

# **Spectrum Mill - Overview**

**Karl Clauser**

**Broad Institute of MIT and Harvard**

**August 2021**

# SM Core Capabilities

Spectrum Mill provides modules for complete analysis of proteomic LC-MS/MS data, including:

- Spectral pre-processing
  - directly from instrument vendor format files: Thermo Fisher, Agilent
- Identification of peptide spectrum matches - PSMs by database search.
  - PSM scoring capitalizes on low ppm product ion mass accuracies with fragment-ion type models optimized for orbitrap HCD, qToF-CID, ion trap CID, and ETD / ETHCD dissociation methods
- PTM-site localization
  - for phosphorylation, acetylation, and ubiquitylation with multi-sample site-level consolidation that handles localization ambiguity
- TMT-10, 11, 16, 18 / iTRAQ 4 reporter ion based quantitation
  - with PSM-level isotope correction and precursor ion purity filtering, and sample-level normalization
- Protein-level assembly of PSMs
  - with attention to isoform-specific and species-specific peptides
- Quality metric calculation/reporting
  - for MS acquisition, chromatography, spectral interpretation and FDR, digestion efficiency, PTM enrichment, and label incorporation
- Automation on an SM server through a service request manager
  - maximizes use of all CPUs, queues and processes tasks for multiple users via workflows with stored parameters

# Topics Covered

- Spectrum Mill architecture – program, files, workflow
- Configuring Extraction, MS/MS search, modifications
- PSM Scoring
- Phosphosite Localization
- Methods to control false discovery rates (FDR) of MS/MS search results
- Collapsing Peptide Spectrum Matches (PSM's) for Quantitation
  - Protein, Phosphosite levels
- Autovalidation
  - False Discovery Rates (FDR) at the PSM, peptide, and protein levels.
- Quantitation features
  - Reporter Ion Correction Factors
- Automation – SRM, workflows
- Quality Metrics
- Process Report
- Future Directions



# SM servers at Broad Institute

Spectrum Mill MS Proteomics Software Rev BI.07.08.214

**Process Automation Tools**

- Workflows** Automate data extraction, search, validation, and summary.
- Status of all SM servers** Request Queue and Completion Log
- Request Queue** View status of request queue for submitted requests, monitor results
- Completion Log** View log of completed extractions, searches, validations, and summ

**Mass Spectral Interpretation Tools**

- Data Extractor** Prepare MS/MS data files for Spectrum Mill processing
- MS/MS Search** Search database with MS/MS spectra
- Autovalidation** Autovalidate MS/MS spectra, calculate false-positive rates
- Spectrum Matcher** Match MS/MS spectra against other MS/MS spectra
- de novo Sequencing** Perform Sherenga de novo MS/MS spectral interpretation
- Manual PMF Search** Search database with MS spectra from Peptide Mass Fingerprint da

**Result Summary Tools**

- Protein/Peptide Summary** Summarize results from MS/MS searches
- Process Report** Run R scripts to parse, plot distributions, and normalize reporter-ion
- Quality Metrics & FDR** Summarize metrics for quality (MS performance, MS/MS interpretati peptide, protein)
- Spectrum Summary** Summarize characteristics of MS/MS spectra

**Utilities**

- Tool Belt** Collection of Spectrum Mill utility tools
- Protein Databases** Manipulate FASTA sequence databases for use with Spectrum Mill s
- Archive Data** Zip and Unzip dataset directories
- Peptide Selector** Select peptides from protein digest likely to produce high quality MS
- Ion Assignments** Generate tab-delimited files of fragment ion assignments for MS/MS
- MRM Selector** Build/Export MRM transition lists based on experimental data in Spe
- Multiple Sequence Aligner** Run Clustal W and highlight non-identical AA's across the alignment
- Peptide List to Masses** Convert a list of peptides to masses and formulas
- MS Product** Predict product ion masses from peptide sequence
- MS Digest** Predict peptide masses for enzymatic digestion of protein
- Peptide String Match** Search database with a list of peptide sequences

**Manuals (HTML):**

- Spectrum Mill Basics
- MS/MS Search
- Personalized Sequence DBs
- Workflow Automation
- PMF Search
- Protein Databases
- Server Administration
- Tips and Tricks
- FAQs
- What's New
- Help Index

**Guides:**

- Installation (.PDF)
- Getting Started (.PDF)
- Custom Modifications (.PDF)
- Report Column Descriptions (.xlsx)

**Slides (.pptx):**

- Spectrum Mill - Overview

**Useful Tables:**

- Mutation Mass Shifts
- Dipeptide Masses

Spectrum Mill MS Proteomics Software Rev BI.07.08.214

**Process Automation Tools**

- Workflows** Automate data extraction, search, validation, and summary.
- Status of all SM servers** Request Queue and Completion Log
- Request Queue** View status of request queue for submitted requests, monitor results
- Completion Log** View log of completed extractions, searches, validations, and summ

**Mass Spectral Interpretation Tools**

- Data Extractor** Prepare MS/MS data files for Spectrum Mill processing
- MS/MS Search** Search database with MS/MS spectra
- Autovalidation** Autovalidate MS/MS spectra, calculate false-positive rates
- Spectrum Matcher** Match MS/MS spectra against other MS/MS spectra
- de novo Sequencing** Perform Sherenga de novo MS/MS spectral interpretation
- Manual PMF Search** Search database with MS spectra from Peptide Mass Fingerprint da

**Result Summary Tools**

- Protein/Peptide Summary** Summarize results from MS/MS searches
- Process Report** Run R scripts to parse, plot distributions, and normalize reporter-ion
- Quality Metrics & FDR** Summarize metrics for quality (MS performance, MS/MS interpretati peptide, protein)
- Spectrum Summary** Summarize characteristics of MS/MS spectra

**Utilities**

- Tool Belt** Collection of Spectrum Mill utility tools
- Protein Databases** Manipulate FASTA sequence databases for use with Spectrum Mill s
- Archive Data** Zip and Unzip dataset directories
- Peptide Selector** Select peptides from protein digest likely to produce high quality MS
- Ion Assignments** Generate tab-delimited files of fragment ion assignments for MS/MS
- MRM Selector** Build/Export MRM transition lists based on experimental data in Spe
- Multiple Sequence Aligner** Run Clustal W and highlight non-identical AA's across the alignment
- Peptide List to Masses** Convert a list of peptides to masses and formulas
- MS Product** Predict product ion masses from peptide sequence
- MS Digest** Predict peptide masses for enzymatic digestion of protein
- Peptide String Match** Search database with a list of peptide sequences

**Manuals (HTML):**

- Spectrum Mill Basics
- MS/MS Search
- Personalized Sequence DBs
- Workflow Automation
- PMF Search
- Protein Databases
- Server Administration
- Tips and Tricks
- FAQs
- What's New
- Help Index

**Guides:**

- Installation (.PDF)
- Getting Started (.PDF)
- Custom Modifications (.PDF)
- Report Column Descriptions (.xlsx)

**Slides (.pptx):**

- Spectrum Mill - Overview

**Useful Tables:**

- Mutation Mass Shifts
- Dipeptide Masses

Spectrum Mill MS Proteomics Software Rev BI.07.08.214

**Process Automation Tools**

- Workflows** Automate data extraction, search, validation, and summary.
- Status of all SM servers** Request Queue and Completion Log
- Request Queue** View status of request queue for submitted requests, monitor results
- Completion Log** View log of completed extractions, searches, validations, and summ

**Mass Spectral Interpretation Tools**

- Data Extractor** Prepare MS/MS data files for Spectrum Mill processing
- MS/MS Search** Search database with MS/MS spectra
- Autovalidation** Autovalidate MS/MS spectra, calculate false-positive rates
- Spectrum Matcher** Match MS/MS spectra against other MS/MS spectra
- de novo Sequencing** Perform Sherenga de novo MS/MS spectral interpretation
- Manual PMF Search** Search database with MS spectra from Peptide Mass Fingerprint da

**Result Summary Tools**

- Protein/Peptide Summary** Summarize results from MS/MS searches
- Process Report** Run R scripts to parse, plot distributions, and normalize reporter-ion
- Quality Metrics & FDR** Summarize metrics for quality (MS performance, MS/MS interpretati peptide, protein)
- Spectrum Summary** Summarize characteristics of MS/MS spectra

**Utilities**

- Tool Belt** Collection of Spectrum Mill utility tools
- Protein Databases** Manipulate FASTA sequence databases for use with Spectrum Mill s
- Archive Data** Zip and Unzip dataset directories
- Peptide Selector** Select peptides from protein digest likely to produce high quality MS
- Ion Assignments** Generate tab-delimited files of fragment ion assignments for MS/MS
- MRM Selector** Build/Export MRM transition lists based on experimental data in Spe
- Multiple Sequence Aligner** Run Clustal W and highlight non-identical AA's across the alignment
- Peptide List to Masses** Convert a list of peptides to masses and formulas
- MS Product** Predict product ion masses from peptide sequence
- MS Digest** Predict peptide masses for enzymatic digestion of protein
- Peptide String Match** Search database with a list of peptide sequences

**Manuals (HTML):**

- Spectrum Mill Basics
- MS/MS Search
- Personalized Sequence DBs
- Workflow Automation
- PMF Search
- Protein Databases
- Server Administration
- Tips and Tricks
- FAQs
- What's New
- Help Index

**Guides:**

- Installation (.PDF)
- Getting Started (.PDF)
- Custom Modifications (.PDF)
- Report Column Descriptions (.xlsx)

**Slides (.pptx):**

- Spectrum Mill - Overview

**Useful Tables:**

- Mutation Mass Shifts
- Dipeptide Masses

## Cibola

24 CPUs 2.0 GHz  
56 GB RAM

## Manzano

24 CPUs 2.0 GHz  
56 GB RAM

## Eldorado

36 CPUs 2.0 GHz  
24 GB RAM

Jurkat QC

HLA - Immunopeptidomics



# Status of all SM servers at Broad Institute

**Spectrum Mill - Request Queue cibola**

[Request Queue](#) [Completion Log](#) [Help](#)

There are 4 requests in the queue. Available memory: 34.0 Gb of 51.5 Gb

[Remove](#)

X #	Task Id	Status	Task Type	Data Directory
<input type="checkbox"/> 1	210723085635.25437	Queued	Autovalidation	CPTAC3/LUAD/Proteome/25CPTAC_LUAD_Proteome_
<input type="checkbox"/> 2	210723085636.25438	Queued	Autovalidation	CPTAC3/LUAD/Proteome/25CPTAC_LUAD_Proteome_
<input type="checkbox"/> 3	210723085638.25439	Queued	Autovalidation	CPTAC3/LUAD/Proteome/25CPTAC_LUAD_Proteome_
<input type="checkbox"/> 4	<a href="#">210723085634.25409P Monitor</a>	Running (535:1924)	MS/MS Search	CPTAC3/LUAD/Proteome/25CPTAC_LUAD_Proteome_

**Spectrum Mill - Request Queue manzano**

[Request Queue](#) [Completion Log](#) [Help](#)

There are 7 requests in the queue. Available memory: 26.5 Gb of 51.5 Gb

[Remove](#)

X #	Task Id	Status	Task Type	Data Directory
<input type="checkbox"/> 1	210722213857.18237	Queued	Autovalidation	CPTAC3/PDAC/Phosphoproteome/10CPTAC_PDA_
<input type="checkbox"/> 2	210722213859.18238	Queued	Autovalidation	CPTAC3/PDAC/Phosphoproteome/10CPTAC_PDA_
<input type="checkbox"/> 3	<a href="#">210722213856.18235P Monitor</a>	Running (986:998)	MS/MS Search	CPTAC3/PDAC/Phosphoproteome/10CPTAC_PDA_
<input type="checkbox"/> 4	210723084219.19241	Queued	Autovalidation	CPTAC3/LUAD/Proteome/01CPTAC_LUAD_Proteom
<input type="checkbox"/> 5	210723084220.19242	Queued	Autovalidation	CPTAC3/LUAD/Proteome/01CPTAC_LUAD_Proteom
<input type="checkbox"/> 6	210723084221.19243	Queued	Autovalidation	CPTAC3/LUAD/Proteome/01CPTAC_LUAD_Proteom

**Spectrum Mill - Request Queue eldorado**

[Request Queue](#) [Completion Log](#) [Help](#)

There are 3 requests in the queue. Available memory: 18.4 Gb of 25.8 Gb

[Remove](#)

X #	Task Id	Status	Task Type	Data Directory
<input type="checkbox"/> 1	210723091357.9	Queued	Autovalidation	HLA/External_Data/PXD01349/PXD019643_humanatlas/DN
<input type="checkbox"/> 2	210723091358.10	Queued	P/P Summary	HLA/External_Data/PXD01349/PXD019643_humanatlas/DN
<input type="checkbox"/> 3	<a href="#">210723091355.3P Monitor</a>	Running (100:1493)	MS/MS Search	HLA/External_Data/PXD01349/PXD019643_humanatlas/DN

**Spectrum Mill - Request Queue shiva**

[Request Queue](#) [Completion Log](#) [Help](#)

There are 3 requests in the queue. Available memory: 66.8 Gb of 90.2 Gb

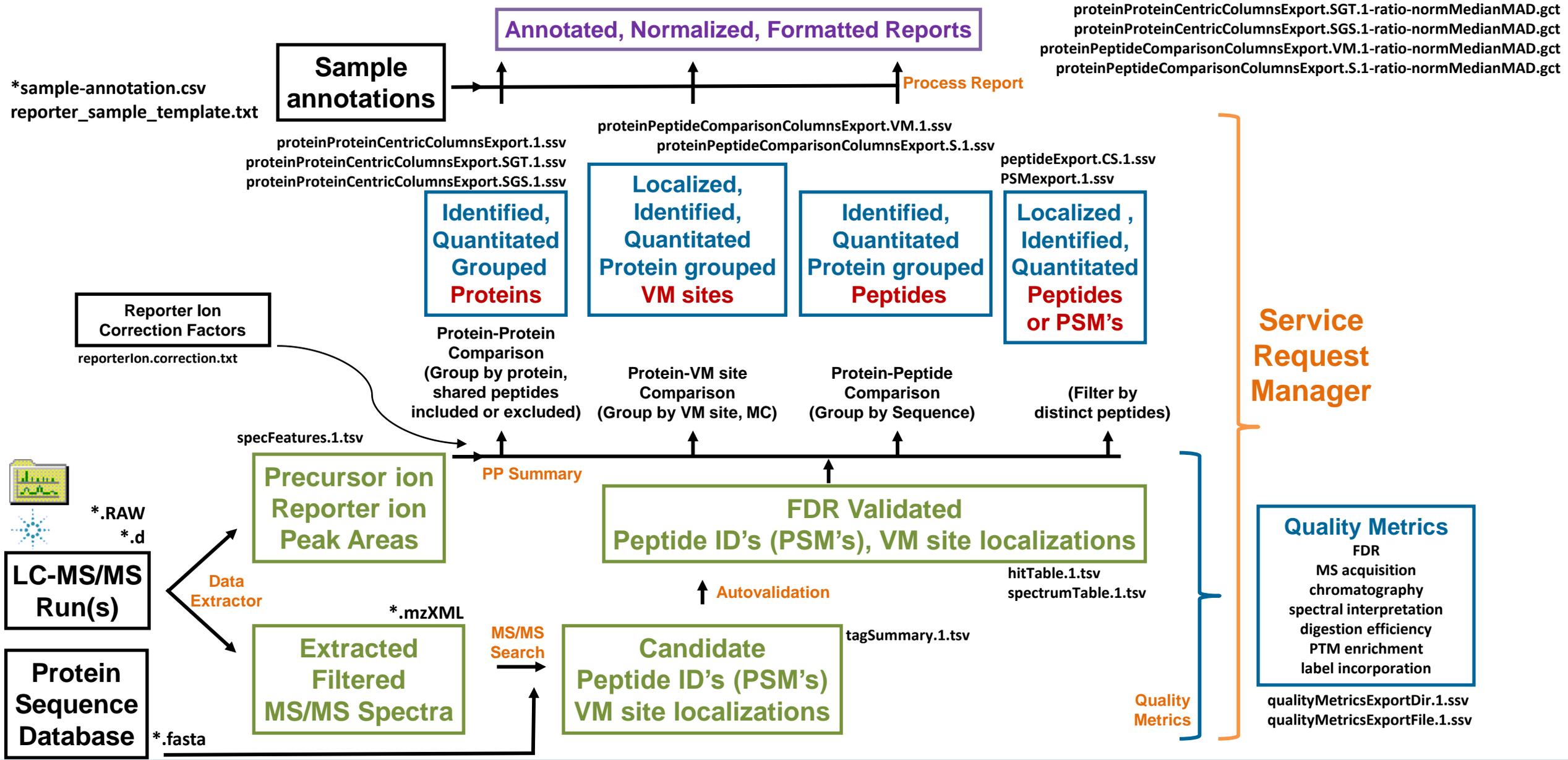
[Remove](#)

X #	Task Id	Status	Task Type	Data Directory
<input type="checkbox"/> 1	210723082331.9793	Queued	Autovalidation	CPTAC3pancancer/PDAC/Phosphoproteome/16CPTAC_F
<input type="checkbox"/> 2	210723082332.9794	Queued	Autovalidation	CPTAC3pancancer/PDAC/Phosphoproteome/16CPTAC_F
<input type="checkbox"/> 3	<a href="#">210723082330.9791P Monitor</a>	Running (384:998)	MS/MS Search	CPTAC3pancancer/PDAC/Phosphoproteome/16CPTAC_F

