



## FRAGARIYO v2.1

### Contents

Overview .....	2
References.....	2
I.    Choose where to save results.....	3
II.   PreRun Items .....	3
III.   Match terminal fragments first.....	5
Terminal Fragment Parameter File .....	6
Terminal Fragment Results.....	7
IV.   Internal Fragmentation.....	8
I.    Input Generator .....	8
Input ready for Internal Fragment matching:.....	8
II.   Running Internal Fragment Search.....	8
I.    Mass Resolution .....	8
II.   Tolerance Error .....	8
III.  Internal Run .....	8
Internal Fragment Results .....	11
V.    Data Analysis.....	13
Sequence Coverage Calculations .....	13
ClipsMS Plotter .....	13
Internal Fragments Outputs .....	13

## Overview

Post-translational modifications are important to understand the connection between protein structure and function. Mass spectrometry (MS)-based proteomics allowed for the high-throughput sequencing and quantitation for proteins in complex samples.<sup>1</sup> The typical bottom-up (BU) workflow requires digestion of proteins into peptides, which are analytes amenable to analysis by the instrumentation available in the early stages of MS. While BU proteomics is well established technique, the required digestion step precludes identification of all the protein forms (proteoforms)<sup>2</sup> that exist in the sample of interest. In top-down (TD) proteomics<sup>3</sup>, intact proteins are analyzed by LC-MS/MS optimized for large ions. Analyzing intact ions allows for the detection of proteoforms, however in typical TD proteomics denaturing conditions are used. Therefore, if there are structural differences between proteoforms that information is lost under denaturing conditions. Protein stoichiometry and sequence information can both be obtained from proteoforms of interest when using native TD (nTD) proteomics.<sup>4</sup> While improved instrumentation has been improving the accuracy and resolution of nTD data, software for the analysis of these datasets is lagging behind (plus a lack of standardization in file types and workflows to use for analysis).

Fragariyo is written in Python 3.7 using PyCharm IDE. It is a set of scripts capable of database search and multi-pass matching for terminals and internal fragments (taking into account modifications and disulfide bonds). Currently, Fragariyo accepts IMTBX\GrpPr .isotopes files (or Mobilatron), and CSV files (m/z, charge, and intensity) Let's hunt some fragments with Fragariyo!

## References

- (1) Aebersold, R.; Mann, M. Mass Spectrometry-Based Proteomics. *Nature* **2003**, *422*, 198–207.
- (2) Aebersold, R.; Agar, J. N.; Amster, I. J.; Baker, M. S.; Bertozzi, C. R.; Boja, E. S.; Costello, C. E.; Cravatt, B. F.; Fenselau, C.; Garcia, B. A.; Ge, Y.; Gunawardena, J.; Hendrickson, R. C.; Hergenrother, P. J.; Huber, C. G.; Ivanov, A. R.; Jensen, O. N.; Jewett, M. C.; Kelleher, N. L.; Kiessling, L. L.; Krogan, N. J.; Larsen, M. R.; Loo, J. A.; Ogorzalek Loo, R. R.; Lundberg, E.; Maccoss, M. J.; Mallick, P.; Mootha, V. K.; Mrksich, M.; Muir, T. W.; Patrie, S. M.; Pesavento, J. J.; Pitteri, S. J.; Rodriguez, H.; Saghatelian, A.; Sandoval, W.; Schlüter, H.; Sechi, S.; Slavoff, S. A.; Smith, L. M.; Snyder, M. P.; Thomas, P. M.; Uhlén, M.; Van Eyk, J. E.; Vidal, M.; Walt, D. R.; White, F. M.; Williams, E. R.; Wohlschlager, T.; Wysocki, V. H.; Yates, N. A.; Young, N. L.; Zhang, B. How Many Human Proteoforms Are There? *Nat Chem Biol* **2018**, *14*, 206–214.
- (3) Toby, T. K.; Fornelli, L.; Kelleher, N. L. Progress in Top-Down Proteomics and the Analysis of Proteoforms. *Annual Review of Analytical Chemistry* **2016**, *9*.
- (4) Zhou, M.; Lantz, C.; Brown, K. A.; Ge, Y.; Paša-Tolić, L.; Loo, J. A.; Lermyte, F. Higher-Order Structural Characterisation of Native Proteins and Complexes by Top-down Mass Spectrometry. *Chem Sci* **2020**, *11*, 12918–12936.

