy | ☐ | Target PSM FDR | (no target) ⌄ |
| | Start Search … | | Reset Form |

**MASCOT** : *Spectral library searching*  © 2017-2023 Matrix Science   MATRIX SCIENCE

Mascot Server can search spectral libraries using MSPepSearch from Steven Stein's group at NIST. When submitting a search, any combination of amino acid FASTA or nucleic acid FASTA databases, and spectral libraries 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

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.

## Protein inference for library matches

- A reference FASTA database must be specified for each library file
  - Library entries mapped to all proteins in reference that contain the sequence
- If library entry not found in reference database, the accession in the library annotations is used
- If no accession, the peptide sequence is used as the accession

MATRIX SCIENCE

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 database 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 database 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.

**MASCOT MS/MS Ions Search**

| | | | |
|---|---|---|---|
| **Your name** | ShirleyJackson ✓ | **Email** | |
| **Search title** | Yeast SL example | | |
| **Database(s)** | FungiDB_EST (NA) | | UniProt_Pig |
| | NIST_S.cerevesiae_IonTrap (SL) | > | Uniprot_Pseudomon_asaerugino |
| | UP2311_S_cerevisiae (AA) | < | Uniprot_Rice |
| | | | UniProt_Rust |
| | | | UP2494_R_norvegicus |
| | | | UP290289_Malus_domestica |
| | | | UP29965_C_sabaeus |
| | | | UP5226_T_rubripes |
| | | | UP5640_H_sapiens |
| | | | UP589_M_musculus |
| | | | UP625_E_coli_K12 |

| | | |
|---|---|---|
| **Peptide tol. ±** 10 ppm ✓ # $^{13}$C 0 ✓ | **MS/MS tol. ±** 0.6 Da ✓ | |
| **Peptide charge** 2+ | **Monoisotopic** ⦿ Average ○ | |
| **Data file** klc_031308p_cptac_study6_6B011.mgf | | |
| **Data format** Mascot generic | | |
| **Instrument** ESI-TRAP ✓ | **Error tolerant** ☐ | |
| **Decoy** ☐ | **Target PSM FDR** (no target) ✓ | |
| **Start Search ...** | **Reset Form** | |

**MASCOT** : *Spectral library searching*    © *2017-2023 Matrix Science*    MATRIX SCIENCE

Let's expand this example and search a mixture of amino acid and nucleic acid databases with a spectral library.

MASCOT : *Spectral library searching* © 2017-2023 Matrix Science

Here is the report from the search. 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 database + library)

- **Protein inference**
  - Library matches are mapped to accessions from the FASTA database
  - 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 database 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

For protein inference, if the peptide sequence can be mapped to one of the FASTA databases 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 database 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 database 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.

**MASCOT** : *Spectral library searching* © 2017-2023 Matrix Science

Here, the top hit has been expanded. You can see that the top ranking PSMs come from both library and FASTA database. 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 database because the enzyme for the search was strict trypsin.