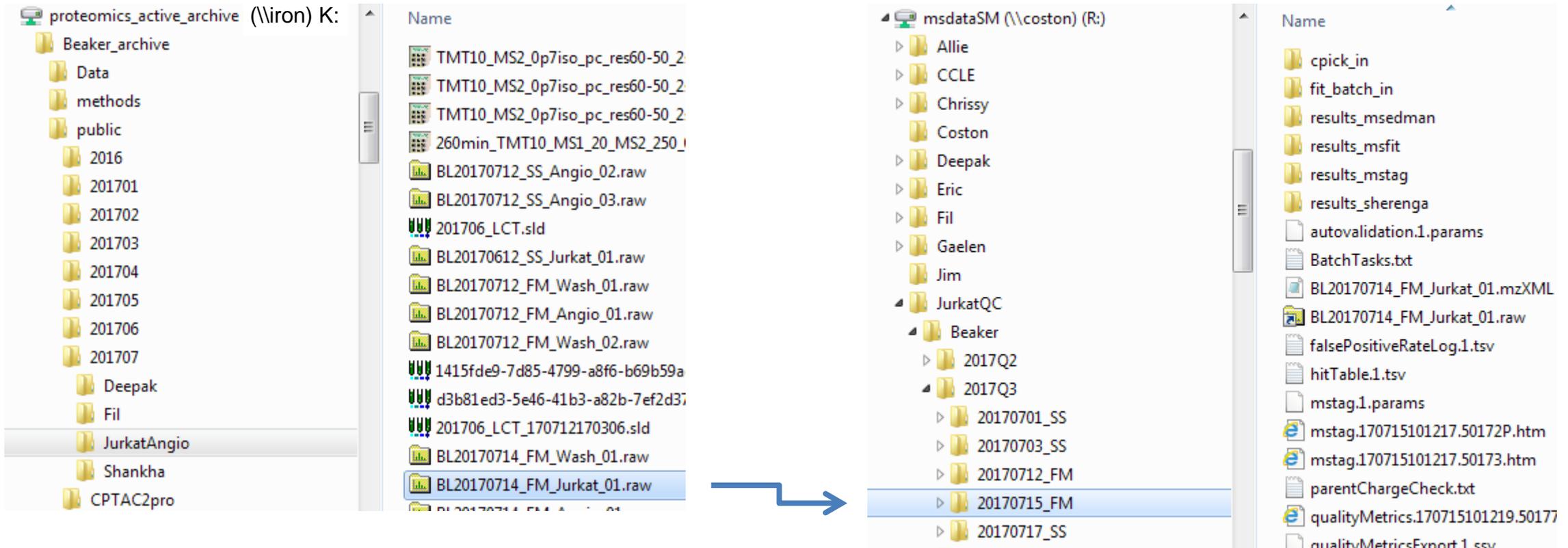
108 CPUs  
Grinding  
out PSMs



# Spectrum Mill Workflow



# Use Windows Explorer to Move data files to the SM server



## At Broad Institute

- Nightly, data files backed up from instrument to archive server.
- Following extraction, raw data files automatically deleted from SM server, replaced by shortcut.

# Maintaining Available SM Server Disk Space for **Active** Data Analysis

## Idle Datasets

Migrate to:

\\iron/proteomics\_storage\_slow/A\_SpectrumMill

Older (read-only access)

\\iron/proteomics\_active\_archive/SpectrumMill\_archive

\\iron/proteomics\_inactive\_archive/SpectrumMill\_archive

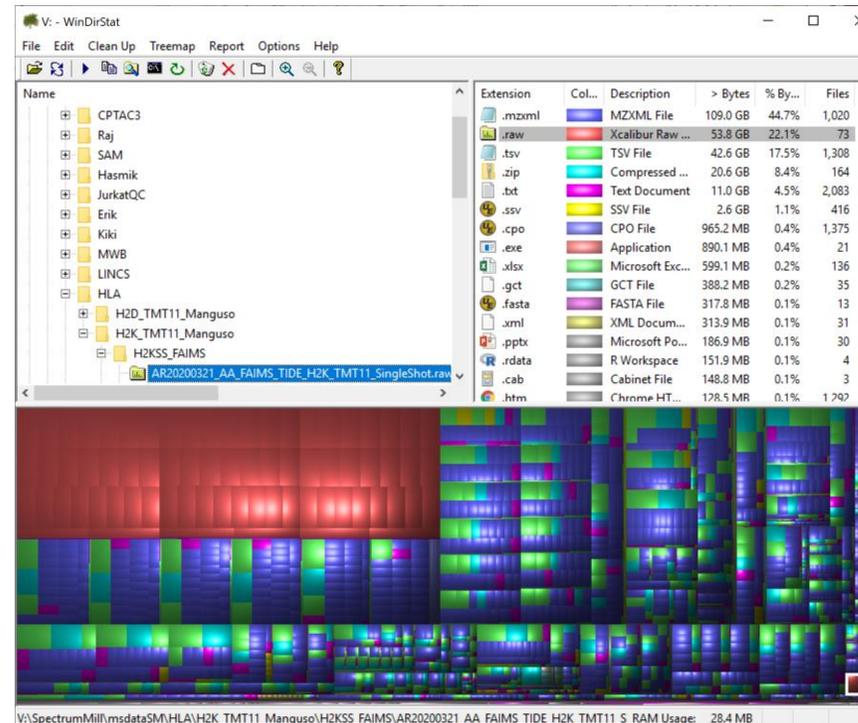
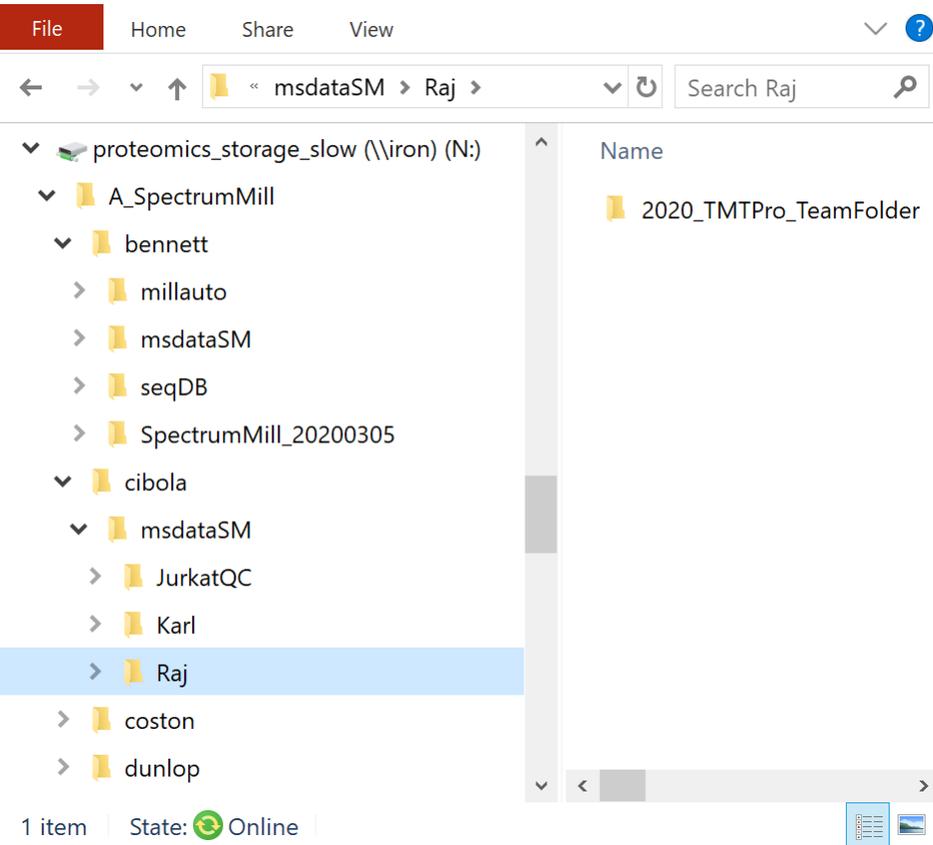
Might Ask,  
Tell via Slack

## Stale .RAW files

Delete & leave behind a placeholder

User can tell what was once there

Don't ask,  
Don't tell,  
Just do it



fresh

stale

Spectrum Mill - Archive Data - (none)

Spectrum Mill Extractor MS/MS Search Auto

Archive Data

Archive  Queue request Save As...

Archive Parameters

Instrument created files

Ignore instrument data files

Delete data files after making placeholder

- Creates a placeholder file for each raw data file, then deletes the raw data files.
- Afterwards, the previously extracted data may be searched and summarized, but the raw data files cannot be re-extracted.

Data Directories

Select ...

HLA/H2K\_TMT11\_Manguso/H2KSS\_FAIMS

# Configuring an Extraction

Spectrum Mill - Data Extractor - Defaults\Thermo\_RAW\_45similar\_STL0\_z5\_600\_6000\_IAA\_HCDv430

Extraction

Extract Save As... Load...  Maximize CPUs  Remove all prior results

Instrument: ESI QExactive HCD v4 30

Data Directories

Select ...  ExampleData/Thermo/Exploris480\_JurkatQC

Modifications

Choose... Fixed: Carbamidomethylation (C)

MS/MS Spectral Feature Filtering

Precursor MH+ 600.0 to 6000.0 Da

Scan time range: 0 to 300 min

Sequence tag length > 0 (For MALDI: Set tag length to -1 and merge secs to total run time.)

Ignore spectra with dissociation mode:  CID,  ETD,  PQD,  HCD

Merge nearby MS<sup>n</sup> scans with same precursor m/z

Retention time & m/z tolerance: +/- 45 secs +/- 1.4 m/z  
(also used for calculating chromatographic peak area of precursor in MS scans)  
(m/z tolerance ignored for high-resolution MS1, dynamically determined based on resolution)

General MS/MS Merging Constraints: Spectral Similarity & RT & m/z

Specialty MS/MS Merging Options:

Same resolution:  Merge CID & HCD MS<sup>n</sup>  Merge CID & PQD MS<sup>n</sup>

Different resolution:  Merge ion trap CID & HCD MS<sup>n</sup> immonium ion region

Merge MS<sup>2</sup> and MS<sup>3</sup> spectra from same precursor: Merge

Precursor m/z & Charge Assignment

Centroiding algorithm for profile mode data: Xcalibur

Precursor Charge: Find Maximum (z): 5 Minimum MS1 S/N: 25

Find <sup>12</sup>C precursor m/z

Results

Thermo Fisher  
Instruments

Spectrum Mill - Data Extractor - Defaults\Agilent\_D\_60similar\_STL0\_z5\_600\_6000\_IAA

Extraction

Extract Save As... Load...  Maximize CPUs  Remove all prior results

Data Directories

Select ...  ExampleData/Agilent/QTOF6546\_plasma

Modifications

Choose... Fixed: Carbamidomethylation (C)

MS/MS Spectral Feature Filtering

Precursor MH+ 600.0 to 6000.0 Da

Scan time range: 0 to 300 min

Sequence tag length > 0 (For MALDI: Set tag length to -1 and merge secs to total run time.)

Ignore spectra with dissociation mode:  CID,  ETD

Merge nearby MS<sup>n</sup> scans with same precursor m/z

Retention time & m/z tolerance: +/- 60 secs +/- 1.4 m/z  
(also used for calculating chromatographic peak area of precursor in MS scans)

General MS/MS Merging Constraints: Spectral Similarity & RT & m/z

Merge MS<sup>2</sup> and MS<sup>3</sup> spectra from same precursor: Merge

Precursor m/z & Charge Assignment

Precursor Charge: Find Maximum (z): 5 Minimum MS1 S/N: 25

Find <sup>12</sup>C precursor m/z

MS Noise threshold: 100 counts (only applies to Agilent Q-TOF)

Results

Agilent  
Instruments

- Convert to .mzXML
- Filter out low quality MS/MS
- Extract reporter ion intensities
- Calculate precursor ion XIC's
- Correct precursor monisotopic m/z assignments
- Merge replicate MS/MS from same chromatographic peak

# Configuring an MS/MS Search

MS/MS Search - Spectrum Mill

Not secure | shiva.broadinstitute.org/millhtml/batchtagpara.htm

Spectrum Mill - MS/MS Search - CPTAC3human\LSCC\_Phospho\_TMT11FullLysOnly\_Phospho\_v4\_CU\_AmqcstynG\_GencodeNuORFs

Spectrum Mill | Extractor | Autovalidation | Quality Metrics & FDR | Protein/Peptide Summary | Workflows | Tool Belt | Databases | Help

**Search**

Start Search | Save As... | Load... |  Remove all prior MS/MS Search results |  Maximize CPUs |  Disable Skipping Repeat Peptides in Database

**Data Directories**

Select ... |  CPTAC3pancancer\LSCC/Phosphoproteome\_GencodeGS/01CPTAC\_LSCC\_Phosphoproteome\_BI\_20190619

**Search Parameters**

Validation filter: spectrum-not-marked-sequence-not-validated | Batch size: 500

Search previous hits | Max reported hits: 5

Database: Gencode\_v34\_3nr.602contams.2043smorfs.nuORFv1.110LSCCgs.fasta | Digest: Trypsin allow P

Species: All | Maximum # missed cleavages: 4

**Modifications**

Choose... | Fixed: Carbamidomethylation (C,U), TMT11 Full Lys only-mix (N-term,K) | Variable: Acetyl (ProtN-term), Oxidized methionine (M), Pyroglutamic acid (N-termQ), Deamidated NG (N), Phosphorylated S (S)

**Search Criteria**

Matching Tolerances	Search Mode	Spectral Quality Filtering
Instrument: ESI QExactive HCD v4 35	Search mode: Variable modifications	<input type="checkbox"/> Max sequence tag length: > 3
Minimum matched peak intensity: 30 %	Precursor mass shift range: -18.0 to 272 Da	<input type="checkbox"/> Precursor isotope quality XIC's (Chi2 vs. Average): > 0.7
Masses are: Monoisotopic	Data Files	<input type="checkbox"/> Precursor Isolation Purity: > 70 %
Precursor mass tolerance: +/- 20 ppm	Fragmentation mode: All	<input type="checkbox"/> Glyco Product Ions Score: > 4.5
Product mass tolerance: +/- 20 ppm	Spectrum files:	
Maximum ambiguous precursor charge: 3	*.pk1	

Acknowledgment of submission to SRM request queue

Select Data Directories

OK | Cancel | Help |  Make Default

To Select, Click or Ctrl-Click on one or more Data Directories

- msdataSM
  - CPTAC3pancancer
    - BRCA
    - COAD
    - GBM
    - HNSCC
    - LSCC
      - Phosphoproteome\_GencodeGS
        - 01CPTAC\_LSCC\_Phosphoproteome\_BI\_20190619**
        - Proteome\_GencodeGS
  - LUAD
  - MEDUL
  - OVHGSC
  - PDAC
  - sample-annotation
  - UCEC
- ExampleData
- Karl
- Susan

# Fixed, Mix and Variable Modifications

## Fixed

Redefine the wild type as

## Mix

TMT 10 full Lys only mix

Search in 2 cycles

Cycle 1: all K , and Nterm TMT10 labeled

Cycle 2: all K TMT10 labeled, Nterm unlabeled

TMT 10 partial mix

Search in 4 cycles

Cycle 1: unlabeled

Cycle 2: all K , and Nterm TMT10 labeled

Cycle 3: all K , and Nterm unlabeled

Cycle 4: all K unlabeled, Nterm TMT10 labeled

Choose Modifications

OK Cancel Reset Details Help

**Fixed/Mix Modifications**

Cysteine  
Carbamidomethylation (C,U)

N-terminus  
TMT10 Full Lys only-mix (N-term,K)

C-terminus  
- unmodified -

Other amino acids

D: - unmodified -

E: - unmodified -

K: TMT10 Full Lys only-mix (N-term,K)

M: - unmodified -

U: Carbamidomethylation (C,U)

Metabolic isotope labels

All: - unmodified -

K: - unmodified -

L: - unmodified -

M: - unmodified -

P: - unmodified -

R: - unmodified -

V: - unmodified -

**Variable Modifications**

- Methylation of RG (R)
- Methylation of PR (R)
- Dimethylation of PR (R)
- Dimethylation of R (R)
- Citrullination of R (R)
- Hydroxylation of PG (P)
- Deamidated (Q)
- Deamidated 18-O (N)
- Methylation of Q (Q)
- Dioxidation of W (W)
- Kynurenin (W)
- iTRAQ-Y (Y)
- iTRAQ (Y,S,T,H)
- TMT0 (Y,S,T,H)
- TMT10 contains His (Y,S,T)
- TMT10 (Y,S,T,H)
- Nitration of Tyr (Y)
- Sulfation of Tyr (Y)
- Biotin Tyramide Labeling (W)
- Biotin Tyramide Labeling (Y)
- Biotin Tyramide Labeling (C)
- Biotin Tyrosine (Y,C,W)
- DAS Scar (Y)
- Carbamidomethyl His (H)
- Methylation of H (H)
- CMK adduct (H,S)
- Pyroglutamic acid (N-termE)
- Gamma Carboxylation of E (E)
- Hydroxylation of E (E)
- Hydroxamic acid (E,D)

## Variable

Allow

2 possibilities

for an AA.

Allow both  
in 1 spectrum

if more than one  
location/AA.