## I. Choose where to save results.

Simply press the “Change Output Directory button”. It can be changed at anytime except while Fragariyo is running processes.

**I. Current Output Directory**  

Please Select a directory to save results!

Change Output Directory

## II. PreRun Items

**II. PreRun Items**  

Modifications Library

IMTBX ouput extraction and renaming

- I. Modifications Library. Once clicked, user can upload a .txt file containing modification information (see ModificationsRepository.txt for a template).

ModificationsRepository.txt

	A	B	C	D	E	F	G	H
1	#Pyteomics ModX name	Printout Mod name	Target Residue	Mass Change	If fixed, which position?	Maximun number of modifications	If terminal, which one?	Chemical Composition
2	acetyl	acetyl	A	42.01056	1	1	N	C:2,H:2,O:1
3	trimethyl	trimethyl	K	42.04695	115	1		C:3,H:6,O:0

- A. Pyteomics ModX name: String with the modification name compliance with ModX

a. ModX = <https://pyteomics.readthedocs.io/en/latest/api/parser.html>

modX is a simple extension of the **IUPAC one-letter peptide sequence representation**.

The labels (or codes) for the 20 standard amino acids in modX are the same as in IUPAC nomenclature. A label for a modified amino acid has a general form of 'modX', i.e.:

- it starts with an arbitrary number of lower-case symbols or numbers (a modification);
- it ends with a single upper-case symbol (an amino acid residue).

The valid examples of modX amino acid labels are: 'G', 'pS', 'oxM'. This rule allows to combine read- and parseability.

- B. Printout Mod name: String with the modification name the user would like the results to print out as.
- C. Target Residue (*Multiple residues support incoming*)
- D. Mass Change (decimal number that can be either positive or negative).
- E. If modifications is fixed, which position
- F. Maximum number of Modifications per peptide
- G. If modifications is a terminal one, which terminus to place it at?
- H. Chemical Composition for scoring (only implemented in internal fragments searches, *terminal fragment support incoming*)

Element#1[colon]number of atoms of Element#1[comma] Element#2[colon]number of atoms of Element#2[comma]...etc...

DON'T ADD SPACES!

## II. IMTBX Output Extraction and Renaming

Script to rename files output from IMTBX/GrpPr. It takes .isotopes files from their original .raw file and renames them according to the .raw folder.

