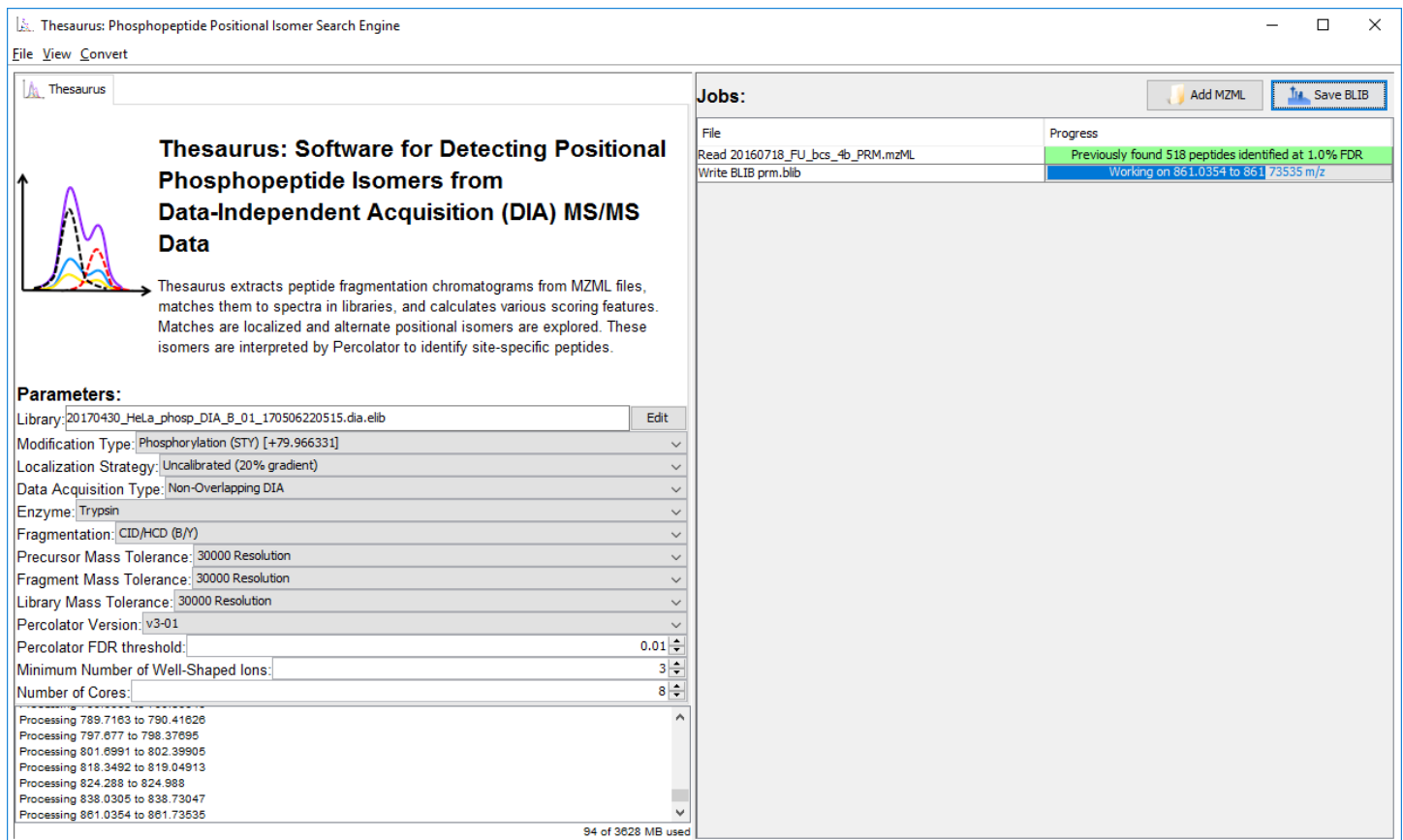# Thesaurus: quantifying phosphoprotein positional isomers

**Tutorial based on Thesaurus version 0.5.7 and Skyline Daily, built by Brian Searle (searleb@uw.edu)**

Thesaurus is a hybrid search engine that detects new positional isomers using site-specific fragment ions from parallel reaction monitoring and data independent acquisition mass spectrometry experiments. Typically mass spectrometers are tuned to make decisions on the fly based on precursor mass and to actively exclude isomers from being sampled again. Two recent publications (PMID: 28604659, 28661500) have demonstrated it is possible to determine the site of phosphorylation from DIA data, and Thesaurus extends this approach to actively identify and quantify additional positional isomers based on previously observed isomers in spectrum libraries. Thesaurus can measure distinct quantitative signaling effects of different positional isomers of the same phosphopeptides, even if those isomers do not separate chromatographically.