## Databases and spectral libraries

| Name | Mode ? | Type ? | Status | | | Latest task |
|---|---|---|---|---|---|---|
| 3UTRrtFull_ncbi | custom | NA | In use | | Deactivate | |
| contaminants | predefined | AA | In use | Get new files | Deactivate | Update succeeded (view log) |
| cRAP | predefined | AA | In use | Get new files | Deactivate | Update succeeded (view log) |
| Mus_musculus_GRCm39_genomic | custom | NA | In use | | Deactivate | |
| NCBIprot | predefined | AA | In use | Get new files | Deactivate | Update succeeded (view log) |
| PRIDE_Contaminants | predefined | SL | In use | Get new files | Deactivate | Update succeeded (view log) |
| SARS-CoV-2 | predefined | AA | In use | Get new files | Deactivate | Update succeeded (view log) |
| SwissProt | predefined | AA | In use | Get new files | Deactivate | Update succeeded (view log) |
| TGAtoAGA_ReadthroughProtein2 | custom | AA | In use | | Deactivate | |
| UP5640_H_sapiens | predefined | AA | In use | Get new files | Deactivate | Update succeeded (view log) |
| UP589_M_musculus | predefined | AA | In use | Get new files | Deactivate | Update succeeded (view log) |
| UP625_E_coli_K12 | predefined | AA | In use | Get new files | Deactivate | Downloading (0.0%) |
| UP6548_A_thaliana | predefined | AA | In use | Get new files | Deactivate | Update succeeded (view log) |
| UP9136_B_taurus | predefined | AA | In use | Get new files | Deactivate | Update succeeded (view log) |
| Vertebrates_EST | predefined | NA | In use | Get new files | Deactivate | Update succeeded (view log) |
| Human_EST | predefined | NA | Offline | Get new files | Activate | Update succeeded (view log) |
| UniRef100 | predefined | AA | Offline | Get new files | Activate | Update succeeded (view log) |

Latest predefined definitions files are from Thu Aug 3 06:04:11 2023 (FASTA databases: databases_20230803T100411.xml) and Thu Oct 13 08:27:38 2022 (spectral libraries: libraries_20221013T122738.xml).

Full database status is available on the database status page.

Refresh

**Database Manager**
Databases (17)
Parse rules (16)
Scheduled updates (0)
Running tasks (1)
Settings

**Fasta**
Enable predefined definition
Synchronise custom definitions
Create new

**Library**
Enable predefined definition
Synchronise custom definitions
Create new
Spectral library filters

**MASCOT** : *Spectral library searching*  © 2017-2023 Matrix Science   MATRIX SCIENCE

Let's turn our attention to administration aspects. Library files in NIST 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 databases and libraries 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.

**Create spectral library from search results**

Library name:
*aardvark*

**Sequence directory**
*/opt/mascot-2.6-dev/sequence*

**Reference database**

Please choose a reference database. Where possible, protein accessions for peptides in the spectral library will be taken from the specified Fasta file (the reference database). This will make protein inference more reliable and allows a Protein View report to be displayed for a library hit.

[ NCBIprot ▼ ]

NCBIprot is larger than 5.0 GB. It is not recommended as a reference database.

**Taxonomy**

If the selected reference database has taxonomy configured, you can optionally choose a taxonomy for reference accessions.

[ . . . . . . . . . . . . . . . Aardvark ▼ ]

**MS/MS tolerance**

Please enter estimates for the absolute and relative tolerances of the fragment masses in the library. The tolerances in the Mascot search form apply to the data being searched. A library contains experimental spectra, also subject to mass measurement error. It is better to enter values that are too large rather than too small.

[ 0.1 ] **Da**

[ 100 ] **ppm**

[ Previous ] [ Create ]

Sidebar navigation:
- Database Manager
- Databases (9)
- Parse rules (18)
- Scheduled updates (0)
- Running tasks (0)
- Settings
- Fasta
- Enable predefined definition
- Synchronise custom definitions
- Create new
- Library
- Enable predefined definition
- Synchronise custom definitions
- Create new
- Spectral library filters

**MASCOT** : *Spectral library searching*   *© 2017-2023 Matrix Science*   MATRIX SCIENCE

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.

| Expect value | < | 0.01 | OR |
| AND |
| Score | > | 50 | OR |
| AND |
| Taxonomy | is / is not | . . . . . . . . . . . . . . . Aardvark | OR |
| AND |

Cancel   Test   Save

**MASCOT** : *Spectral library searching*   © 2017-2023 Matrix Science

**MATRIX SCIENCE**

There is a lot of flexibility here. 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.