### III. Match terminal fragments first

- a. Choose Peaklist Files (File types acceptable are):
  - i. .isotopes files from IMTBX/Grppr for Waters and from Mobilitron for Agilent files
  - ii. mMass files (<http://www.mmass.org/>)
  - iii. csv files (mass-to-charge ratio, charge, intensity with headers m/z,z,int or sample name)
  - iv. unmatched files (produced from running terminal fragments. If ions are not matched, the unmatched ones will be placed in that file)
- b. Upload .ions (a previously produced theoretical ion database) files when prompted to, if a theoretical database has been created.
- c. Load the Parameter file (see TerminalFragments\_template.csv to learn how to fill a parameter file).
- d. Program will run (progress can be followed thru the pop-up terminal window).

```
Current Analysis = 1 of 3
PassName = CaM_xcz
~~~~~Fragmentation starts~~~~~
The disulfide list =
Protein Length: 148
Total prediction time: 1.8378021717071533
Total theoretical ions 5800
```

*Terminal Matching starts. Fragariyo GUI will freeze, but progress can be followed in the terminal window running with Fragariyo. Fragmentation will start. The analysis number will appear as well as the current pass (PassName). An analysis is a group of passes.*

Last update: October 18<sup>th</sup>, 2023

## Terminal Fragment Parameter File

**DON'T USE SPACES WHEN INPUTTING ITEMS SEPARATED BY SEMICOLONS “;”**

#Analysis	N	#PassTag (any type of text nu	protein seq	Fragmenta	ions_types	maxcl	Neutral Lo	Considerin	modificati	number of	Uniprot_o	Number of	Disulfides	naturally_reduce	c	Considerin	noncys_mods (separate with ";"	init_tol (p	final_tol (p	AUTO_CAL
1		NISTmAB_pGHC-xcz		QVTLRESG ECD	c;z	7	FALSE	FALSE								TRUE	pyrogluQ	150	5	TRUE
1		NISTmAB_pGHC-G2F-xcz		QVTLRESG ECD	c;z	7	FALSE	FALSE								TRUE	pyrogluQ;GiiF	150	5	TRUE
1		NISTmAB_pGHC-G2Fss0-xcz		QVTLRESG ECD	c;z	7	FALSE	TRUE	sshl;shl;chl	10	0	0	22-97;147	223;229;232		TRUE	pyrogluQ;GiiF	150	5	TRUE
1		NISTmAB_pGHC-G2Fss1-xcz		QVTLRESG ECD	c;z	7	FALSE	TRUE	sshl;shl;chl	10	0	1	22-97;147	223;229;232		TRUE	pyrogluQ;GiiF	150	5	TRUE

- A. Analysis number
- B. Name for the pass (it can be something informative about the pass e.g. ion type, mutant, etc)
- C. One-letter-code protein sequence
- D. Fragmentation type, it selects the correct c and z ion types to use. Currently CID and ECD are allowed. If another type is needed, fell free to request it.
- E. Ion types (do you need both c and z ions or only one?)
- F. Maximum charge state to be used
- G. Neutral losses bool. TRUE will considered neutral losses.
- H. Disulfide bond bool. TRUE makes the script to determine the amount of disulfide bonds and where they are and to separate them if the two cysteine partners end up in different fragments
- I. What modifications to use for the reduced cysteines
  - a. sshl;shl;chhsshl;hl;sh;h
- J. What is the maximum number of modifications to be used for the cysteine modifications (maximum number of cysteines that can be reduced)
- K. Uniprot\_offset, int. Number that the sequence index must be offset to match the disulfide bond positions , if these positions were obtain from the uniprot database. Put “0” is no offset needed