## PREREQUISITES

Thesaurus is a cross-platform Java application that has been tested for Windows, Macintosh, and Linux. The Thesaurus GUI can be opened by double clicking on the Thesaurus .JAR (e.g. thesaurus-0.5.4-executable.jar). Alternatively, Thesaurus can be used through the command line.

Thesaurus requires 64-bit Java 1.8. If you don't already have it, you can download "Windows x64 Offline" from: http://www.oracle.com/technetwork/java/javase/downloads/jre8-downloads-2133155.html

# THESAURUS QUICKSTART WALKTHROUGH

Here is the download link for the demo library:

https://bitbucket.org/searleb/thesaurus/downloads/20160718_FU_bcs_4b_PRM.mzML.dlib

And here is the download link for a zipped mzML from a PRM study using a Thermo Fusion:

https://bitbucket.org/searleb/thesaurus/downloads/20160718_FU_bcs_4b_PRM.mzML.zip


After you have downloaded the demo files, unzip the mzML and run Thesaurus.



Complete the following steps:

1. Load the library using the "Edit" button next to the pink Library option on the left parameter panel:



2. Set the rest of the settings to match the above screen shot. Notably, make sure the "Localization Strategy" is set to "Across entire window" and the "Data Acquisition Type" is set to "Non-Overlapping DIA".

3. Only after all the settings have been made, load the mzML by clicking on the "Add MZML" button in the upper right corner:

This will start a new processing job that will work through several algorithmic steps, including indexing the mzML, creating the null fragment distribution, processing each window individually, and then running percolator across the experiment. This may take several minutes to finish:



The job will turn green and indicate the number of detected peptides after it has completed:



In addition to several TXT and PDF reports, Thesaurus creates a report file (in this case "20160718_FU_bcs_4b_PRM.mzML.thesaurus.elib"), which is an open format SQLite database that contains all the localization and quantitation data. These reports will be saved in the same folder that contains the original mzML.

To visualize the report file, launch the browser by selecting the "View/Launch Thesaurus Browser" menu option. In the browser you must select both the report library file (XXX.thesaurus.elib), and the mzML raw file:



This will load a table with peptides that you can select. Select the first peptide will display:



Here the upper panel has the precursor chromatograms and the bottom panel has the fragmentation chromatograms.

This PRM experiment was mainly meant to confirm the localization of the peptide KGSGDYMPMSPK from IRS1. There is a search bar in the lower left. Type in "KGSGDYMPMSPK" in the bar to find that peptide and select it:



Here there are three significant localizations for KGSGDYMPMSPK: KG**Sp**GDYMPMSPK (RT=50.3 min), KGSGD**Yp**MPMSPK (RT=45 min), and KGSGDYMPM**Sp**PK (RT=49.5 min). You can view the site localizing ions (bold fragment chromatograms) and shared ions (dashed fragment chromatograms) for each of these forms by selecting the appropriate tab above the fragment chromatogram panel.

# THESAURUS PARAMETER SETTINGS

**Parameters:**

| | | |
|---|---|---|
| Library: | 20170430_HeLa_phosp_DIA_B_01_170506220515.dia.elib | Edit |
| Modification Type: | Phosphorylation (STY) [+79.966331] | ⌄ |
| Localization Strategy: | Uncalibrated (20% gradient) | ⌄ |
| Data Acquisition Type: | Non-Overlapping DIA | ⌄ |
| Enzyme: | Trypsin | ⌄ |
| Fragmentation: | CID/HCD (B/Y) | ⌄ |
| Precursor Mass Tolerance: | 30000 Resolution | ⌄ |
| Fragment Mass Tolerance: | 30000 Resolution | ⌄ |
| Library Mass Tolerance: | 30000 Resolution | ⌄ |
| Percolator Version: | v3-01 | ⌄ |
| Percolator FDR threshold: | 0.01 | |
| Minimum Number of Well-Shaped Ions: | 3 | |
| Number of Cores: | 8 | |

**Library:** You must load a .DLIB or .ELIB library before progressing.

**Modification Type:** Here you can specify a variety of modifications to localize.

**Localization Strategy:** For most DIA experiments, you should specify "Recalibrated (20% gradient)". For PRM experiments you should specify "Across entire window". For paired DDA/DIA experiments (e.g. SWATH) you should specify "Uncalibrated (20% gradient)". If you're only interested in localizing peaks detected via DDA or in a previous library search experiment you can specify "peak width only" searches.