**Journal of proteome research**

---

---

*Journal of* **proteome** .research

Article — pubs.acs.org/jpr

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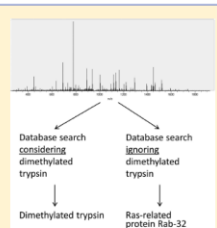
[‡]Omics Center Graz, BioTechMed-Graz, 8010 Graz, Austria

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

**S** 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** : *Spectral library searching*    © 2017-2023 Matrix Science    **MATRIX SCIENCE**

---

Let's look at a practical example of how these new features might be used. This JPR paper 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.
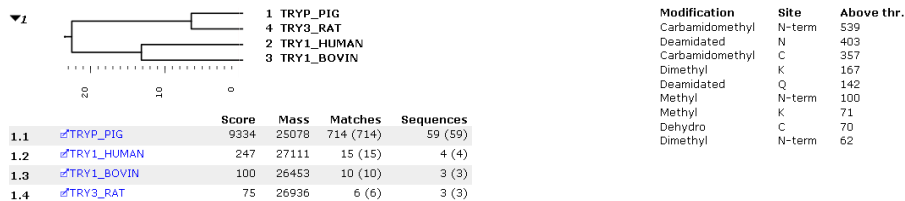
22

# Creating a trypsin library

- 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 |

**MASCOT**    **: *Spectral library searching***    *© 2017-2023 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'.

# Creating a trypsin library

- Large search space, low sensitivity, but many matches to Trypsin
- Import TRYP_PIG matches as new spectral library "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) |

| 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 |

**MASCOT** **: *Spectral library searching*** *© 2017-2023 Matrix Science*

**MATRIX SCIENCE**

This makes the search space very large, but we do get many matches to trypsin and many modified peptides. 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.

| Search title | : Benchmark small |
| MS data file | : C:\ProgramData\Matrix Science\Mascot Daemon\MGF\40 Benchmark small\mascot_daemon_merge.mgf |
| Databases | : **1:** SwissProt 2021_04 (565,928 sequences; 204,173,280 residues) |
| | **2:** Trypsin 20221216 (113 library entries) |
| Timestamp | : 16 Dec 2022 at 12:43:40 GMT |

Re-search   ● All ○ Non-significant ○ Unassigned   [help]   Export   As  XML

**▼Search parameters**

| Type of search | : MS/MS Ion Search |
| Enzyme | : Trypsin/P |
| Fixed modifications | : Carbamidomethyl (C) |
| Variable modifications | : Oxidation (M) |
| Mass values | : Monoisotopic |
| Protein mass | : Unrestricted |
| Peptide mass tolerance | : ± 20 ppm |
| Fragment mass tolerance | : ± 0.5 Da |
| Max missed cleavages | : 1 |
| Instrument type | : ESI-TRAP |
| Number of queries | : 99,299 |

▶ Score distribution

**▼Modification statistics for all protein families**

| Modification | Delta | Type | Site | Total matches |
|---|---|---|---|---|
| Carbamidomethyl | 57.021464 | fixed | C | 1357 |
| Oxidation | 15.994915 | variable | M | 505 |
| Deamidated | 0.984016 | SL | N | 127 |
| Dimethyl | 28.0313 | SL | K | 8 |
| Deamidated | 0.984016 | SL | Q | 5 |
| Methyl | 14.01565 | SL | K | 4 |
| Carbamidomethyl | 57.021464 | SL | E | 2 |
| Carbamidomethyl | 57.021464 | SL | N-term | 2 |
| Dehydrated | -18.010565 | SL | C | 2 |
| Carbamidomethyl | 57.021464 | SL | C | 1 |

**MASCOT** : *Spectral library searching*   *© 2017-2023 Matrix Science*   MATRIX SCIENCE

Here, we search Swissprot 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 776 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 uses NIST MSPepSearch for spectral library searches
- You can search any combination of FASTA databases and spectral libraries
- Results are presented using the protein family summary report
- A reference FASTA database 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 database searches agree
- MSP files are configured and updated just like FASTA databases
- Libraries can be created by importing results from searches against FASTA databases

To summarise.