- L. Number of disulfide bonds to be broken even if their cysteines are in the same fragment
- M. Disulfide bond list
  - a. Cys1\_index-Cys1partner\_index; Cys2\_index-Cys2partner\_index (e.g NISTmAb disulfide bonds: 22-97;147-203;223-1000;229-1000;232-1000;264-324;370-428).
- N. Index for cysteines that do not participate in any disulfide bonds.
- O. Considering modifications (not involved in disulfide bonds)
- P. Modifications list (separate modifications with semi colons (e.g. NISTmAb modification for the n-term and glycan: pyrogluQ;GiiF)
  - a. See Modifications.py section for more information
- Q. Initial error tolerance
- R. Final error tolerance
- S. Auto Calibration bool.If TRUE, Auto calibration will be performed using

Last update: October 18<sup>th</sup>, 2023

## Terminal Fragment Results

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V
1	Protein se	DIQMTQSPSTLSASVGD	RVITITCSASSRVGYMHWYQKPGKAPKLLIYDTSKLASGVPSRFSGSGSGTEFTLTISSLQPD	FATYYCFQGS	GYPTFTGGG	KVEIKRTVAAPSVFIFPPSDEQLKSGTASV	VCLNNFYPREAKVQWKVDNALQSGNSQESVTEQDSK	DSTYLS	SSTLT													
2	Pass num	cal mz	mc mz	mono cal error(p	charge	ion type	mods	losses	neutral mz	free_disul	cys_locati	disulfide r	DT mono	(Ht (mono)	Area (mon	mz_top	DT (top)	Ht (top)	Area (top)	#Peaks in	top pk	inde charge
3	N	5	DIQMT																			
4	NISTmAB	606.2908	606.2916	0.598443	1	(c)5			605.2843	0	et(		3.19E-06			1		606.291	0.949329			
5	N	6	DIQMTQ																			
6	NISTmAB	734.3517	734.3502	-2.86253	1	(c)6			733.3429	0	et(		1.12E-05			1		734.352	-2.51165			
7	N	81	DIQMTQSPSTLSASVGD	RVITITCSASSRVGYMHWYQKPGKAPKLLIYDTSKLASGVPSRFSGSGSGTEFTLTISSLQPD																		
8	NISTmAB	1725.259	1725.259	-0.99642	5	(c-dot)81			8621.258	1	23 'hl';		1.24E-06			5		1725.26	-0.64554			
9	N	91	DIQMTQSPSTLSASVGD	RVITITCSASSRVGYMHWYQKPGKAPKLLIYDTSKLASGVPSRFSGSGSGTEFTLTISSLQPD	FATYYCFQGS																	
10	NISTmAB	2447.904	2447.937	4.712009	4	(c-dot)91			9787.72	0	23; 87		2.90E-06			4		2447.807	53.19886			
11	N	92	DIQMTQSPSTLSASVGD	RVITITCSASSRVGYMHWYQKPGKAPKLLIYDTSKLASGVPSRFSGSGSGTEFTLTISSLQPD	FATYYCFQGS																	
12	NISTmAB	3282.543	3282.588	4.732689	3	(c-dot)92			9844.741	0	23; 87		2.79E-06			3		3282.413	53.21954			
13	N	94	DIQMTQSPSTLSASVGD	RVITITCSASSRVGYMHWYQKPGKAPKLLIYDTSKLASGVPSRFSGSGSGTEFTLTISSLQPD	FATYYCFQGS	GYP																
14	NISTmAB	1685.189	1685.15	-4.58456	6	(c-dot)94			10104.86	0	23; 87		1.63E-06			6		1685.122	16.7159			
54	C	135	PDDFATYYCFQGS	GYPTFTGGG	KVEIKRTVAAPSVFIFPPSDEQLKSGTASV	VCLNNFYPREAKVQWKVDNALQSGNSQESVTEQDSK	DSTYLS	SSTLTLSKADY	EKHKVYACEVTHQGLSSPVT	KSFNRGEC												
55	NISTmAB	2419.715	2419.717	0.529907	6	(z)135			14512.26	4	193; 213; 'sshl'; 'chh'		2.55E-06			6		2419.716	0.43596			
56	C	127	CFQGS	GYPTFTGGG	KVEIKRTVAAPSVFIFPPSDEQLKSGTASV	VCLNNFYPREAKVQWKVDNALQSGNSQESVTEQDSK	DSTYLS	SSTLTLSKADY	EKHKVYACEVTHQGLSSPVT	KSFNRGEC												
57	NISTmAB	1935.13	1935.131	-0.4778	7	(z-dot)127			13538.86	4	193; 213; 'sshl'; 'chh'		2.50E-06			7		1935.131	-0.12691			