## Overlabeled

(TMT on S,T,Y

only when

Peptide

Contains His)

smaller search space

# Database Search: General Considerations

- Choice of database
  - reference vs. personalized + contams + non Canonical ORFs + tagged bait protein
    - reference proteome (UniProt, RefSeq, Ensembl, Gencode)
  - decoy databases – estimation of false discovery rate
- Choice of enzyme
  - full vs. partial enzyme specificity, no enzyme for immunopeptidomics
- Choice of fixed and variable modifications
  - fixed modifications: target AA always considered as modified
  - variable modifications: target AA may or may not be modified
  - *expansion of search space!*
- Precursor ion (peptide) mass tolerance
  - depends on the measurement mass accuracy
  - lower tolerance decreases search space (fewer candidate masses)
  - higher tolerance helps FDR statistics, requires post-search tolerance filtering (3 std dev of mean)
- MS/MS fragment ion mass tolerance
  - may not be the same as precursor mass tolerance (hybrid instruments)



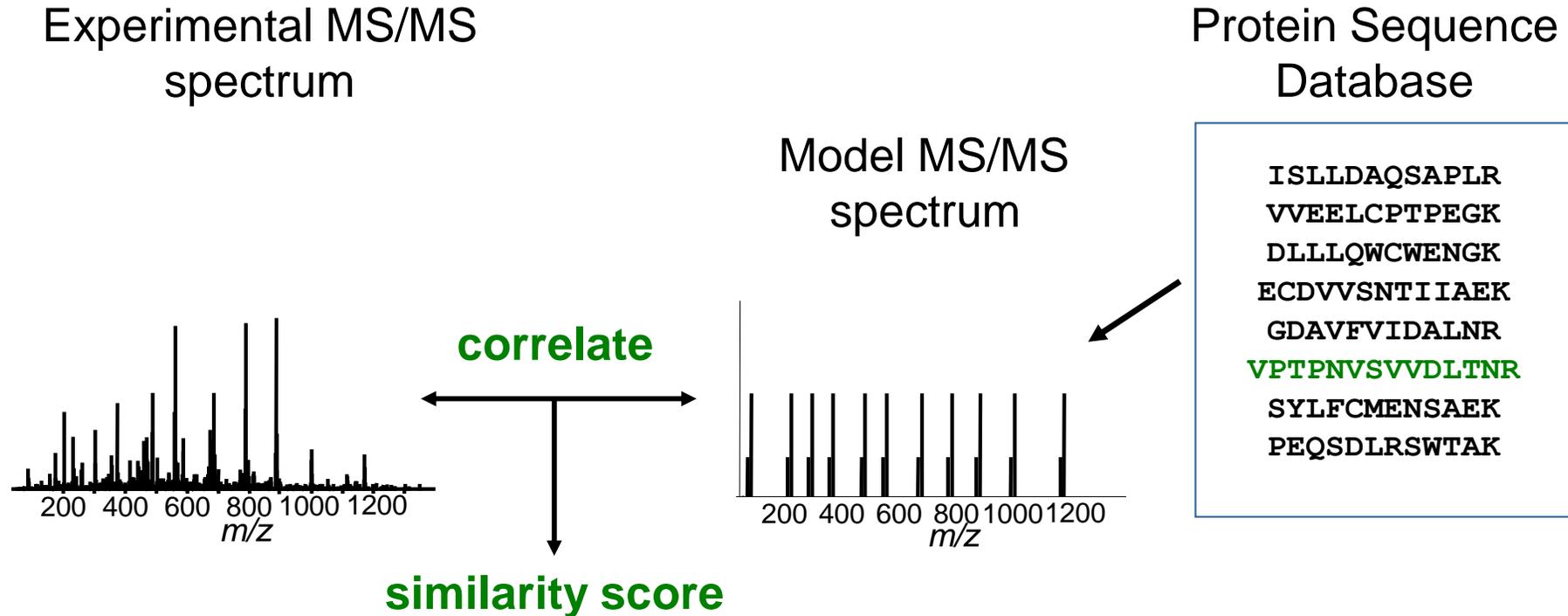
# Key Features of Spectrum Mill Database Searching with MS/MS Spectra

- Scoring specific to instrument types, dissociation modes
  - Ion types allowed and scoring weights
  - On board peak detection (de-isotoping, S/N thresholding, fragment z assignment).
- Product ion mass tolerance in ppm or Da.
- Spectra grouped in batches of similar precursor m/z, z for fast precursor mass filtering.
- Target-decoy FDR accomplished by on-the-fly peptide inner reversal
  - SAMPLER becomes SELPMAR
  - Search time is only ~1.5x as long as a target only database, instead of 2x for a concatenated forward/reverse database. 1x digestion of all proteins, 1x precursor mass filtering of all peptides, 2x MS/MS matching of every sequence (fwd & rev)
  - Search results do not comingle target and decoy hits.
    - For each PSM report top target hit, and delta Fwd – Rev score of top decoy hit
    - A false positive PSM had delta Fwd-Rev < 0. (Top Decoy hit has higher score)
- Multiple cycles for fixed/mix modifications or unknown charge, with combined single result.
- Variable modifications constrained by precursor mass shift, not by # of mods/peptide.
- Second pass search possible with leftover spectra from 1<sup>st</sup> pass.
- Homology modes available
  - Unassigned single mass gap
  - Allow single mutation per peptide



# PSM Scoring

# MS/MS Search Engines: look up answer in back of book



- Find best matching database peptide
- Can miss unanticipated modifications
  - post-translational, sample handling
- Determine peptide FDR by searching reversed DB

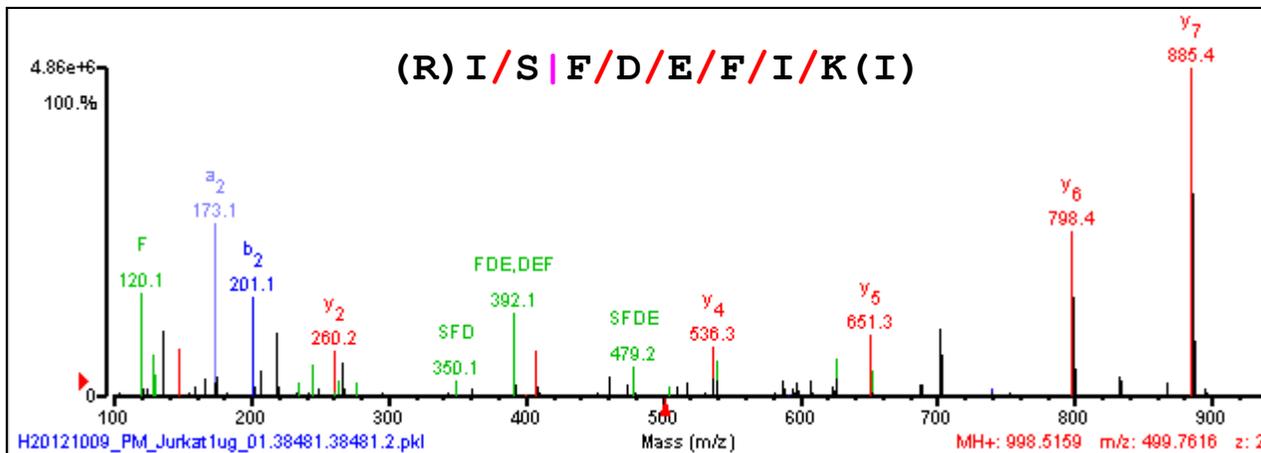
Algorithms: Mascot, Sequest, MaxQuant, Spectrum Mill, MSGF+, MSFragger, PEAKS, X-Tandem...

# PSM scoring considerations

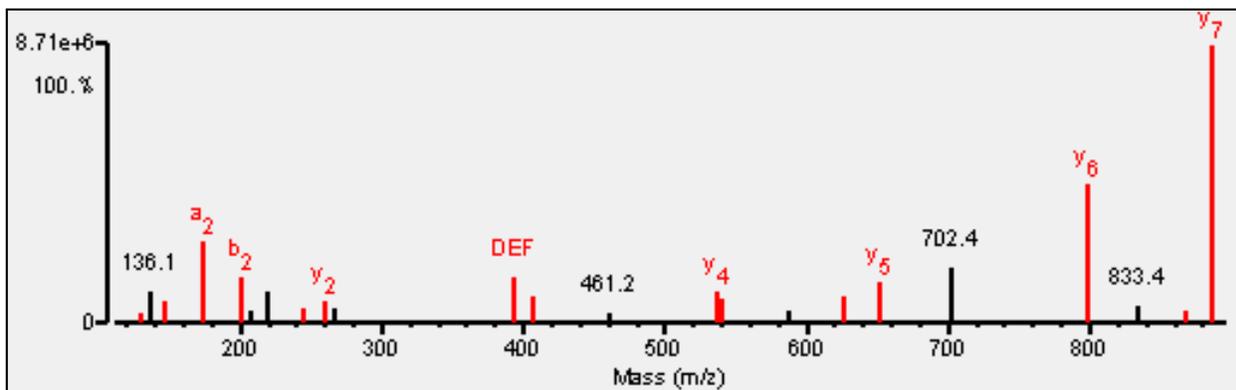
- Signal pre-processing
  - deisotoping, fragment charge assignment, noise removal
- Fragment ion types allowed and relative scoring weights
  - b/y, neutral loss (-H<sub>2</sub>O, -NH<sub>3</sub>, -H<sub>3</sub>PO<sub>4</sub>), internal, immonium
  - specific to instrument dissociation method (HCD, CID, ETD)
- Peak intensity
  - Contribute to score or just peak detection
  - Intensity used in model spectra (fragment ion type, mass range)
- Fragment mass tolerance units appropriate to mass analyzer
  - Da (ion trap, quadrupole), ppm (Orbitrap, ToF)
- Use of secondary scores
  - Gap between rank1 rank2 peptide
  - Sequence coverage of peptide
  - Unexplained peaks in spectrum
  - Ion series continuity
- Size of database may effect score



# Spectrum Mill Scoring of MS/MS Interpretations



Peak Selection: De-Isotoping, S/N thresholding,  
Parent - neutral removal, Charge assignment  
Match to Database Candidate Sequences



Score  
=  
Assignment Bonus  
(Ion Type Weighted)  
+  
Marker Ion Bonus  
(Ion Type Weighted)  
-  
Non-assignment Penalty  
(Intensity Weighted)

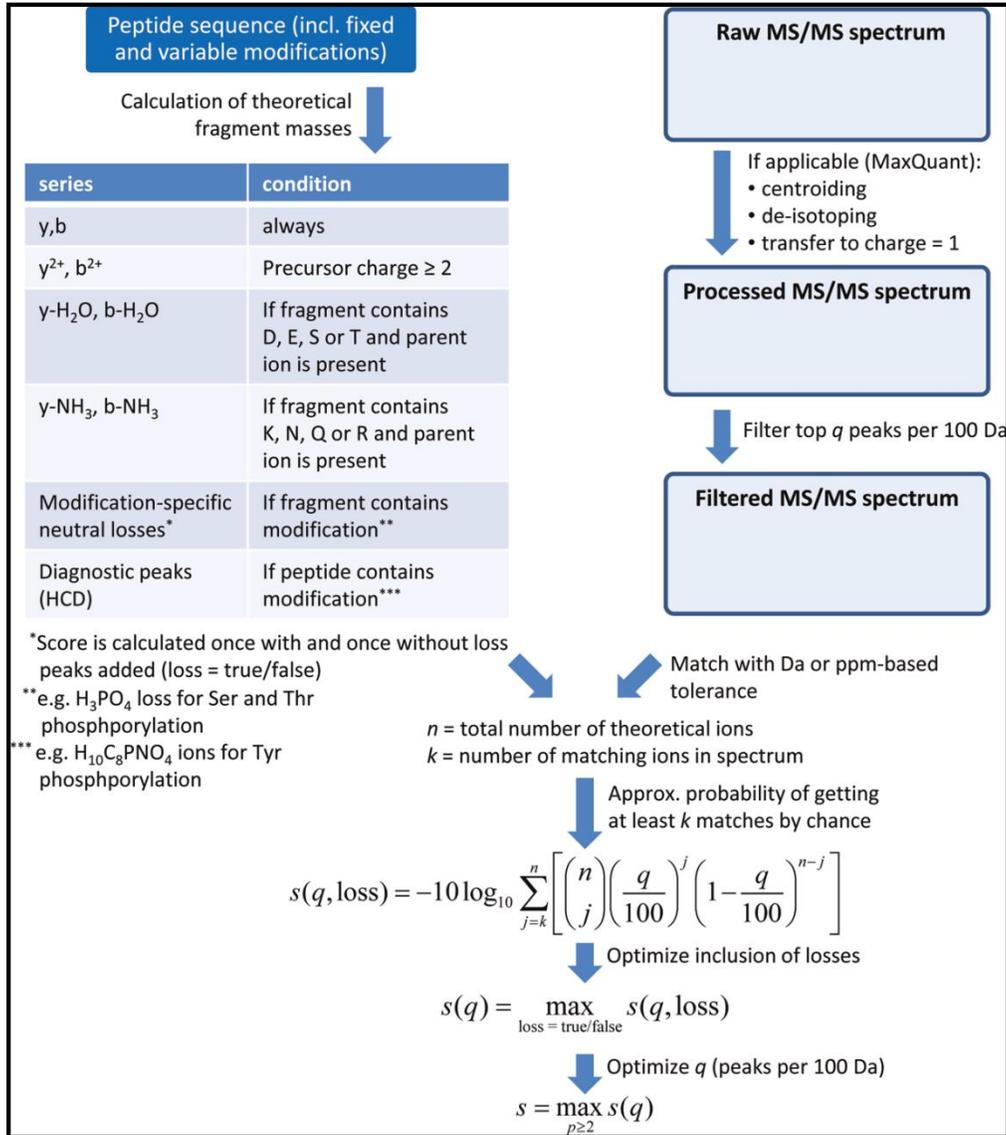
SPI (%)  
Scored Peak Intensity

PIP(%)  
Precursor Ion Purity

11.52  
81.5%  
80.3%

Fragment-ion (m/z)	120.080	129.102	130.086	136.075	147.112	173.128	201.122	207.112	219.112	245.075	260.195	267.115	392.143	407.263	461.163	536.304	539.210	587.366	626.242	651.332	702.390	798.400	833.430	867.416	885.430
Frac. Inten.(% of TIC)	-1.00	-1.00	-1.00	-1.00	-1.00	-1.00	-1.00	-1.00	-1.00	-1.00	-1.00	-1.00	-1.00	-1.00	-1.00	-1.00	-1.00	-1.00	-1.00	-1.00	-1.00	-1.00	-1.00	-1.00	-1.00
Rel. Inten.(% of BP)	0.53	0.22	3.83	11.45	8.14	29.49	16.93	4.54	11.12	5.41	8.07	5.92	16.63	9.57	3.61	11.65	9.01	4.52	10.00	15.16	20.21	50.16	6.43	4.73	100.00
Signal/Noise	17.8	6.4	2.7	11.1	7.4	30.8	17.0	3.5	10.7	4.4	7.4	5.0	16.7	9.0	2.5	45.2	33.6	14.0	38.0	60.6	82.6	213.7	22.4	14.9	431.7
Score	1.00	0.20	0.250	-0.115	1.000	0.250	1.000	-0.045	-0.111	0.500	1.000	-0.059	0.500	1.000	-0.036	1.000	0.500	-0.045	1.000	1.000	-0.202	1.000	-0.064	0.500	1.000
Ion-type	F	K	y <sub>1</sub> -NH <sub>3</sub>		y <sub>1</sub>	a <sub>2</sub>	b <sub>2</sub>			DE	y <sub>2</sub>		DEF	y <sub>3</sub>		y <sub>4</sub>	FDEF		SFDEF	y <sub>5</sub>		y <sub>6</sub>		y <sub>7</sub> -H <sub>2</sub> O	y <sub>7</sub>
Delta ppm	-4.1	-5.0	-4.1		-6.2	-5.5	-5.0			-7.7	-7.9		-5.9	-5.6		-7.2	-6.7		-6.2	-4.0		-4.6		-9.6	-5.7

# PSM Identification Scoring – Andromeda, MaxQuant



## Peak selection assumptions

- Expects to use same # of peaks for each 100 Da region of spectrum
- Tall and short peak intensities equally diagnostic
- All regions of spectrum treated as equally likely
  - multiply charged fragments below precursor
  - some 100-300 m/z values not possible, di-pep AA combinations
- Tolerance in Da, not ppm

## Scoring

- Scoring implies +/- 0.5 Da fragment tolerance
  - Matches separately constrained by user specified ppm or Da fragment tolerance
- Scoring implies all ion types have same value
- All ion types used for scoring

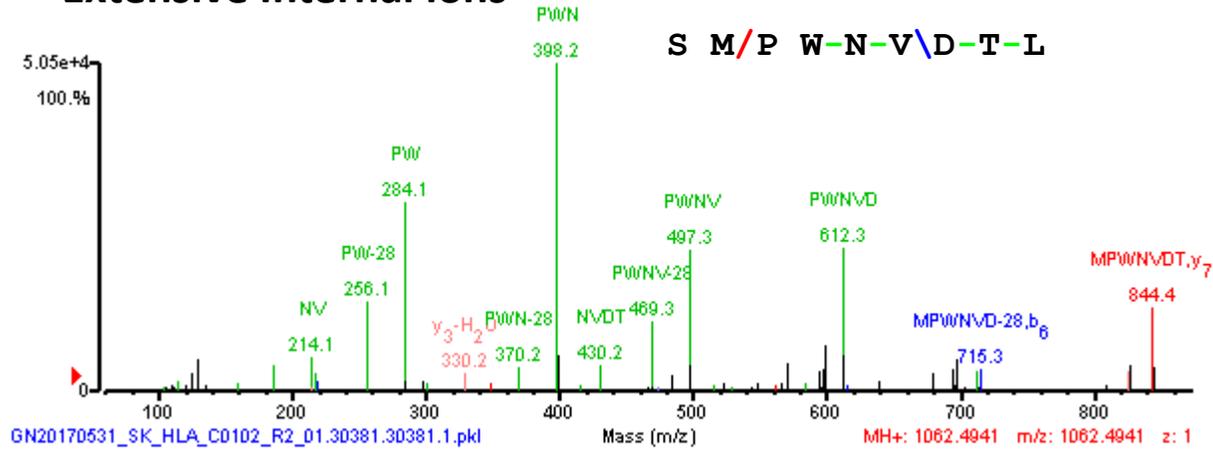
## Final score optimizes

- # of highest intensity peaks used per 100 m/z
- Inclusion/exclusion of neutral loss ion types

Karl's take: Effective score, despite some awkward assumptions made to squeeze MS/MS interpretation into a binomial probability framework

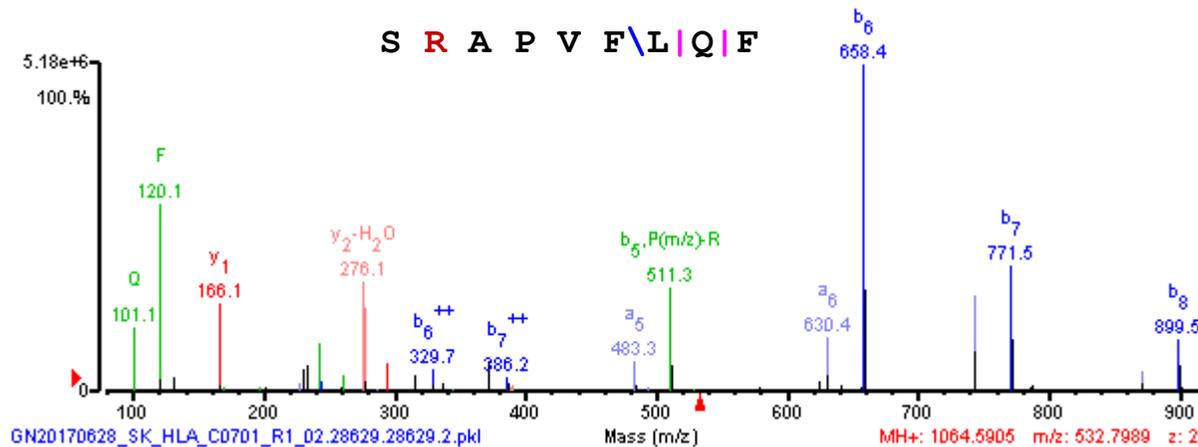
# Example HCD MS/MS of range of scoring challenges

## Extensive Internal Ions



This peptide with no basic residues and a singly charged precursor ion fragments to yield mostly internal sequence ions.