**Data Acquisition Type**: You can specify "Overlapping DIA" or "Non-Overlapping DIA". In general, we recommend processing overlapping DIA experiments with the PRISM algorithm to deconvolute overlapping windows in MSConvert and specifying "Non-Overlapping DIA" in Thesaurus. However, it is possible to load overlapping window data directly and Thesaurus will use a heuristic-based deconvolution scheme to deconvolute the data.

**Enzyme**: Several digestion enzymes are supported.

**Fragmentation**: In general, we recommend using CID/HCD (B/Y) fragmentation for most CID or HCD experiments. However, if your library is particularly large or messy you may get improved results with "HCD (y-only)".

**Precursor/Fragment/Library Mass Tolerance:** Tolerances can be specified in PPM, AMU, or resolution.

**Percolator Version:** Percolator 3.1 is recommended for most experiments.

**Percolator FDR threshold**: You can specify the q-value for localization, e.g. 0.01 or 0.05. This value is used as an FDR threshold for detection of new forms and as a p-value threshold for localization thresholds.

**Minimum Number of Well-Shaped Ions**: This setting requires a minimum number of ions that fit the peak-shape of localizing ions. In general, we recommend setting this to 0.

**Number of Cores**: This is the number of CPU cores you allow Thesaurus to use.

## BUILDING A LIBRARY FROM SKYLINE/NIST

You can build .DLIB libraries using the "Convert" menu. Thesaurus supports BLIB from Skyline or MSP from NIST. All libraries require a FASTA database to link peptides to proteins. Skyline supports retention times and Thesaurus will use those if they are added to the BLIB. Additionally, large multi-file Skyline libraries can be retention time-aligned using iRT standards using an optional .IRTDB database.



## RUNNING THESAURUS

Once parameters are set up, you can run Thesaurus by adding an mzML using the button in the jobs panel.



You can queue up as many jobs as you like and the settings will be saved with the job. Jobs will be executed one at a time. Thesaurus will automatically generate several reports, including a .ELIB SQLite3 database for visualization. Alternatively, you can also execute the "Save BLIB" job to export a library for further visualization in Skyline.

## THESAURUS VISUALIZATION

Alternatively, you can visualize results directly in Thesaurus using the "View" menu.



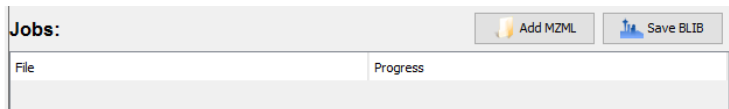First you must specify both a library result file (.ELIB) and a mzML file. This will populate a table on the left showing the detected positional isomers for each peptide, and a series of chromatogram charts on the right.

The table is sortable by clicking on the headers. The table headers are:

**Peptide sequence**: all forms of this peptide are displayed together

**Number of Mods**: the number of modifications in this peptide form

**Number of Forms**: the number of positional isomers detected for this peptide

**Number of Sites**: the number of potentially modifiable residues by the localization modification

**Delta RT (sec)**: if there are multiple forms, this is the difference in retention time between those forms in seconds

**Protein**: the protein accession numbers associated with this peptide.

Clicking on a table entry displays the chromatograms for that peptide.



The top panel is the precursor chromatogram, where monoisotopic ions are blue, M+1 ions are purple, and M+2 ions are red. If there are multiple positional isomers detected, they are indicated by dashed lines and listed as multiple tabs and bottom panel shows the fragment ion chromatogram corresponding to the selected tab:



Rounded parentheses () indicate where the modification is localized. If the form can only be partially localized, then the parentheses may extend to multiple residues. The localization score is also displayed at the tab header. Finally, in the fragment ion chromatograms site-specific ions are bold, while shared ions are dashed.