58	C	118	TFTGGG	KVEIKRTVAAPSVFIFPPSDEQLKSGTASV	VCLNNFYPREAKVQWKVDNALQSGNSQESVTEQDSK	DSTYLS	SSTLTLSKADY	EKHKVYACEVTHQGLSSPVT	KSFNRGEC													
59	NISTmAB	2540.281	2540.285	3.983459	5	(z)118			12696.39	3	193; 213; 'sshl'; 'chh'		2.66E-06			5		2540.282	1.156402			
60	C	117	TFTGGG	KVEIKRTVAAPSVFIFPPSDEQLKSGTASV	VCLNNFYPREAKVQWKVDNALQSGNSQESVTEQDSK	DSTYLS	SSTLTLSKADY	EKHKVYACEVTHQGLSSPVT	KSFNRGEC													
61	NISTmAB	1822.937	1822.904	0	7	(z)117			12753.28	1	193; 213; 'hl';		2.56E-06			7		1822.865	21.30047			
62	C	116	TFTGGG	KVEIKRTVAAPSVFIFPPSDEQLKSGTASV	VCLNNFYPREAKVQWKVDNALQSGNSQESVTEQDSK	DSTYLS	SSTLTLSKADY	EKHKVYACEVTHQGLSSPVT	KSFNRGEC													
63	NISTmAB	3128.587	3128.58	0.297673	4	(z)116			12510.29	3	193; 213; 'sshl'; 'sshl'		1.30E-06			4		3128.588	-2.52938			
64	NISTmAB	3128.587	3128.58	0.297274	4	(z)116			12510.29	3	193; 213; 'sshl'; 'sshl'		1.30E-06			4		3128.588	-2.52978			
65	C	110	EIKRTVAAPSVFIFPPSDEQLKSGTASV	VCLNNFYPREAKVQWKVDNALQSGNSQESVTEQDSK	DSTYLS	SSTLTLSKADY	EKHKVYACEVTHQGLSSPVT	KSFNRGEC														
66	NISTmAB	1995.008	1995.011	0.477797	6	(z-dot)110			11964.02	3	193; 213; 'sshl'; 'chh'		4.03E-06			6		1995.009	0.828682			
67	C	109	IKRTVAAPSVFIFPPSDEQLKSGTASV	VCLNNFYPREAKVQWKVDNALQSGNSQESVTEQDSK	DSTYLS	SSTLTLSKADY	EKHKVYACEVTHQGLSSPVT	KSFNRGEC														
68	NISTmAB	2009.644	2009.651	2.399016	6	(x)109			12051.86	3	193; 213; 'hl'; 'hl'; 'st'		2.67E-06			6		2009.645	2.749901			
69	C	105	VAAPSVFIFPPSDEQLKSGTASV	VCLNNFYPREAKVQWKVDNALQSGNSQESVTEQDSK	DSTYLS	SSTLTLSKADY	EKHKVYACEVTHQGLSSPVT	KSFNRGEC														
70	NISTmAB	1685.189	1685.15	-4.58456	6	(c-dot)94			10104.86	0	23; 87		1.63E-06			6		1685.122	16.7159			
NISTmAB_320V-9_hi-4_hits																						

The results are organized by termini and fragment site. The results are saved as.csv files and in binary .hits files (created with Pickle package).

## IV. Internal Fragmentation

### I. Input Generator

Internal fragment matches are scored by comparing experimental with theoretical isotope envelopes. Therefore, experimental isotope information is extracted from files containing the xy coordinates of the spectrum of interest. First a .csv file with experimental ions (mass-to-charge ratio, charge, intensity with headers m/z,z,int) is requested. Next Fragariyo will ask for the xy coordinates (currently only accepting .csv files from Agilent and xy files from Breuker. (select filetype in the lower right part of the pop-up window).

.cvs format example:

	A	B	C
1	m/z	z	int
2	202.1153	1	4.08628E-05
3	218.1355	1	2.00339E-05
4	239.0934	1	0.00012455
5	261.078	1	2.15709E-05
6	277.0535	1	2.95209E-05
7	298.148	1	0.071021523
8	303.1738	1	8.33549E-05
9	319.1933	1	4.79452E-05
10	373.209	1	8.42767E-05
11	381.1521	1	1.92355E-05
12	396.1439	1	4.75986E-05
13	403.1347	1	7.26961E-06
14	434.2228	1	4.86089E-05
15	450.2428	1	2.67029E-05
16	477.2106	1	7.10965E-06
17	478.2367	1	4.31501E-06

For IM-MS data obtained in Waters instrument use IMTBX-Grppr (<https://pubs.acs.org/doi/10.1021/acs.analchem.7b04999>) to produced .isotopes files and use the following columns in the .csv format:

.isotopes files > .csv  
mz\_mono > m/z  
charge > z(charge)  