## Low Sequence Coverage



This peptide fragments to yield many different ion types, but they are derived from only 3 positions along the peptide backbone.

# Key aspects of Spectrum Mill DB-search scoring revisions for HLA peptides

## Added PSM quality threshold separate from primary score:

Minimum Backbone Cleavage Score (BCS) of 5

- Enforces uniformly higher minimum sequence coverage for each PSM, at least 4 or 5 residues of unambiguous sequence.
- Decreases false-positives associated with PSMs having fragmentation in a limited portion of the peptide that yield multiple ion types.
- BCS is a sequence coverage metric. The score is 1 or 0 for cleavage between adjacent AA's in the sequence (max score is peptide length-1). To receive a score, cleavage of the peptide backbone must be supported by presence of a primary ion type for HCD: b, y, or internal ion C-terminus (i.e. if the internal ion is for PWN then BCS is credited only for the backbone bond after the N).

**BCS Example** (top spectrum)

S M/P W-N-V\D-T-L

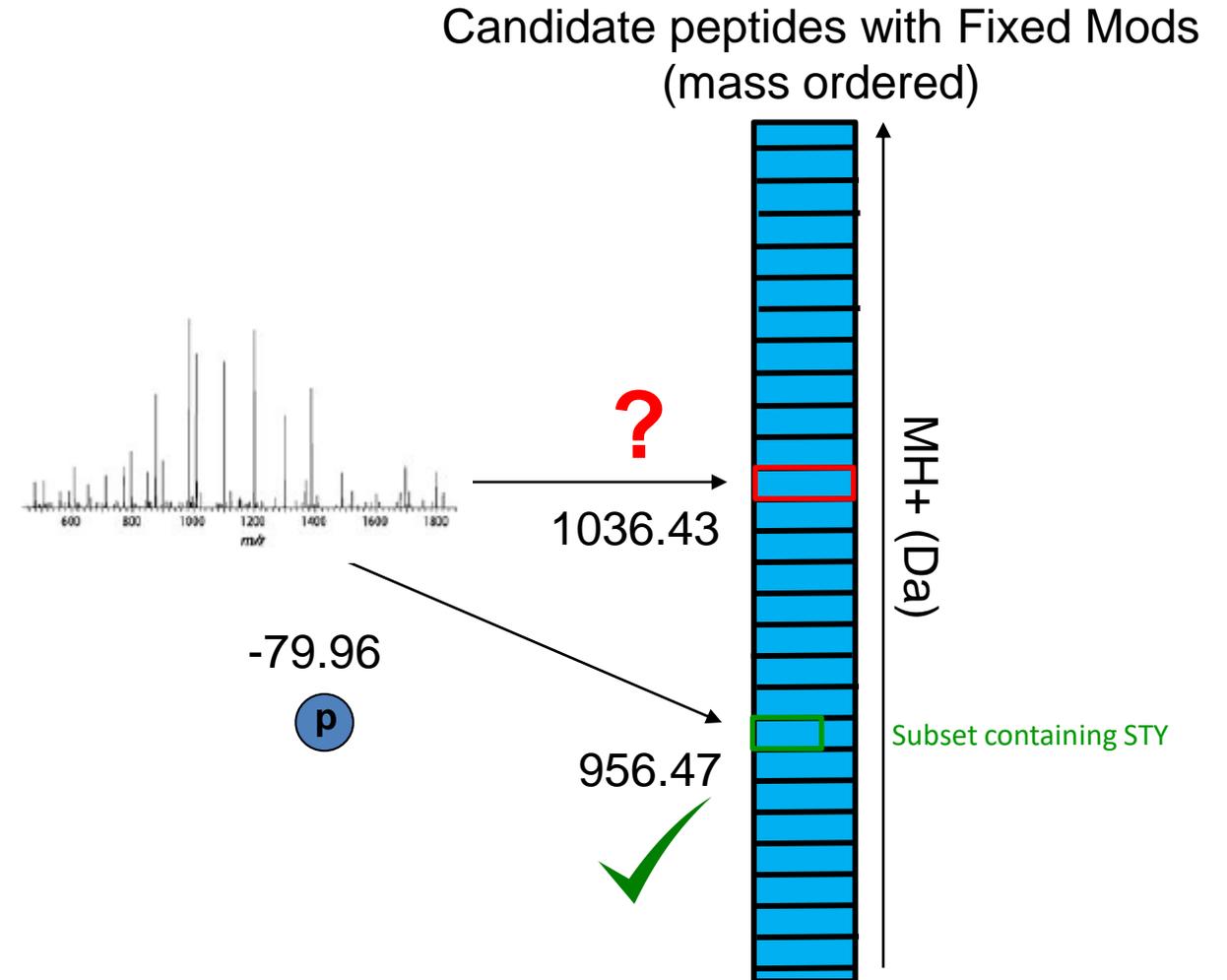
BCS: 6 /: y<sub>7</sub> \: b<sub>6</sub> ion -: internal ions (PW, PWN, PWNVD, PWNVDT)

**Enabled additional internal ion types:** neutral loss of NH<sub>3</sub>, H<sub>2</sub>O, and CO.

- Minimal bonus score (10% of the score of b or y ions) for presence accompanying an internal b-ion (75% of the score of b or y-ions).
- Prevents intensity-based penalty score when the ion types are not enabled.

# Search Space Expansion for Identification of Variably Modified Peptides

- Modifications alter the mass of the corresponding protein/peptide
- Have to be considered in database search
- Fixed modification: corresponding amino acid is always modified (e.g. carbamidomethylation of C)
- Variable modification: corresponding amino acid might or might not be modified (e.g. phosphorylation of S/T/Y)



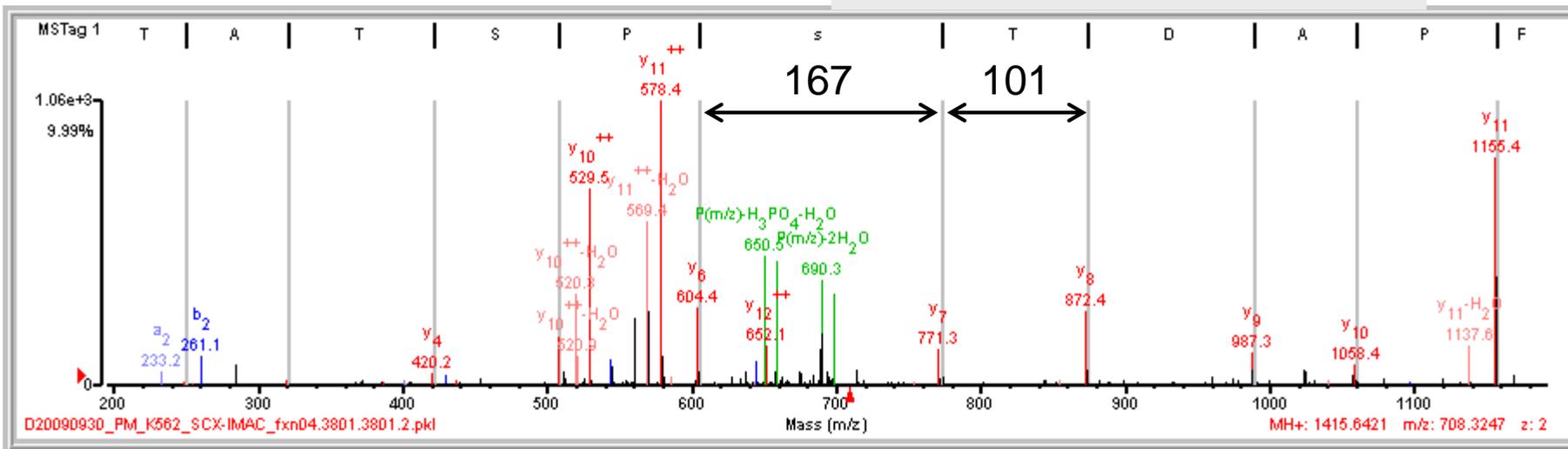
Variable mods require the search engine to evaluate candidate sequences with a precursor mass **offset** from the measured mass

# Phosphosite Localization

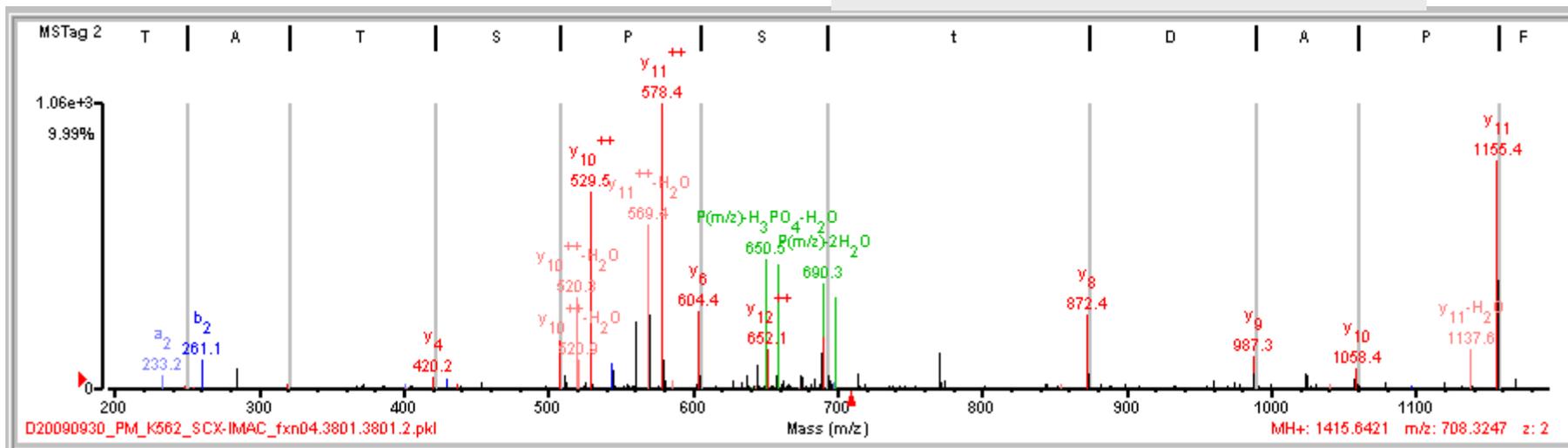
# Localizing a Phosphorylation Site

same spectrum  
2 different interpretations

L/F | P/A/D | T/s/P/S T A\T K



L/F | P/A/D | t S/P/S T A\T K



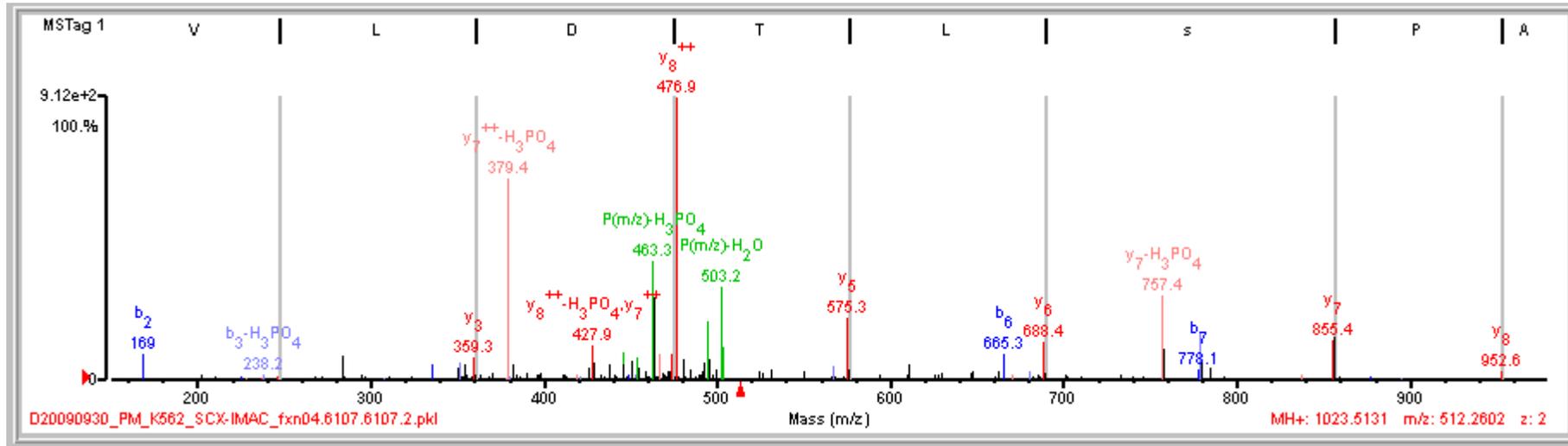
# PTM Site Localization

Test all Locations, Examine Score Gaps

	<u>Locations Tested</u>	<u>Conclusion</u>
No possible ambiguity	AV <b>s</b> EEQQPALK	AV <b>S(1.0)</b> EEQQPALK # PO <sub>4</sub> sites = # S,T, or Y
Single Site	AP <b>s</b> LTDLVK *	AP <b>S(0.99)</b> LT(0.0)DLVK
	APSL <b>t</b> DLVK -	
	<b>s</b> SSAGPEGPQLDVPR *	
	<b>S</b> sSAGPEGPQLDVPR *	
	SS <b>s</b> AGPEGPQLDVPR -	<b>S(0.50)</b> <b>S(0.50)</b> <b>S(0.0)</b> AGPEGPQLDVPR
Multiple Sites	VTNDI <b>s</b> PE <b>s</b> SPGVGR *	VT(0.0)NDI <b>S(0.99)</b> PE <b>S(0.50)</b> <b>S(0.50)</b> PGVGR
	VTNDI <b>s</b> PE <b>S</b> SPGVGR *	
	VTNDI <b>S</b> PE <b>s</b> SPGVGR -	
	V <b>t</b> NDI <b>s</b> PE <b>S</b> SPGVGR -	
	V <b>t</b> NDI <b>S</b> PE <b>s</b> SPGVGR -	
	V <b>t</b> NDI <b>S</b> PE <b>S</b> SPGVGR -	



# PTM Site Localization – Confident Localization

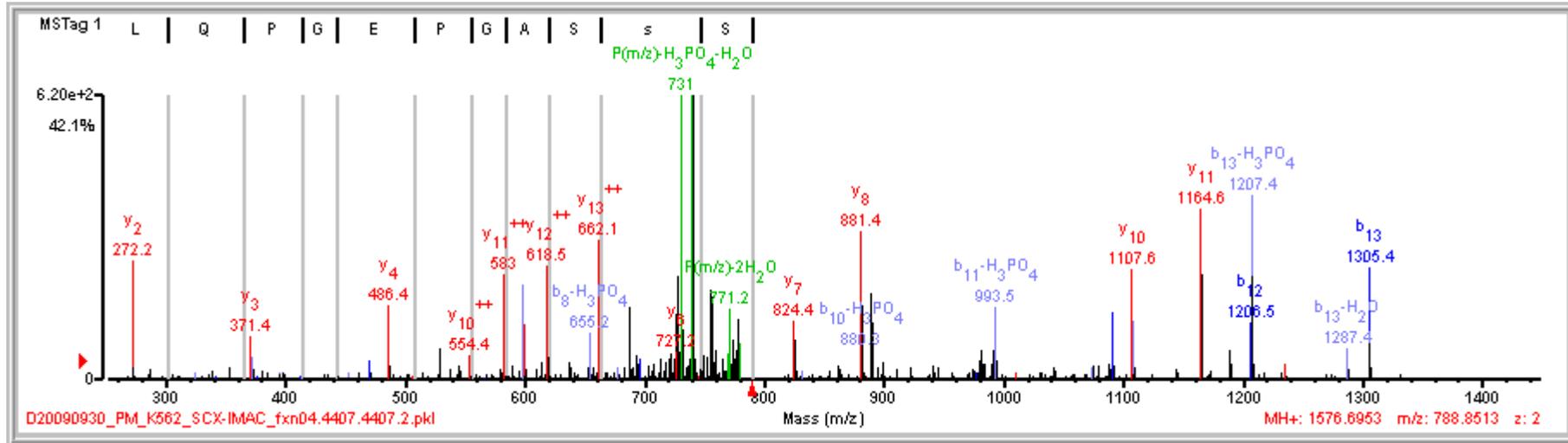


(K) A/P | s/L/T D | L/V K(S)