# THESAURUS COMMAND LINE

```
Command Prompt                                                    —   □   ×

D:\thesaurus>java -Xmx10g -jar thesaurus-0.5.4-executable.jar -h
[16:05:42] Thesaurus Help
Thesaurus is software for detecting positional phosphopeptide isomers from DIA data.
You should prefix your arguments with a high memory setting, e.g. "-Xmx8g" for 8gb
Required Parameters:
    -i      input .DIA or .MZML file
    -l      library .ELIB file
Other Parameters:
    -o      output report file (default: [input file].thesaurus.txt)
    -acquisition                       (default: Non-Overlapping DIA)
    -enzyme                            (default: trypsin)
    -expectedPeakWidth                 (default: 25)
    -fixed                             (default: C=57.0214635)
    -foffset                           (default: 0)
    -frag                              (default: CID)
    -ftol                              (default: 10)
    -ftolunits                         (default: ppm)
    -lftol                             (default: 10)
    -lftolunits                        (default: ppm)
    -localizationModification          (default: none)
    -minNumOfQuantitativePeaks         (default: 3)
    -numberOfExtraDecoyLibrariesSearche (default: 0.0)
    -numberOfQuantitativePeaks         (default: 5)
    -percolatorThreshold               (default: 0.01)
    -percolatorVersionNumber           (default: 3)
    -poffset                           (default: 0)
    -precursorIsolationMargin          (default: 0)
    -ptol                              (default: 10)
    -ptolunits                         (default: ppm)
    -scoringBreadthType                (default: window)

D:\thesaurus>_
```

Thesaurus can be executed from the command line. For example:

java –Xmx10g –jar thesaurus-0.5.4-executable.jar –h

Here, you must specify the maximum amount of memory available to Thesaurus (e.g. –Xmx10g for 10 GB of RAM). You should allocate at least 2 GB of RAM for most jobs, up to a maximum of your total physical memory minus 2 GB.

Thesaurus requires that you provide an input mzML file (-i) and an .ELIB or .DLIB library (-l). You can create those libraries from Skyline or NIST libraries using the GUI. You must also specify a modification type (-localizationModification), which can be "Phosphorylation", "Acetylation", "Oxidation", "Methylation", "Ubiquitination", and "OHexNAc". Finally, you must specify a scoring type (-scoringBreadthType), which can be "window" (for PRM), "recal" or "recal20" (DIA), "uncal" or "uncal20" (paired DDA/DIA). "Window" searches the entire acquisition window for positional isomers, while "recal" (search peak-width only) and "recal20" (search 20% of the gradient) perform an initial search to align each library entry to the DIA file. "Uncal" and "uncal20" expect the library retention times to be precisely aligned to the DIA file.

Data acquisition options include:

-acquisition (either overlappingDIA, DIA, or PRM)

-ftol/-ptol/-lftol (tolerances for fragments, precursors, or library fragments)

-ftolunits/-ptolunits/-ltolunits (units for mass accuracy, can be PPM, AMU, or RESOLUTION)

-foffset/-poffset (offset in PPM to combat poor instrument tuning)

-expectedPeakWidth (in seconds)

-precursorIsolationMargin (size of DIA precursor isolation margins in M/Z)

Search parameter options include:

-enzyme (trypsin, lys-c, lys-n, arg-c, cnbr, chymotrypsin, pepsin a, elastase, thermolysin, no enzyme)

-fixed (X=####, where X is an amino acid and #### is the added mass. Use commas for multiple mods)

-minNumOfQuantitativePeaks (the number of well-formed peaks required for localization)

-percolatorThreshold (q-value for localization, e.g. 0.01 or 0.05)

# THESAURUS SQLITE SCHEMA

The structure of Thesaurus' SQLite report library file is simple. There are two information tables and four data tables that are connected by three unique keys: PrecursorCharge, PeptideModSeq, and SourceFile. The primary table of interest is "peptidelocalizations", which contains quantification data (LocalizedIntensity and TotalIntensity) for all localized peptides, as well as localization data (IsLocalized, LocalizationScore, LocalizationIons, and LocalizationPeptideModSeq).