ab\_top\_total > int (intensity )



**Input ready for Internal Fragment matching:**

	A	B	C	D	E	F	G	H
1	#neutral_r/z		mz	int	isoenv_mz	isoenv_int		
2	201.1072	1	202.115	1.05E-05	201.11529	2086;1499;1209;810;511;42		
3	217.1282	1	218.136	9.84E-06	217.14085	263;197;83;98;49;43;45;44;		
4	238.0852	1	239.093	2.31E-06	238.09680	0;0;5;0;1;0;0;10;0;0;21;50;2		
5	297.1402	1	298.148	4.38E-07	297.1521;	6;40;8;6;6;0;4;7;0;0;7;0;0;5		
6	302.1662	1	303.174	2.70E-05	302.17880	11612;8478;5518;4057;3014		

\_interexpion.csv file

- A. Neutral mass of experimental ion obtained from charge (z) and m/z values
- B. Charge
- C. Mass-to-charge ratio
- D. Intensity
- E. Values *m/z* corresponding to the peak isotope envelope
- F. Values *intensity* corresponding to the peak isotope envelope

## II. Running Internal Fragment Search

### II. Mass Resolution:

Approximate the mass resolution of your mass spectra (used to generate theoretical isotope envelopes)

### III. Tolerance Error:

What error to use (for internal fragments is recommended to internally calibrate data first to obtain the most accurate results). Typical, error used for internal fragments = 1 – 2- ppm.

### IV. Internal Run:

Once you click run, the program will request a batch file that contains the locations of other files. In order to set the internal run three files will be used: \_interexpion.csv (product of IV.I Input Generator), a parameter file to create the

Last update: October 18<sup>th</sup>, 2023

database (or a .theofrags file if the database has been created and saved previously), and the batch file. Below you will find a description of each file and what their columns mean. **DON'T USE SPACES WHEN INPUTTING ITEMS SEPARATED BY SEMICOLONS “;”**

Last update: October 18<sup>th</sup>, 2023

parameter file (for internal fragments, not the one use for terminals):

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q
1	#AnalysisN	#AnalysisPass	sequence	Fragmenta	maxcharge	minimum_f	maximum_g	noncys_m	mods_arra	number of	Considerin	Uniprot_o	Number of	disulfides (separate bonds with ; but don't end list in :)	naturally-reduced cysteines		
2	1	NISTmAb_HC_cz_z4	QVTLRESG ECD		4	80	150			10	FALSE	0	0				
3	1	NISTmAb_HC-G0F_cz_z4	QVTLRESG ECD		4	80	150	GoF		10	FALSE	0	0				
4	1	NISTmAb_HC-G0Fss0_cz_z4	QVTLRESG ECD		4	80	150	GoF	sshl;shl;chl	10	TRUE	0	0	22-97;147-203;264-324;370-429	223;229;232		
5	1	NISTmAb_HC-G0Fss1_cz_z4	QVTLRESG ECD		4	80	150	GoF	sshl;shl;chl	10	TRUE	0	1	22-97;147-203;264-324;370-430	223;229;233		
6	1	NISTmAb_HC-G1F_cz_z4	QVTLRESG ECD		4	80	150	GiF		10	FALSE	0	0				
7	1	NISTmAb_HC-G1Fss0_cz_z4	QVTLRESG ECD		4	80	150	GiF	sshl;shl;chl	10	TRUE	0	0	22-97;147-203;264-324;370-429	223;229;232		
8	1	NISTmAb_HC-G1Fss1_cz_z4	QVTLRESG ECD		4	80	150	GiF	sshl;shl;chl	10	TRUE	0	1	22-97;147-203;264-324;370-430	223;229;233		
9	1	NISTmAb_HC-G2F_cz_z4	QVTLRESG ECD		4	80	150	GiiF		10	FALSE	0	0				
10	1	NISTmAb_HC-G2Fss0_cz_z4	QVTLRESG ECD		4	80	150	GiiF	sshl;shl;chl	10	TRUE	0	0	22-97;147-203;264-324;370-429	223;229;232		
11	1	NISTmAb_HC-G2Fss1_cz_z4	QVTLRESG ECD		4	80	150	GiiF	sshl;shl;chl	10	TRUE	0	1	22-97;147-203;264-324;370-430	223;229;233		