APS (0.99) LT (0.0) DLVK

$y_6$ ,  $y_7$  ions provide  
confident localization  
to the Ser

# PTM Site Localization – Ambiguous Localization



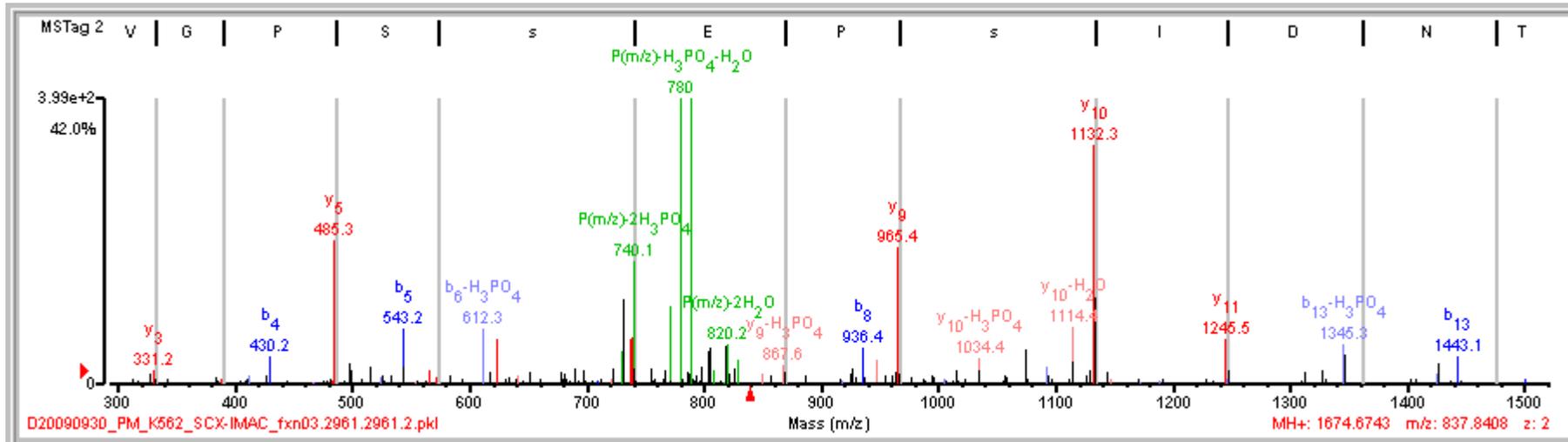
(R)S s/S/A/G/P E/G/P Q L|D|V|P R(E)

S(0.50)S(0.50)S(0.0)AGPEGPQLDVPR

$y_{13}^{++}$  ion excludes  
localization to Ser 3

# PTM Site Localization – Ambiguous Localization

## 2 sites: 1 confident, 1 ambiguous



(R)V T N D | I | s / P E | s S / P G V \ G R (R)

VT (0.0) NDIS (0.99) PES (0.50) S (0.50) PGVGR

y<sub>9</sub>, y<sub>10</sub> ions provide  
confident localization to  
the Ser -6

y<sub>9</sub>, y<sub>10</sub> ions provide  
ambiguous localization to  
Ser-9, Ser-10

# Spectrum Mill Variable Modification Localization Score

VML score = Difference in Score of same identified sequences with different variable modification localizations

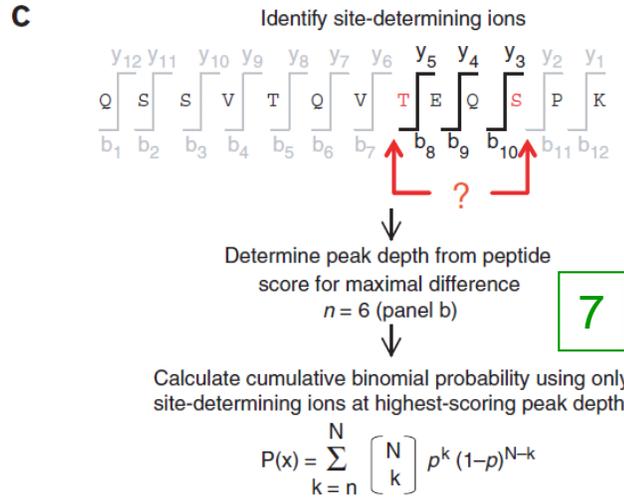
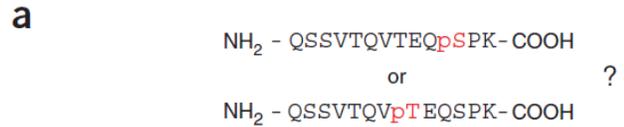
VML score > 1.1 indicates confident localization

Why a threshold value of 1.1?

1 implies that there is a distinguishing ion of b or y ion type

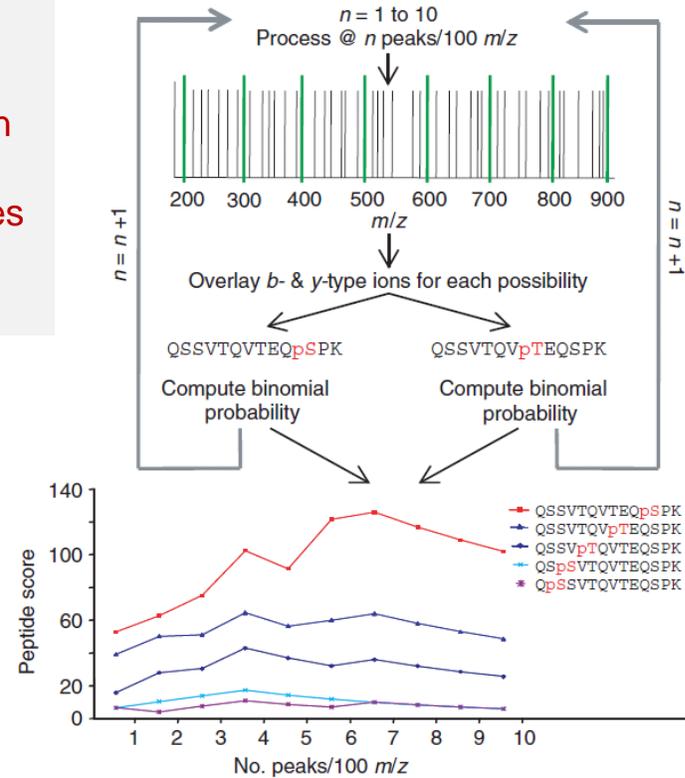
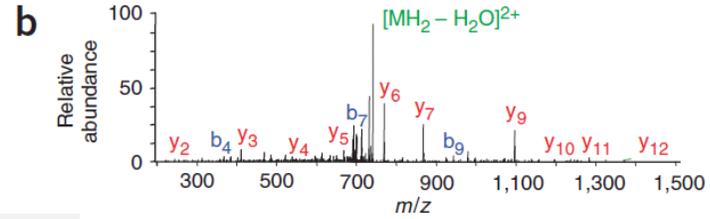
0.1 means that when unassigned, the peak is 10% the intensity of the base peak

# Phosphosite Localization Scoring - Ascore



- Typical score thresholds are equivalent to 2 peaks supporting localization
- Peak selection probability term implies +/- 0.5 Da tolerance instead of ppm

Phosphopeptide	QSSVTQVTEQ <b>p</b> SPK	QSSVTQ <b>Vp</b> TEQSPK
Trials ( $N$ )	6 ( $y_3, y_4, y_5, b_8, b_9, b_{10}$ )	6 ( $y_3, y_4, y_5, b_8, b_9, b_{10}$ )
Successes ( $n$ )	5 ( $y_3, y_4, y_5, b_9, b_{10}$ )	0
$p$ (6 peaks / 100 $m/z$ )	0.06	<b>0.07</b>
$P$	0.0000044	1.0
Score [ $-10 \times \log(P)$ ]	53.57	0
Ascore = ambiguity score (difference of the top two candidates)	53.57 - 0 = 53.57	



# VML Score Features Reported

Filename	z	Score	Fwd-Rev Score	Rank1-Rank2 Score	SPI (%)	# Backbone Cleavages	Missed Cleavages	#KR	Score VML	# STY	# sty Sites	# Loc sty Sites	# Ambig sty Sites	Variable Sites	VM site Flanks	Sequence Map VML sequence
BL20160819_FM_Medullo_Phosphoproteome_Plex02_Fxn06.0005.0005.3	3	12.94	12.94	12.94	68.5	8	0	1	0.460	5	1	0	1	S99s	RRHSHSsPMSTRRR	(R)H\S\H\S\H-s/P/M S T/R(R) HS(0.0)HS(0.50)HS(0.50)PMS(0.0)T(0.0)R
BL20160819_FM_Medullo_Phosphoproteome_Plex02_Fxn06.1390.1390.3	3	9.47	9.47	9.47	41.4	6	0	1	99	1	1	1	0	S623s	AHHAAGQsVRSGRLLG	(K)A\H H\A\A\G Q s/V/R(S) AHHAAGQS(1.0)VR
BL20160819_FM_Medullo_Phosphoproteome_Plex02_Fxn06.3872.3872.3	3	11.65	11.65	11.65	63.1	7	1	2	2.135	2	1	1	0	S437s	VLTKSGRsAHQVARY	(K)S G\I\s/A/H/Q\I V A/R(Y) S(0.0)GRS(0.99)AHQVAR
BL20160819_FM_Medullo_Phosphoproteome_Plex02_Fxn06.3908.4284.3	3	11.77	8.39	8.29	68.4	6	1	2	0.277	4	1	0	1	S409s	PPPKRHsPTPQSN	(K)T R H\s/P T/P/Q/Q S N/R(T) T(0.33)RHS(0.33)PT(0.33)PQQS(0.0)NR
BL20160819_FM_Medullo_Phosphoproteome_Plex02_Fxn06.4548.4548.3	3	9.98	6.78	4.81	71.6	6	1	2	1.214	4	1	1	0	S491s	EKGSSSRsPGPHRS	(K)S S\I S R\s/P G/P H/P/R(S) S(0.0)S(0.0)S(0.0)RS(0.99)PGPHPR
BL20160819_FM_Medullo_Phosphoproteome_Plex02_Fxn06.4548.4548.3	3	9.55	5.39	5.04	50.2	8	1	2	4.086	3	1	1	0	S61s	PASSHREsPRGSGGA	(R)G\I P A\S\H R\I\s/P/R(G) GPAS(0.0)S(0.0)HRES(0.99)PR
BL20160819_FM_Medullo_Phosphoproteome_Plex02_Fxn06.4638.4638.3	3	9.93	9.05	8.54	68.0	6	2	3	0.751	5	1	0	1	S298s	SRRSTTKsPGPSRRS	(R)S T\T K-s/P/G/P S R/R(S) S(0.0)T(0.0)T(0.50)KS(0.50)PGPS(0.0)RR
BL20160819_FM_Medullo_Phosphoproteome_Plex02_Fxn06.4673.4673.3	3	8.74	7.18	7.35	67.5	5	0	2	1.945	2	1	1	0	S694s	PNKRHSPsPRRAPQ	(R)H\S\I P s/P R/P/R(A) HS(0.0)PS(0.99)PRPR
BL20160819_FM_Medullo_Phosphoproteome_Plex02_Fxn06.4881.4881.3	3	8.33	3.69	2.24	57.1	5	2	3	1.379	3	1	1	0	T310t	GTARRTG\PSDPRRR	(R)R T G\I/P S D\I-P-R(R) RT(0.0)GT(0.99)PS(0.0)DPRR
BL20160819_FM_Medullo_Phosphoproteome_Plex02_Fxn06.5095.5095.3	3	13.13	13.61	10.25	78.1	7	0	2	7.333	4	1	1	0	S544s	GSSAPEHsPRTSGLG	(R)S\I P R\I P G S S\A\I P E\H s/P/R(T) S(0.0)PRPGS(0.0)S(0.0)APEHS(0.99)PR
BL20160819_FM_Medullo_Phosphoproteome_Plex02_Fxn06.5077.5077.3	3	11.46	7.21	6.63	69.7	7	0	2	99	1	1	1	0	S287s	LAQLLARsPPPHPRP	(R)s/P/P/P/P/H R/P/R(L) S(1.0)PPPHPRP

Site flanks  
+/- 7 AA  
in protein

**Summarize** Save As... Load...

Queue request  Excel  AMRT export

Mode: Peptide - Spectrum Match

Filter by Proteogenomic Features: Off

Data directories: Select ...

CPTAC3pancancer/MEDUL/Phosphoproteome\_...

Search result files: \*Fxn06\*.3.pk1.spo

Search result files exclude:

**Validation and Sorting**

Filter results by: valid

Validation preset: none

Sort peptides by: Retention Time

Filter peptides by: Score: > 8 % SPI: > 0

Required AAs: any

Disallowed AAs: none

Peptide pI: from 3.0 to 10.0  All

Accession #'s:

**Review Fields**

Filename  Score  Fwd-Rev score  Rank 1-2 score  SPI (%)  Backbone Cleavage Score  Glyco Product Ions Score  Unmatched ions  Var mod sites  VML score s|t|y  Solution charge  Start AA position  Proteogenomic features

Sequence  Rev Sequence  VML sequence  Prec Av Chi<sup>2</sup>  Ret time, width  Ion mobility  Precursor m/z  Delta mass  # Sequences per MH+  Protein MW  Species  Accession #  Protein name

b/y map  Rank 2  Isol Pur  MH+  Pep pI  Prot pI

Precursor XIC Intensity  Precursor S/N  DEQ ratios  Invert  Reporter Ratios: TMT 10  126  127N  127C  Intensities  Correction Method: None

Modification names  N-term  C-term  Cysteines  Fragmentation mode  Max tag length  Longest tag  # b/y pairs

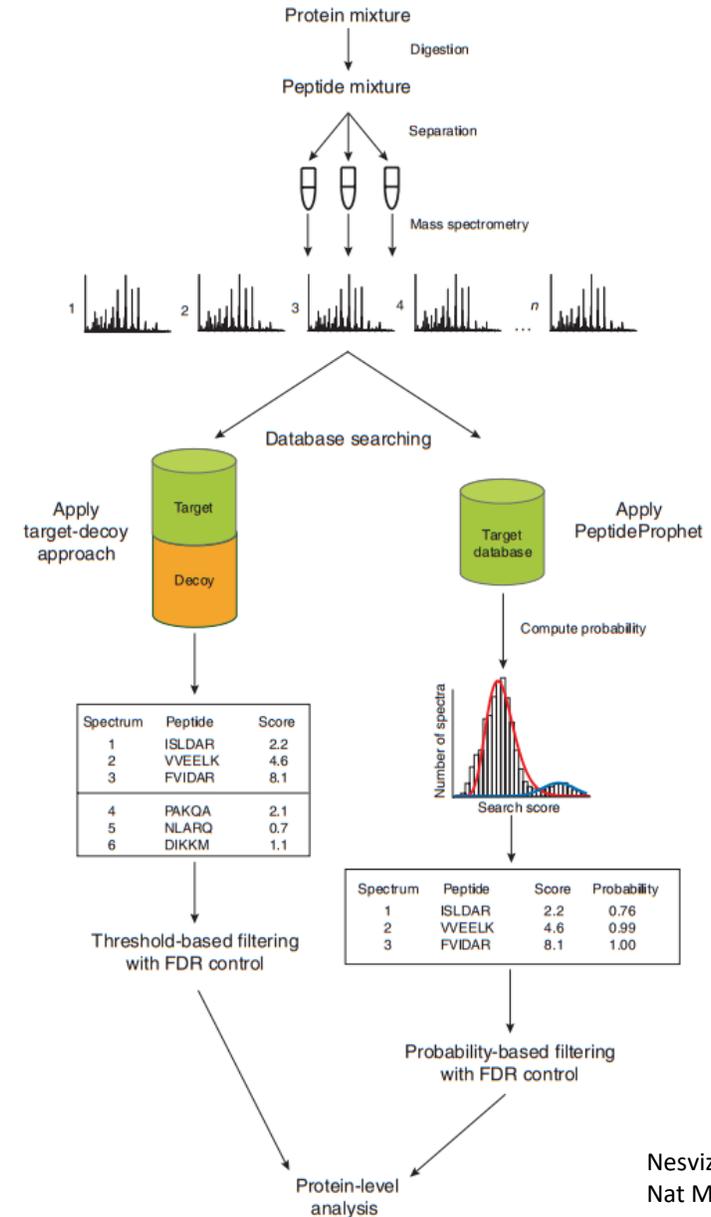
Category: Select ...

Genocode\_v34\_3nr.fasta.categories.tsv:geneSymb|

## Methods to control false discovery rates (FDR) of MS/MS search results

# Estimating False Discovery Rate (FDR)

- Basic Approaches
  - Target-decoy approach
  - Bayesian posterior probability
- Goals
  - Accurate error estimation
  - Accurate sensitivity estimation
  - Comparable across labs and experiments



Nesvizhskii AI, Vitek O, Aebersold R.  
Nat Methods. 2007 Oct;4(10):787-97.