| entries | |
|---|---|
| **PrecursorMz** | DOUBLE |
| **PrecursorCharge** | INT |
| **PeptideModSeq** | STRING |
| **PeptideSeq** | STRING |
| **Copies** | INT |
| **RTInSeconds** | DOUBLE |
| **Score** | DOUBLE |
| **MassEncodedLength** | INT |
| **MassArray** | BLOB |
| **IntensityEncodedLength** | INT |
| **IntensityArray** | BLOB |
| CorrelationEncodedLength | INT |
| CorrelationArray | BLOB |
| RTInSecondsStart | DOUBLE |
| RTInSecondsStop | DOUBLE |
| MedianChromatogramEncodedLength | INT |
| MedianChromatogramArray | BLOB |
| **SourceFile** | STRING |

| peptidequants | |
|---|---|
| **PrecursorCharge** | INT |
| **PeptideModSeq** | STRING |
| **PeptideSeq** | STRING |
| **SourceFile** | STRING |
| **RTInSecondsCenter** | DOUBLE |
| **RTInSecondsStart** | DOUBLE |
| **RTInSecondsStop** | DOUBLE |
| **TotalIntensity** | DOUBLE |
| **NumberOfQuantIons** | INT |
| **QuantIonMassLength** | INT |
| **QuantIonMassArray** | BLOB |
| **BestFragmentCorrelation** | DOUBLE |
| **BestFragmentDeltaMassPPM** | DOUBLE |
| **MedianChromatogramEncodedLength** | INT |
| **MedianChromatogramArray** | BLOB |
| **IdentifiedTICRatio** | DOUBLE |

| peptidelocalizations | |
|---|---|
| **PrecursorCharge** | INT |
| **PeptideModSeq** | STRING |
| **PeptideSeq** | STRING |
| **SourceFile** | STRING |
| LocalizationPeptideModSeq | STRING |
| LocalizationScore | DOUBLE |
| LocalizationIons | STRING |
| NumberOfMods | INT |
| NumberOfModifiableResidues | INT |
| IsSiteSpecific | BOOLEAN |
| IsLocalized | BOOLEAN |
| RTInSecondsCenter | DOUBLE |
| LocalizedIntensity | DOUBLE |
| TotalIntensity | DOUBLE |

| peptidescores | |
|---|---|
| **PrecursorCharge** | INT |
| **PeptideModSeq** | STRING |
| **PeptideSeq** | STRING |
| **SourceFile** | STRING |
| QValue | DOUBLE |
| PosteriorErrorProbability | DOUBLE |
| IsDecoy | BOOLEAN |

| peptidetoprotein | |
|---|---|
| PeptideSeq | STRING |
| ProteinAccession | STRING |

| metadata | |
|---|---|
| **Key** | STRING |
| **Value** | STRING |

An example query to generate quantification reports would be:

```sql
select
    pep.PrecursorCharge, pep.PeptideModSeq, pep.PeptideSeq, pep.SourceFile,
    max(pep.LocalizedIntensity), max(pep.TotalIntensity), pep.IsSiteSpecific,
    pep.RTInSecondsCenter, pep.localizationScore,
    group_concat(p.ProteinAccession, ';') as ProteinAccessions
from
    peptidelocalizations pep
    left join peptidetoprotein p
where
    pep.PeptideSeq = p.PeptideSeq
group by pep.rowid;
```

Which will produce a table with the following columns:

| Column Name | Row 1 | Row 2 | Row 3 | ... |
|---|---|---|---|---|
| PrecursorCharge | 3 | 2 | 2 | ... |
| PeptideModSeq | VPQAEGPPKRVS[+79.966331] | KLS[+79.966331]SAMSAAK | S[+79.966331]ATFPNAGPR | ... |
| PeptideSeq | VPQAEGPPKRVSLVGADDLR | KLSSAMSAAK | SATFPNAGPR | ... |
| SourceFile | 20160718_FU_bcs_4b_PRM.n | 20160718_FU_bcs_4b_PRM.n | 20160718_FU_bcs_4b_PRM.n | ... |
| max(pep.LocalizedIntensity) | 7885.420898 | 31973.8125 | 34802.86328 | ... |
| max(pep.TotalIntensity) | 4290.157227 | 37123.98828 | 32597.57031 | ... |
| IsSiteSpecific | 1 | 0 | 0 | ... |
| RTInSecondsCenter | 2774.026855 | 2153.177002 | 2749.969971 | ... |
| LocalizationScore | 3.766996622 | 24.30858421 | 9.050885201 | ... |
| ProteinAccessions | sp\|Q14160\|SCRIB_HUMAN;sp | sp\|P40925-3\|MDHC_HUMAN | sp\|Q86X27-3\|RGPS2_HUMAN | ... |