- A. Analysis number
- B. Pass name (something descriptive about the pass (e.g ions searched for, mods searched for, etc))
- C. Sequence = Protein sequence one-letter code
- D. Fragment chemistry (it selects internal fragments allowed based on the chemistry selected). Currently CID and ECD are allowed. If another type is needed, fell free to request it.
- E. Maximum charge to be considered (right now max and min are the same to keep the passes somehow smaller. In a future update hopefully this will be simplified).
- F. Minimum fragment length in amount of residues
- G. Maximun fragment length in amount of residues
- H. Modifications for non-cysteiens residues separate by semicolons (examples shoes an analysis with several passes where a different NISTmAb glycan is searched for).
- I. What modifications to use for the reduced cysteines
- J. What is the maximum number of modifications to be used for the cysteine modifications (maximum number of cysteines that can be reduced)

Last update: October 18<sup>th</sup>, 2023

- L. Disulfide bond bool. TRUE makes the script to determine the amount of disulfide bonds and where they are and to separate them if the two cysteine partners end up in different fragments
- M. Uniprot\_offset, int. Number that the sequence index must be offset to match the disulfide bond positions , if these positions were obtain from the uniprot database. Put "0" is no offset needed
- N. Number of disulfide bonds to be broken even if their cysteines are in the same fragment
- O. Disulfide bond list
  - a. Cys1\_index-Cys1partner\_index; Cys2\_index-Cys2partner\_index (e.g NISTmAb disulfide bonds: 22-97;147-203;223-1000;229-1000;232-1000;264-324;370-428).
- P. Index for cysteines that do not participate in any disulfide bonds.

### Batch file:

You will be prompted to upload a batch file to run several internal fragment searches automatically or to just run one search (it helps to keep track if the right database file got run with the right ions).

	A	B	C	D
1	C:\Users\caror\Dropbox (University of Michigan)\Doctoral\CIU-ECD Manuscripts\CIU-ECD_methods_manuscript\NISTmAb\Peak Matching\2023_analysis\internals_inputandtheo			
2	#Expionsfile	Ions File	Param File	
3	NISTmAb_250V-45m-3_hi-4_hit_Unmatched_interexpion.csv		NISTmAb_HC_z4.csv	
4	NISTmAb_250V-45m-3_hi-4_hit_Unmatched_interexpion.csv		NISTmAb_HC_z5.csv	
5	NISTmAb_250V-45m-3_hi-4_hit_Unmatched_interexpion.csv		NISTmAb_HC_z6.csv	
6	NISTmAb_320V-3_hi-4_hit_Unmatched_interexpion.csv		NISTmAb_HC_z4.theofrags	
7	NISTmAb_320V-3_hi-4_hit_Unmatched_interexpion.csv		NISTmAb_HC_z5.theofrags	
8	NISTmAb_320V-3_hi-4_hit_Unmatched_interexpion.csv		NISTmAb_HC_z6.theofrags	

- A1: Path to the folder where all the files o be used are found.
- Row 2: Headers
- Column A: Place \_interexpion.csv file name, including extension as seen in image above.
- Column B: Name of the .theofrags file (if exists or it it was created to a previous search in the same batch).
- Column C: Name of the Parameter file to be used for the \_interexpion.csv file in the same row.