# Creating a Target/Decoy Database

## Reverse each protein sequence

- Amino acid composition conserved
- #, length of proteins conserved
- **# of peptides conserved**

## Shuffle each protein sequence

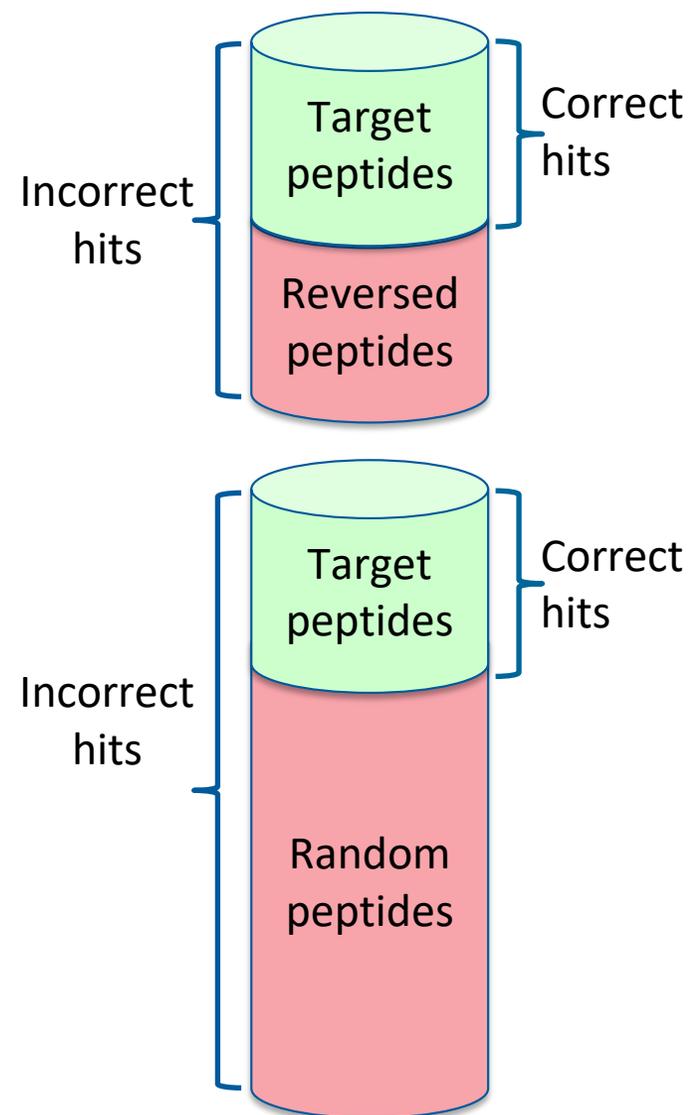
- Amino acid composition conserved
- #, length of proteins conserved
- **# of peptides increases**
  - Protein homology is lost
  - Yields more distinct peptides

ProtA: ACDERFGHIK -> ACDER FGHIK **3**

ProtB: LMNPRFGHIK -> LMNPR FGHIK  
shuffled

RandA: CHERIGDAEK -> CHER IGDAEK **5**

RandB: PLKMHFRIGN -> PLK MHFR IGN



## SM

### Reverse each peptide

**SAMPLER becomes SELPMAR**

Faster, skips digestion overhead

SM ignores peptides if decoy = target

from James Eddes

# Target-Decoy Thresholding – 1D Example

- Sort hits in descending score order
- Count target hits above chosen score
- Count decoy hits above chosen score
- For each score calculate FDR
  - $FDR = (2 * FP) / (FP + TP)$
- Filter list at score corresponding to chosen FDR

Rank	Xcorr	Sequence	Reference	Est. FP	Est TP	Est FN	Est TN	precision	sensitivity	accuracy
1	8.43	R.QLLLTADDRVNPICIGGVILFHETLYQK.A	IPI:PI00549682.2	0	1	2152	456	1.000	0.000	0.175
2	8.15	R.SGPFQIIFRPDNFVFGQSGAGNNWAK.G	IPI:PI00007752.1	0	2	2151	456	1.000	0.001	0.176
3	7.32	K.FEGEPATHIQPGVQLQSNTYDLQESNVR.L	IPI:PI00216139.1	0	3	2150	456	1.000	0.001	0.176
4	7.30	K.GHYTEGAEVLDSVLDVVR.K	IPI:PI00007752.1	0	4	2149	456	1.000	0.002	0.176
5	7.12	K.DYPVWVSIEDPFDQDDWGAWQK.F	IPI:PI00465248.4	0	5	2148	456	1.000	0.002	0.177
6	7.08	K.EKLCYVALDFEQEM*ATAASSSSLEK.S	IPI:PI00642997.1	0	6	2147	456	1.000	0.003	0.177
7	7.08	R.DLKPENLLDQGGYIQVTFDFGFAK.R	IPI:PI00217960.1	0	7	2146	456	1.000	0.003	0.177
8	7.08	R.SHSTEPGLVLTGGQDVGQLGLGENVM*ER.K	IPI:PI00001661.2	0	8	2145	456	1.000	0.004	0.178
9	7.07	R.TTGIVMDSGDGVHTHTVPIYEGYALPHAILR.L	IPI:PI00642997.1	0	9	2144	456	1.000	0.004	0.178
10	6.99	K.KFDLGGQVDFTGHALALYR.T	IPI:PI00514395.1	0	10	2143	456	1.000	0.005	0.179
741	3.64	K.VILALGDYM*GATCHACIGGTNVR.N	IPI:PI00328328.3	0	741	1412	456	1.000	0.344	0.459
742	3.64	K.EIDDSVLGGTGPYR.R	IPI:PI00306516.1	0	742	1411	456	1.000	0.345	0.459
743	3.64	K.DVTNNVHYENR.S	IPI:PI00033025.7	0	743	1410	456	1.000	0.345	0.460
744	3.64	K.EHWNPAIVALVYVNLK.A	IPI:PI00556528.1	0	744	1409	456	1.000	0.346	0.460
745	3.63	R.YISPDQLADLYK.G	IPI:PI00013769.1	0	745	1408	456	1.000	0.346	0.460
746	3.63	K.HLREIDRSLLEGRPEEDSDIEIYK.R	##IPI:PI00479822.2	2	744	1409	454	0.997	0.346	0.459
747	3.63	K.IPEGLFDPSPNVK.G	IPI:PI00216622.1	2	745	1408	454	0.997	0.346	0.460
748	3.63	K.IDLGVHVDGFIANVAHTFVVDVAGQTQVTR.K	IPI:PI00299000.3	2	746	1407	454	0.997	0.346	0.460
749	3.63	R.KGETIFVERPLVAAQFLWNALYR.Y	IPI:PI00013789.4	2	747	1406	454	0.997	0.347	0.460
750	3.62	K.SM*AASGNLGHPTFVDEL.-	IPI:PI00305383.1	2	748	1405	454	0.997	0.347	0.461
1251	2.93	K.ELAEQLGLSTGEK.E	IPI:PI00290460.3	14	1237	916	442	0.989	0.575	0.644
1252	2.92	R.IDAVNAETIR.E	IPI:PI00289535.6	14	1238	915	442	0.989	0.575	0.644
1253	2.92	R.LFLLPQAHIM*GTIGAIQSALAPPDSSTLLK.L	##IPI:PI00444490.1	16	1237	916	440	0.987	0.575	0.643
1254	2.92	R.FLLSLPEHR.G	IPI:PI00014151.3	16	1238	915	440	0.987	0.575	0.643
1255	2.92	R.SVDPDSPAESGLR.A	IPI:PI00003527.4	16	1239	914	440	0.987	0.575	0.644
1256	2.92	R.M*TDLLPNPAM*KGAPHAQASLR.R	##IPI:PI00658084.1	18	1238	915	438	0.986	0.575	0.642
1257	2.92	R.VGQAVALR.A	IPI:PI00012835.1	18	1239	914	438	0.986	0.575	0.643
1258	2.92	K.M*LLYTEVTR.Y	IPI:PI00010154.3	18	1240	913	438	0.986	0.576	0.643
1259	2.92	R.LIVENLSSR.V	IPI:PI00555857.1	18	1241	912	438	0.986	0.576	0.644
1260	2.91	K.WLNENAVEK.V	IPI:PI00012007.5	18	1242	911	438	0.986	0.577	0.644
1801	2.27	R.FM*EQVIFK.Y	IPI:PI00028091.1	76	1725	428	380	0.958	0.801	0.807
1802	2.27	K.EIRRVLEVLMEAM*K.A	##IPI:PI00020915.1	78	1724	429	378	0.957	0.801	0.806
1803	2.26	R.M*VLAAGVHEHQQLDLDAQK.H	IPI:PI00013847.3	78	1725	428	378	0.957	0.801	0.806
1804	2.26	R.LEODEYALR.S	IPI:PI00022793.4	78	1726	427	378	0.957	0.802	0.806
1805	2.26	R.ALKIPQVRSTDM*JAAR.R	##IPI:PI00457056.1	80	1725	428	376	0.956	0.801	0.805
1806	2.26	R.GGPNITLADIVKDPVSR.T	IPI:PI00016613.2	80	1726	427	376	0.956	0.802	0.806
1807	2.26	R.TQGVSSTDLVGR.M	IPI:PI00015285.6	80	1727	426	376	0.956	0.802	0.806
1808	2.26	K.HFVLDEC DK.M	IPI:PI00644431.1	80	1728	425	376	0.956	0.803	0.806
1809	2.26	K.KGEEALFTTRE	IPI:PI00019004.1	80	1729	424	376	0.956	0.803	0.807
1810	2.26	R.ISPELSAEAM*GSGTSVWNR.R	##IPI:PI00027003.1	82	1728	425	374	0.955	0.803	0.806
2101	1.94	R.CAHVARTYSIGR.S	IPI:PI000396391.7	162	1939	214	294	0.923	0.901	0.856
2102	1.94	R.VKIERQEVSK.E	##IPI:PI00028214.3	164	1938	215	292	0.922	0.900	0.855
2103	1.93	R.VSPQASVLEHR.E	IPI:PI00167300.3	164	1939	214	292	0.922	0.901	0.855
2104	1.93	K.FCLDNGAK.S	IPI:PI00169383.2	164	1940	213	292	0.922	0.901	0.856
2105	1.93	R.VQAVRHLRSLSEAEATR.K	##IPI:PI00305267.2	166	1939	214	290	0.921	0.901	0.854
2106	1.93	K.EKEKASWSSLSMDEK.V	IPI:PI00006579.1	166	1940	213	290	0.921	0.901	0.855
2107	1.93	R.YEEIVK.E	IPI:PI00180730.1	166	1941	212	290	0.921	0.902	0.855
2108	1.93	K.LYENFISEFEHR.V	IPI:PI00549672.2	166	1942	211	290	0.921	0.902	0.856
2109	1.93	K.ALAPTWEQLALGLEHSETVK.I	IPI:PI00171438.2	166	1943	210	290	0.921	0.902	0.856
2110	1.93	R.CEPIAM*TVPR.K	IPI:PI00010133.1	166	1944	209	290	0.921	0.903	0.856
2581	1.04	R.KSSLVTSK.L	IPI:PI00644717.1	430	2151	2	26	0.833	0.999	0.834
2582	1.03	-.M*DWCCSM*PTK.G	IPI:PI00457291.1	430	2152	1	26	0.833	1.000	0.835
2583	1.02	K.YVPPRR.V	IPI:PI00304809.6	430	2153	0	26	0.834	1.000	0.835
2584	1.01	K.LDRHIIPTRTHLFLGK.L	##IPI:PI00478412.2	432	2152	1	24	0.833	1.000	0.834
2585	1.01	K.NLSYDERSISIM*K.V	IPI:PI00045496.1	432	2153	0	24	0.833	1.000	0.834
2586	0.99	R.GPGSLGRRL	##IPI:PI00552079.1	434	2152	1	22	0.832	1.000	0.833
2587	0.99	R.LSSQTRLPR.K	##IPI:PI00655766.1	436	2151	2	20	0.831	0.999	0.832
2588	0.98	R.APQRLKR.A	##IPI:PI00305305.3	438	2150	3	18	0.831	0.999	0.831
2589	0.97	R.LSVLQVRRAR.G	##IPI:PI00456163.1	440	2149	4	16	0.830	0.998	0.830
2590	0.91	R.TIPLTFK.D	IPI:PI00002841.1	440	2150	3	16	0.830	0.999	0.830

score ↓

$$FDR = \frac{n_{decoy}}{n_{decoy} + n_{forward}} \times 2$$

Decoy hits (reversed database)

$$n_{reversed} = 13$$

$$n_{forward} = 2590 - 13 = 2577$$

$$FDR = \frac{13}{13 + 2577} \times 2 = 0.0100386 \approx 1\%$$

Target/Decoy FDR Tutorial  
Elias & Gygi, *Nature Methods*, 4, 207-214, 2007.

# Target-Decoy Thresholding – 3D to 7D possible in SM

In SM to achieve FDR goal:

1. Filter - sequence length range (min, max)
2. Filter - min backbone cleavage score threshold is applied.
3. Filter - limit by modifications (i.e. phospho sty or not)
4. Filter - precursor mass error distributions calculated for each LC-MS/MS run, retain PSMs within 3 stddev of mean
- 3D** 5. Optimize – thresholds for score, delta Rank1 – Rank2; separately for each precursor charge
  - a. enforces the notion that a credible ID should have both a suitable absolute score and be separated from the rank2 match by a suitable margin

**Example: for z=2 PSMs only**

delta Rank1 – Rank2 (rows)

Score(columns)

Optimum: 15,010 PSMs @ FDR 1.0

score  
dR1-R2

4.4  
1.0

	8.0	7.9	7.8	7.7	7.6	7.5	7.4	7.3	7.2	7.1	7.0	6.9	6.8	6.7	6.6	6.5	6.4	6.3	6.2	6.1	6.0	5.9	5.8	5.7	5.6	5.5	5.4	5.3	5.2	5.1	5.0	4.9	4.8	4.7	4.6	4.5	4.4	4.3	4.2	4.1
1.1	0.53 12904	0.54 12975	0.57 13038	0.59 13118	0.64 13189	0.63 13266	0.63 13343	0.64 13415	0.64 13486	0.65 13551	0.66 13612	0.66 13682	0.65 13748	0.65 13812	0.68 13856	0.69 13917	0.70 13984	0.70 14034	0.71 14090	0.74 14142	0.75 14187	0.74 14237	0.74 14288	0.75 14346	0.75 14379	0.75 14418	0.75 14464	0.77 14514	0.78 14553	0.78 14582	0.78 14630	0.78 14665	0.80 14705	0.83 14741	0.83 14781	0.84 14809	0.85 14834	0.85 14858	0.85 14884	0.87 14908
1.0	0.63 13030	0.64 13104	0.67 13170	0.69 13253	0.74 13327	0.73 13404	0.73 13483	0.74 13560	0.73 13632	0.76 13700	0.77 13761	0.77 13835	0.76 13902	0.76 13966	0.79 14011	0.80 14072	0.81 14141	0.80 14192	0.81 14249	0.85 14304	0.86 14350	0.86 14401	0.86 14453	0.87 14512	0.87 14548	0.86 14587	0.86 14634	0.89 14686	0.91 14727	0.92 14757	0.92 14805	0.93 14841	0.95 14881	0.98 14917	0.98 14957	0.99 14985	1.00 15010	1.01 15035	1.01 15061	1.03 15088
0.9	0.66 13070	0.67 13147	0.70 13214	0.72 13298	0.76 13372	0.76 13450	0.75 13530	0.76 13607	0.76 13679	0.79 13749	0.80 13810	0.79 13884	0.79 13952	0.78 14017	0.81 14064	0.82 14127	0.83 14197	0.83 14248	0.84 14306	0.88 14361	0.89 14407	0.90 14460	0.90 14512	0.91 14573	0.90 14610	0.90 14649	0.90 14696	0.92 14749	0.95 14790	0.96 14820	0.96 14869	0.97 14906	0.99 14946	1.01 14983	1.01 15027	1.02 15055	1.03 15080	1.05 15105	1.04 15132	1.07 15160
0.8	0.72 13108	0.73 13185	0.75 13252	0.78 13336	0.82 13410	0.82 13488	0.81 13569	0.82 13646	0.82 13718	0.84 13789	0.85 13850	0.85 13924	0.84 13993	0.84 14058	0.88 14105	0.88 14168	0.90 14239	0.90 14291	0.91 14349	0.94 14405	0.95 14451	0.97 14506	0.96 14559	0.97 14620	0.97 14658	0.97 14697	0.96 14744	0.99 14798	1.02 14842	1.04 14873	1.03 14922	1.04 14960	1.07 15000	1.09 15038	1.10 15084	1.11 15113	1.12 15138	1.13 15163	1.13 15190	1.16 15219
0.7	0.73 13141	0.74 13218	0.77 13285	0.79 13370	0.85 13445	0.84 13525	0.85 13607	0.86 13684	0.86 13756	0.88 13828	0.89 13891	0.89 13966	0.88 14035	0.88 14100	0.90 14147	0.91 14211	0.95 14284	0.95 14336	0.96 14394	1.00 14450	1.01 14496	1.02 14552	1.01 14605	1.02 14666	1.03 14706	1.03 14745	1.03 14792	1.05 14847	1.09 14891	1.10 14922	1.10 14972	1.11 15010	1.13 15050	1.17 15089	1.18 15135	1.19 15164	1.22 15191	1.24 15216	1.23 15244	1.27 15274
0.6	0.79 13166	0.80 13243	0.83 13310	0.85 13395	0.91 13470	0.90 13551	0.91 13635	0.92 13712	0.91 13784	0.94 13858	0.95 13921	0.94 13997	0.94 14066	0.93 14131	0.96 14178	1.01 14242	1.01 14315	1.01 14367	1.05 14425	1.06 14481	1.07 14527	1.07 14583	1.07 14637	1.07 14699	1.09 14739	1.10 14780	1.09 14827	1.12 14882	1.15 14928	1.16 14959	1.16 15011	1.17 15049	1.19 15090	1.23 15129	1.24 15175	1.25 15204	1.29 15231	1.31 15258	1.31 15286	1.34 15316
0.5	0.85 13188	0.86 13267	0.88 13335	0.91 13421	0.96 13496	0.96 13579	0.97 13664	0.98 13741	0.97 13813	0.99 13887	1.00 13951	1.00 14027	1.00 14096	0.99 14161	1.01 14209	1.02 14273	1.06 14346	1.06 14398	1.07 14456	1.10 14512	1.13 14559	1.14 14615	1.13 14669	1.14 14731	1.14 14771	1.15 14812	1.16 14859	1.18 14914	1.22 14960	1.23 14992	1.22 15044	1.25 15084	1.27 15125	1.31 15164	1.31 15210	1.34 15240	1.38 15267	1.40 15294	1.41 15324	1.45 15355
0.4	0.91 13242	0.92 13321	0.94 13391	0.96 13479	1.02 13554	1.03 13638	1.03 13724	1.04 13802	1.04 13874	1.06 13949	1.07 14014	1.08 14092	1.10 14163	1.11 14230	1.13 14278	1.14 14343	1.18 14416	1.17 14469	1.18 14529	1.22 14586	1.24 14634	1.25 14691	1.25 14746	1.26 14809	1.27 14850	1.29 14893	1.31 14943	1.33 14998	1.37 15045	1.38 15078	1.39 15133	1.41 15174	1.43 15215	1.47 15254	1.49 15301	1.51 15333	1.55 15360	1.57 15389	1.58 15420	1.63 15453
0.3	0.95 13278	0.96 13357	0.98 13427	1.01 13515	1.06 13591	1.07 13675	1.08 13761	1.08 13840	1.08 13913	1.10 13988	1.11 14053	1.12 14131	1.14 14202	1.15 14271	1.17 14319	1.18 14385	1.22 14458	1.21 14511	1.22 14571	1.26 14628	1.29 14679	1.33 14738	1.32 14793	1.33 14856	1.34 14897	1.37 14940	1.39 14990	1.41 15045	1.44 15093	1.45 15122	1.46 15181	1.48 15222	1.51 15263	1.54 15303	1.56 15350	1.59 15382	1.62 15409	1.65 15439	1.65 15470	1.70 15503
0.2	0.99 13302	1.00 13381	1.03 13452	1.05 13540	1.10 13617	1.11 13701	1.12 13787	1.12 13867	1.12 13940	1.14 14015	1.16 14081	1.17 14161	1.19 14232	1.22 14302	1.24 14350	1.25 14416	1.28 14490	1.29 14543	1.32 14603	1.36 14660	1.39 14711	1.39 14770	1.39 14825	1.41 14889	1.42 14930	1.44 14973	1.46 15024	1.49 15079	1.52 15127	1.53 15160	1.54 15215	1.56 15256	1.58 15297	1.62 15337	1.64 15384	1.66 15416	1.70 15445	1.72 15475	1.73 15506	1.78 15539

↑ (score, dR1-R2: with FDR > 1.0)

FDR  
# PSMs

HLA class I: B5703\_20180828

# Machine Learning to Combine Multiple Scores- Percolator

Not yet done with SM



© The Author(s), 2016. This article is published with open access at Springerlink.com

J. Am. Soc. Mass Spectrom. (2016) 27:1719–1727

DOI: 10.1007/s13361-016-1460-7



FOCUS: BIOINFORMATICS, SOFTWARE, AND  
MS-BASED "OMICS": RESEARCH ARTICLE

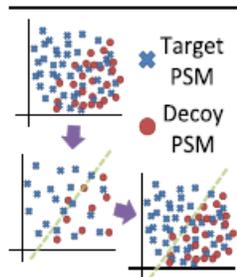
## Fast and Accurate Protein False Discovery Rates on Large-Scale Proteomics Data Sets with Percolator 3.0

Matthew The,<sup>1</sup> Michael J. MacCoss,<sup>2</sup> William S. Noble,<sup>2,3</sup> Lukas Käll<sup>1</sup>

<sup>1</sup>Science for Life Laboratory, School of Biotechnology, KTH – Royal Institute of Technology, Box 1031, 17121, Solna, Sweden

<sup>2</sup>Department of Genome Sciences, School of Medicine, University of Washington, Seattle, WA 98195, USA

<sup>3</sup>Department of Computer Science and Engineering, University of Washington, Seattle, WA 98195, USA



**Abstract.** Percolator is a widely used software tool that increases yield in shotgun proteomics experiments and assigns reliable statistical confidence measures, such as  $q$  values and posterior error probabilities, to peptides and peptide-spectrum matches (PSMs) from such experiments. Percolator's processing speed has been sufficient for typical data sets consisting of hundreds of thousands of PSMs. With our new scalable approach, we can now also analyze millions of PSMs in a matter of minutes on a commodity computer. Furthermore, with the increasing awareness for the need for reliable statistics on the protein level, we compared several easy-to-understand protein inference methods and implemented the best-performing method—grouping proteins by their corresponding sets of theoretical peptides and

then considering only the best-scoring peptide for each protein—in the Percolator package. We used Percolator 3.0 to analyze the data from a recent study of the draft human proteome containing 25 million spectra (PM:24870542). The source code and Ubuntu, Windows, MacOS, and Fedora binary packages are available from <http://percolator.ms/> under an Apache 2.0 license.

### SM scores/features to percolate

#### SM scores

##### Score

score (w/o penalty portion)

penalty portion

Delta Rank1-Rank2 score

SPI – percent scored peak intensity

BCS – backbone cleavage score

Sequence coverage

Fragmentation category

# sequences past parent filter (DB size dependency)

Parent mass error

Sequence length

#### SM Spectral Features

$z$

MH+

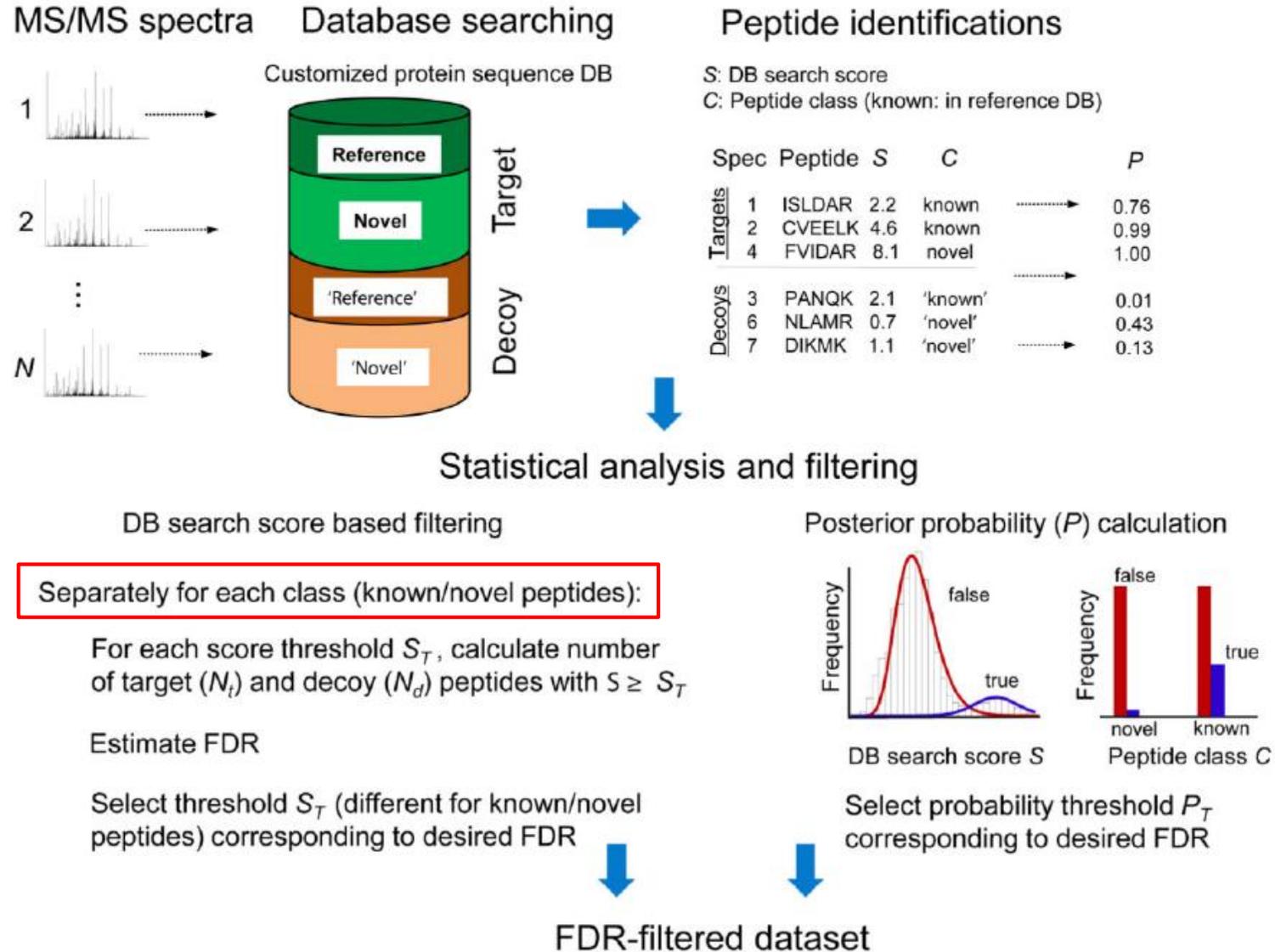
RT – retention time

Max sequence tag length

PIP – precursor ion purity

Precursor average  $\chi^2$

# Calculate FDR separately for rare peptide subsets



**Collapsing Peptide Spectrum Matches (PSM's) for Quantitation at  
the  
Protein level  
and Phosphosite level**

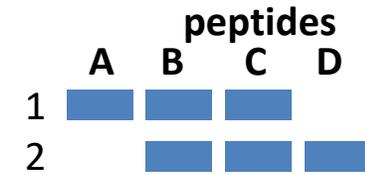
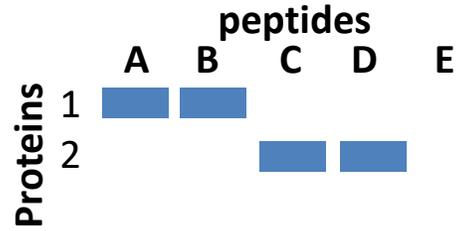
# Protein Grouping and Shared peptides

# groups    # subgroups

2

1, 1

2  
Distinct  
proteins



2  
differentiable  
proteins

# groups    # subgroups

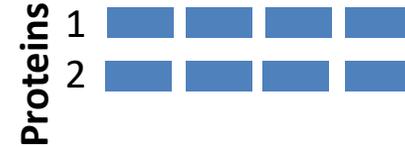
1

2

1

1

2  
Indistinguishable  
proteins



subset  
Protein  
#2

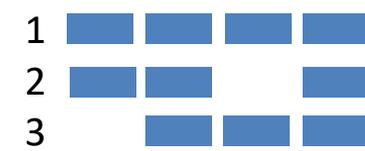
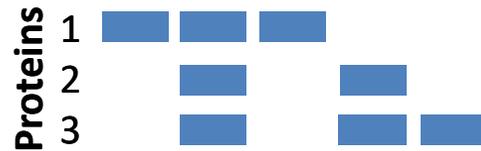
1

1

1

2

Subsumed  
Protein  
#2



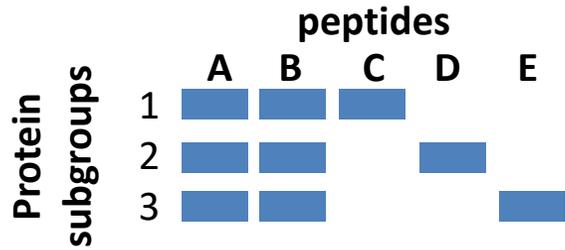
A group  
with some  
Shared peptides

1

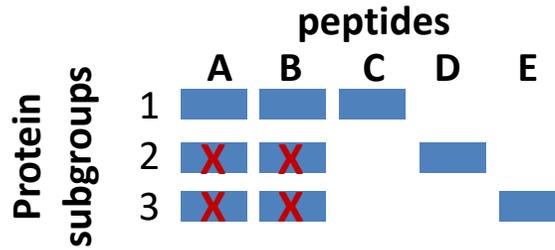
1

Protein Inference Tutorial  
Nesvizhskii, *Mol Cell Proteomics*, 4, 1419-1440, 2005.

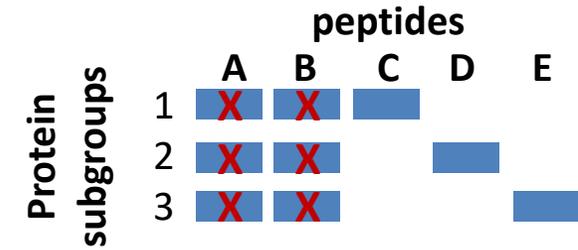
# Shared Peptides and Protein Grouping



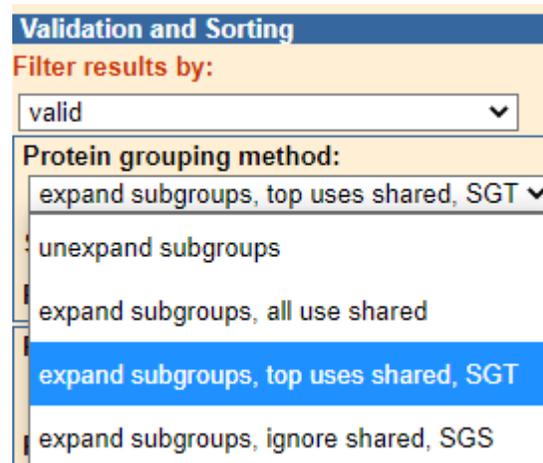
All subgroups used shared



Top subgroups uses shared  
Subgroup Top (SGT)



Ignore shared  
Subgroup Specific (SGS)



Ranking of subgroups is by protein score.  
Protein score is  $\Sigma$  distinct peptide scores.