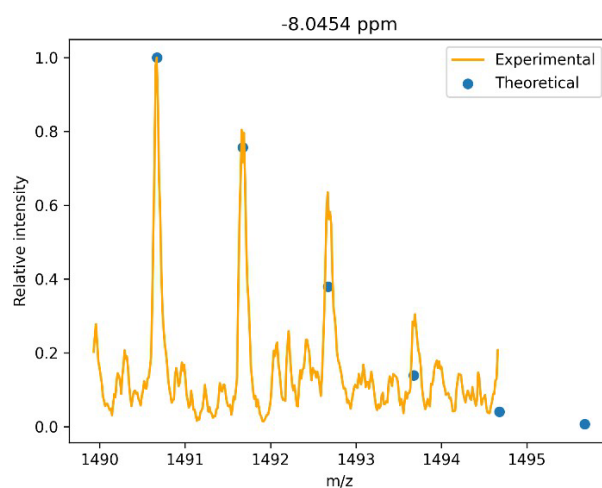
Script will run ( a protein a large as BSA took several days and breaking the analysis down to several pieces)

Last update: October 18<sup>th</sup>, 2023

### Internal Fragment Results

The results will be saved in a folder created with the name of the experimental ion .csv file: The folder contains:

1. A .theofrags file (a.k.a. .ions file) a binary version of the theoretical database
2. A.txt.frags. A human readable version of the theoretical databased (they can be quite large). To be removed.
3. PNG files showing the matching of the theoretical isotopic envelope and experimental isotopic envelope



#### 4. TSV files (tab-delimited files) containing the details of the matched peaks

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U
1	neutral ex	neutral the seq		charge	mz_mono	mods	ion_type	cysteine lc	ss_count	cysteines	cysteine m	indexstart	indexend	reverse_bc	cyclic dens	error	chemical_	isomz_sco	isoint_sco	fragment_score	
2	1489.661	1489.673	GEKLTDEE	1	1490.68	[]	c-z	0	0	0	None	113	126	FALSE	0	-8.04541	{'C': 62, 'H'	100	25	62.5	
3	1758.8	1758.798	LTDEEVDE	1	1759.805	[]	c-zdot	0	0	0	None	116	131	FALSE	0	1.159883	{'C': 73, 'H'	100	25	62.5	

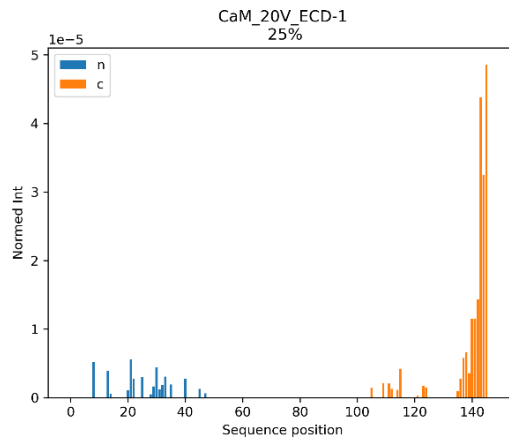
- A. Neutral mass (experimental ion)
- B. Neutral mass (theoretical ion)
- C. Fragment sequence
- D. Charge
- E. Monoisotopic  $m/z$  value
- F. Modifications
- G. Ion type
- H. Indices where the cysteines are located
- I. Number of disulfide bonds in internal fragment
- J. How many cysteines are being modified
- K. Cystine modification if disulfide bonds were reduced
- L. Index where the fragment starts
- M. Index where the fragment ends
- N. Bool. If TRUE the sequence is a reverse sequence. To be used as FDR (False Discovery Rate) calculations in the future.
- O. Cyclic density (number of cyclic regions form by disulfide bonds/ total residue number)
- P. Matching error
- Q. Chemical compositions (in python dictionary form = {'C': 44, 'H': 72, 'N': 14, 'O': 18, 'S': 0, 'Fe': 0})
- R. Scoring in  $m/z$
- S. Scoring in *intensity*
- T. Compound score

Last update: October 18<sup>th</sup>, 2023

## V. Data Analysis

They are divided into terminal and internal sections

### Sequence Coverage Calculations

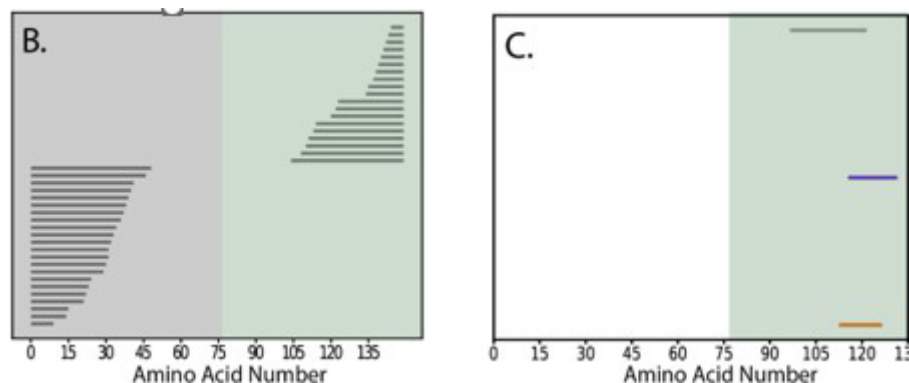


- Outputs example: A graph showing sequence coverage (coordinate information to replicate this graph is saved as a file: samplename\_analysis.csv) It can be done per file or for a combination of two files.

### ClipsMS Plotter

Based on ClipsMS (Software from the Loo lab) graphing code. It will transform terminal fragment input to be compatible with ClipsMS plotter. There is a terminal and an internal fragments version.

- ❖ Output: B. Terminal Fragments C. Internal Fragments (color has been added using Illustrator).



### Internal Fragments Outputs

In the Internal fragments analysis tools there is the option of merging the output into one file. The other option is to “average” them: Only retain fragments that appear across all files.

Last update: October 18<sup>th</sup>, 2023