# Protein Groups and Subgroups, Genes

#1.3																luad-v1.0-proteome-ratio-norm-NArm.gct																								
12205		215		16		26																																		
id	id	id.description	geneSymbol	numColumns	numSpecTraProtei	protein_mw	percentCoverage	numPepsUnique	scoreUnique	species	orfCategory	accession_number	accession_number	subgroup	entry_name	GeneSymbol	C3N.0	C3L.0																						
Sample.ID	na	na	na	na	na	na	na	na	na	na	na	na	na	na	na	na	C3N.0	C3N.0	C3L.0	C3L.0	C3N.0	C3N.0	C3N.0	C3N.0	C3L.0	C3L.0	C3N.0	C3N.0	C3N.0	C3N.0	C3L.0	C3L.0	C3N.0	C3N.0	C3N.0	C3N.0	C3L.0	C3L.0		
Experiment	na	na	na	na	na	na	na	na	na	na	na	na	na	na	na	na	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	
Type	na	na	na	na	na	na	na	na	na	na	na	na	na	na	na	na	Tumor	Norm																						
Smoking.Stage	na	na	na	na	na	na	na	na	na	na	na	na	na	na	na	na	smoke	smoke	smoke	smoke	non-s	non-s																		
Region.of.Origin	na	na	na	na	na	na	na	na	na	na	na	na	na	na	na	na	weste	weste	weste	weste	asian	asian	asian	asian	asian	weste	weste	asian	weste											
Age	na	na	na	na	na	na	na	na	na	na	na	na	na	na	na	na	81.00	81.00	58.00	58.00	48.00	48.00	70.00	70.00	45.00	45.00	79.00	79.00	64.00	64.00	48.00	48.00	59.00	59.00	71.00					
Gender	na	na	na	na	na	na	na	na	na	na	na	na	na	na	na	na	male	male	femal	femal	femal	femal	male	femal	femal	male														
Weight.kg	na	na	na	na	na	na	na	na	na	na	na	na	na	na	na	na	82.00	82.00	52.90	52.90	44.00	44.00	79.00	79.00	60.00	60.00	80.46	80.46	48.00	48.00	64.00	64.00	50.00	50.00	86.82					
BMI	na	na	na	na	na	na	na	na	na	na	na	na	na	na	na	na	27.37	27.37	22.60	22.60	18.08	18.08	27.99	27.99	20.05	20.05	22.77	22.77	18.75	18.75	22.15	22.15	21.93	21.93	29.10					
Cigarettes.per.Day	na	na	na	na	na	na	na	na	na	na	na	na	na	na	na	na	5.00	5.00	20.00	20.00	NA	NA	40.00	40.00	NA	NA	10.00	10.00	NA	NA	5.00	5.00	NA	NA	20.00					
Smoking.History	na	na	na	na	na	na	na	na	na	na	na	na	na	na	na	na	Currer	Currer	Currer	Currer	Lifelo	Lifelo																		
QC.status	na	na	na	na	na	na	na	na	na	na	na	na	na	na	na	na	QC.pa																							
normalization	na	na	na	na	na	na	na	na	na	na	na	na	na	na	na	na	2comp	2comp																						
NP_001317422.1	NP_001317	protein O	OSCP1	25	192	44,643	75.5	25	320.89	HOMO SAPIENS		NP_001317	NP_001317	4385.1	protein O	OSCP1	-2.45	2.09	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-			
NP_001106273.1	NP_001106	sodium/c	SLC8A1	23	148	105,600	31.3	22	320.75	HOMO SAPIENS		NP_001106	NP_001106	4386.1	sodium/c	SLC8A1	-1.27	-0.12	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-		
NP_620156.1	NP_620156	aldose 1-	GALM	25	709	37,993	66	20	320.67	HOMO SAPIENS		NP_620156	NP_620156	4387.1	aldose 1-	GALM	-0.64	-0.53	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-		
NP_001446.1	NP_001446	forkhead	FOXO3	25	216	71,561	46.8	21	320.63	HOMO SAPIENS		NP_001446	NP_001446	4388.1	forkhead	FOXO3	0.19	1.17	-0.68	0.88	-0.22	1.34	-0.35	1.46	-1.19	1.61	0.90	0.86	-0.62	1.23	-0.25	0.57	-0.63	1.22	0.72					
NP_002006.2	NP_002006	forkhead	FOXO1	25	198	70,060	38.1	18	280.85	HOMO SAPIENS		NP_002006	NP_002006	4388.2	forkhead	FOXO1	0.29	0.79	-0.45	0.21	0.15	2.48	-0.05	1.49	-0.14	2.61	0.30	0.26	-0.77	-0.04	-1.15	1.10	-1.59	2.18	-0.70					
NP_005929.2	NP_005929	forkhead	FOXO4	25	83	53,969	19.2	8	96.07	HOMO SAPIENS		NP_005929	NP_005929	4388.3	forkhead	FOXO4	-0.41	0.19	-1.30	0.47	-1.02	2.75	-1.74	1.74	1.29	2.84	-0.89	-0.16	-1.24	-1.88	-1.59	2.43	-1.57	2.09	-0.25					
NP_001278210.2	NP_001278	forkhead	FOXO6	24	79	57,486	18.4	8	92.14	HOMO SAPIENS		NP_001278	NP_001278	4388.4	forkhead	FOXO6	1.42	1.00	-0.97	0.53	0.71	1.37	-0.25	1.11	-0.90	2.78	1.88	0.26	-0.77	-1.96	-2.02	0.38	-1.19	1.93	0.21					
NP_057250.1	NP_057250	E3 SUMO-	PIAS1	25	159	72,805	46.3	23	320.61	HOMO SAPIENS		NP_057250	NP_057250	4389.1	E3 SUMO-	PIAS1	0.33	0.64	0.50	0.87	0.08	0.89	-0.34	0.77	0.56	1.16	0.74	0.60	-1.13	1.59	-0.08	0.85	-0.23	0.76	0.67					
NP_006090.2	NP_006090	E3 SUMO-	PIAS3	20	85	69,043	35.8	18	191.87	HOMO SAPIENS		NP_006090	NP_006090	4389.2	E3 SUMO-	PIAS3	NA	NA																						
NP_001243764.1	NP_001243	queuine	QTRT2	25	258	49,116	60.8	19	320.61	HOMO SAPIENS		NP_001243	NP_001243	4390.1	queuine	QTRT2	0.12	-0.29	0.42	-0.45	0.44	-0.36	0.82	-0.35	0.59	-0.33	-0.16	0.27	0.34	0.11	0.41	-0.51	-0.23	-0.39	0.07					
NP_001243463.1	NP_001243	tricarbox	SLC25A1	25	669	35,394	55.6	17	317.88	HOMO SAPIENS		NP_001243	NP_001243	4429.1	tricarbox	SLC25A1	0.87	-1.26	1.24	0.29	0.12	-1.44	-0.91	-1.49	3.86	-0.92	-2.61	-2.94	0.43	-0.23	-1.59	-1.15	-0.17	-0.98	-1.50					
NP_001274316.1	NP_001274	tricarbox	SLC25A1	25	500	23,438	66.3	13	251.51	HOMO SAPIENS		NP_001274	NP_001274	4429.2	tricarbox	SLC25A1	0.95	-1.26	1.51	0.40	0.11	-1.33	-0.91	-0.91	4.73	-1.11	-2.61	-3.07	0.16	-0.23	-1.59	-1.47	0.02	-1.16	-1.72					
NP_006811.2	NP_006811	interferon	IFI44L	24	206	52,006	63.9	23	317.83	HOMO SAPIENS		NP_006811	NP_006811	4430.1	interferon	IFI44L	-5.89	0.6																						
NP_003013.1	NP_003013	SH3 dom	SH3BGRL	25	1131	12,774	91.2	16	317.77	HOMO SAPIENS		NP_003013	NP_003013	4431.1	SH3 dom	SH3BGRL	-0.17	1.3																						
NP_116193.2	NP_116193	2-aminoe	ADO	25	355	30,093	78.5	21	317.74	HOMO SAPIENS		NP_116193	NP_116193	4432.1	2-aminoe	ADO	-0.36	0.2																						

Protein group with multiple genes

4388.1 forkhead FOXO3  
 4388.2 forkhead FOXO1  
 4388.3 forkhead FOXO4  
 4388.4 forkhead FOXO6

Protein group with isoforms of 1 gene

4429.1 tricarbox SLC25A1  
 4429.2 tricarbox SLC25A1



# Link from Group Number in Protein/Peptide Summary HTML Results

## Highlights Identified Peptides

Subgroup-specific Peptides

- 1
- 2
- 3

ID	Sub Group #	Length	Identical AA's	%ID	Species	Protein Name
<a href="#">54696354</a>	1	327	327	100.0	Homo sapiens	protein phosphatase 1, catalytic subunit, beta isoform
<a href="#">46249376</a>	1	327	326	99.7	Homo sapiens	protein phosphatase 1, catalytic subunit, beta isoform 1
<a href="#">45827798</a>	2	286	254	88.8	Homo sapiens	protein phosphatase 1, catalytic subunit, alpha isoform 2
<a href="#">4506003</a>	2	330	291	88.2	Homo sapiens	protein phosphatase 1, catalytic subunit, alpha isoform 1
<a href="#">56790945</a>	2	341	290	85.0	Homo sapiens	protein phosphatase 1, catalytic subunit, alpha isoform 3
<a href="#">5668560</a>	3	294	269	91.5	Homo sapiens	serine/threonine phosphatase 1 gamma
<a href="#">4506007</a>	3	323	285	88.2	Homo sapiens	protein phosphatase 1, catalytic subunit, gamma isoform
<a href="#">54696354</a>	(1)	MADG-ELNVDSLITRLLLEVR-----GCRPGKIVQMTAEAVRGLCLIKSREIFLSQPIPLELEAPLKICGDIHGQYTDLLRRLFYEGGFPPEANYLF				
<a href="#">46249376</a>	(1)	MADG-ELNVDSLITRLLLEVR-----GCRPGKIVQMTAEAVRGLCLIKSREIFLSQPIPLELEAPLKICGDIHGQYTDLLRRLFYEGGFPPEANYLF				
<a href="#">45827798</a>	(2)	MSDSIKLNLDIIIGRLLEVC-----DIIHGQYDILLRRLFYEGGFPPEANYLF				
<a href="#">4506003</a>	(2)	MSDSIKLNLDIIIGRLLEVC-----GSRPGKKNVQLTENEIRGLCLIKSREIFLSQPIPLELEAPLKICGDIHGQYDILLRRLFYEGGFPPEANYLF				
<a href="#">56790945</a>	(2)	MSDSIKLNLDIIIGRLLEVC-----GSRPGKKNVQLTENEIRGLCLIKSREIFLSQPIPLELEAPLKICGDIHGQYDILLRRLFYEGGFPPEANYLF				
<a href="#">5668560</a>	(3)	MADLDKLNIDSIIQRLLEVR-----GSKPGKKNVQLQENEIRGLCLIKSREIFLSQPIPLELEAPLKICGDIHGQYDILLRRLFYEGGFPPEANYLF				
<a href="#">4506007</a>	(3)	MADLDKLNIDSIIQRLLEVR-----GSKPGKKNVQLQENEIRGLCLIKSREIFLSQPIPLELEAPLKICGDIHGQYDILLRRLFYEGGFPPEANYLF				
<a href="#">54696354</a>	(1)	LGDYVDRGRKQSLETICLLLAYKIKYPENFFLLRCNMHECASINRIYGFYDECKRRFNKIKLWKTFTDCFNCLPIAAIVDEKIFCCHGGLSPDLQSMEQIRRI				
<a href="#">46249376</a>	(1)	LGDYVDRGRKQSLETICLLLAYKIKYPENFFLLRCNMHECASINRIYGFYDECKRRFNKIKLWKTFTDCFNCLPIAAIVDEKIFCCHGGLSPDLQSMEQIRRI				
<a href="#">45827798</a>	(2)	LGDYVDRGRKQSLETICLLLAYKIKYPENFFLLRCNMHECASINRIYGFYDECKRRYNIKIKLWKTFTDCFNCLPIAAIVDEKIFCCHGGLSPDLQSMEQIRRI				
<a href="#">4506003</a>	(2)	LGDYVDRGRKQSLETICLLLAYKIKYPENFFLLRCNMHECASINRIYGFYDECKRRYNIKIKLWKTFTDCFNCLPIAAIVDEKIFCCHGGLSPDLQSMEQIRRI				
<a href="#">56790945</a>	(2)	LGDYVDRGRKQSLETICLLLAYKIKYPENFFLLRCNMHECASINRIYGFYDECKRRYNIKIKLWKTFTDCFNCLPIAAIVDEKIFCCHGGLSPDLQSMEQIRRI				
<a href="#">5668560</a>	(3)	LGDYVDRGRKQSLETICLLLAYKIKYPENFFLLRCNMHECASINRIYGFYDECKRRYNIKIKLWKTFTDCFNCLPIAAIVDEKIFCCHGGLSPDLQSMEQIRRI				
<a href="#">4506007</a>	(3)	LGDYVDRGRKQSLETICLLLAYKIKYPENFFLLRCNMHECASINRIYGFYDECKRRYNIKIKLWKTFTDCFNCLPIAAIVDEKIFCCHGGLSPDLQSMEQIRRI				
<a href="#">54696354</a>	(1)	MRPTDVPDTGLLCDLLWSDPKDQVQGWGENDRGVSFTFGADVVSFKFLNRHDLDLICRAHQVVEDGYEFFAKRQLVTLFSAPNYCGEFDNAGGMMSVDETL				
<a href="#">46249376</a>	(1)	MRPTDVPDTGLLCDLLWSDPKDQVQGWGENDRGVSFTFGADVVSFKFLNRHDLDLICRAHQVVEDGYEFFAKRQLVTLFSAPNYCGEFDNAGGMMSVDETL				
<a href="#">45827798</a>	(2)	MRPTDVPDQGLLCDLLWSDPKDQVQGWGENDRGVSFTFGAEVVAKFLHKHDLDLICRAHQVVEDGYEFFAKRQLVTLFSAPNYCGEFDNAGAMMSVDETL				
<a href="#">4506003</a>	(2)	MRPTDVPDQGLLCDLLWSDPKDQVQGWGENDRGVSFTFGAEVVAKFLHKHDLDLICRAHQVVEDGYEFFAKRQLVTLFSAPNYCGEFDNAGAMMSVDETL				
<a href="#">56790945</a>	(2)	MRPTDVPDQGLLCDLLWSDPKDQVQGWGENDRGVSFTFGAEVVAKFLHKHDLDLICRAHQVVEDGYEFFAKRQLVTLFSAPNYCGEFDNAGAMMSVDETL				
<a href="#">5668560</a>	(3)	MRPTDVPDQGLLCDLLWSDPKDQVQGWGENDRGVSFTFGAEVVAKFLHKHDLDLICRAHQVVEDGYEFFAKRQLVTLFSAPNYCGEFDNAGAMMSVDETL				
<a href="#">4506007</a>	(3)	MRPTDVPDQGLLCDLLWSDPKDQVQGWGENDRGVSFTFGAEVVAKFLHKHDLDLICRAHQVVEDGYEFFAKRQLVTLFSAPNYCGEFDNAGAMMSVDETL				
<a href="#">54696354</a>	(1)	MCSFQILKPSEK-KAKY-QYGGNS-GRPVTTPRT--ANPPKRR				
<a href="#">46249376</a>	(1)	MCSFQILKPSEK-KAKY-QYGGNS-GRPVTTPRT--ANPPKRR				
<a href="#">45827798</a>	(2)	MCSFQILKPADKKNKGYQFSGLNPGGRPITPPRN--SAKAKK-				
<a href="#">4506003</a>	(2)	MCSFQILKPADKKNKGYQFSGLNPGGRPITPPRN--SAKAKK-				
<a href="#">56790945</a>	(2)	MCSFQILKPADKKNKGYQFSGLNPGGRPITPPRN--SAKAKK-				
<a href="#">5668560</a>	(3)	MCSFQ-----				
<a href="#">4506007</a>	(3)	MCSFQILKPAEKKKPNA-----TRPVTTPRGMITKQAKK-				

ClustalW aligns proteins in group SM colors identified peptides and labels subgroups



# Collapsing PSM's to Phosphosite level – Show site-grouped PSMs - HTML

		#Am		#Loc		big		iTRAQ		iTRAQ		iTRAQ		Gene		PO4	
m/z	ppm	ID	Filename	VML score	sty sites	sty sites	114/1	115/1	116/1	PIP (%)	Sequence	VML sequence	Gene name	PO4 Site			
552.960	5.1	13.83	1050.10050.3	0.00	0	1	0.96	0.08	1.27	76.7	RSStLSQLPGDK	RS(0.33)S(0.33)T(0.33)LS(0.0)QLPGDK	SPAG9	T581t			
643.689	-0.1	10.84	7464.7464.3	0.00	0	1	0.93	0.06	1.45	94.5	KRSsTLSQLPGDK	KRS(0.25)S(0.25)T(0.25)LS(0.25)QLPGDK					
500.923	-1.4	8.06	558.11558.3	0.00	0	1	1.48	0.31	0.88	78.7	SsTLSQLPGDK	S(0.25)S(0.25)T(0.25)LS(0.25)QLPGDK					
750.882	0.8	7.23	174.11174.2	0.00	0	1	1.84	0.16	0.99	87.1	SStLSQLPGDK	S(0.0)S(0.50)T(0.50)LS(0.0)QLPGDK					
750.881	-0.2	7.09	572.11572.2	0.00	0	1	1.96	0.14	0.71	62.8	SStLSQLPGDK	S(0.0)S(0.50)T(0.50)LS(0.0)QLPGDK					
678.833	3.6	7.07	720.11720.2	1.04	0	1	0.99	0.10	0.96	82.9	SStLSQLPGDK	S(0.33)S(0.33)T(0.33)LS(0.0)QLPGDK					
551.943	0.1	10.81	1410.20410.3	4.43	1	0	0.40	0.70	1.84	75.1	DGGsVVGASVIFYK	DGGS(0.99)VVGAS(0.0)VIFY(0.0)K	SPAG9	S691s			
755.361	2.1	7.91	1717.20717.2	7.91	1	0	1.03	0.89	1.60	94.0	DGGsVVGASVIFYK	DGGS(0.99)VVGAS(0.0)VIFY(0.0)K					
827.408	-2.0	7.73	1423.20423.2	7.73	1	0	0.57	0.90	1.63	87.7	DGGsVVGASVIFYK	DGGS(0.99)VVGAS(0.0)VIFY(0.0)K					
721.043	0.7	16.86	757.16757.3	2.29	1	0	0.48	0.98	1.77	69.8	SAsQSSLDKLDQELK	S(0.0)AS(0.99)QS(0.0)S(0.0)LDKLDQELK	SPAG9	S716s			
645.813	1.0	12.89	4977.4977.2	2.66	1	0	0.74	0.08	1.45	91.6	SAsQSSLDK	S(0.0)AS(0.99)QS(0.0)S(0.0)LDK					
541.035	2.5	12.36	777.16777.4	0.00	0	1	0.42	1.08	1.74	74.7	SASQsSLDKLDQELK	S(0.0)AS(0.50)QS(0.50)S(0.0)LDKLDQELK					
1081.057	-2.6	13.71	1073.15073.2	1.21	1	0	0.44	0.76	1.56	87.9	SASQsSLDKLDQELK	S(0.0)AS(0.0)QS(0.0)S(0.99)LDKLDQELK	SPAG9	S719s			
541.032	-1.7	8.73	1034.15154.4	0.00	0	1	0.48	0.74	1.56	75.6	SASQsSLDKLDQELK	S(0.0)AS(0.0)QS(0.50)S(0.50)LDKLDQELK					
721.040	-2.5	7.37	1039.15039.3	0.00	0	1	0.35	0.75	1.52	72.3	SASQsSLDKLDQELK	S(0.0)AS(0.0)QS(0.50)S(0.50)LDKLDQELK					
430.878	0.6	4.40	4593.4593.3	0.00	0	1	0.65	0.06	1.52	69.7	SASQsSLDK	S(0.0)AS(0.0)QS(0.50)S(0.50)LDK					

Spectrum Mill - Protein/Peptide Summary - KarlTMT10\_131\_VMsite\_sty\_PIP50\_africa\_showAll

Spectrum Mill | Extractor | MS/MS Search | Autovalidation | Quality Metrics & FDR | Workflows | Tool Belt | Process Report | MRM Selector | Spectrum Summary | Help

Summarize | Save As... | Load...

Queue request  Excel export

Mode: Protein - Var Mod Site Comparison

Group columns by:  File  Directory

Filter by Proteogenomic Features: Off

Data directories: Select ...

CPTAC3pancancer/MEDUL/Phosphoproteome\_...

Search result files: \*.spp

Validation and Sorting

Filter results by: valid

Protein grouping method: unexpand subgroups

Sort proteins by: Score

Filter by protein score: > 0.0

Filter peptides by: Score: > 0 % SPI: > 0.0

Required AAs: s|t|y

Disallowed AAs: none

Accession #'s:

Review Fields

Filename  Sequence  b|y map  Precursor XIC Intensity

Score  Rev Sequence  Rank 2  DEQ ratios  Invert Median

Fwd-Rev score  VML sequence  Reporter Ratios: TMT 10

Rank 1-2 score  Prec Av Chi<sup>2</sup>  Isol Pur

SPI (%)  Ret time, width  Ion mobility

Backbone Cleavage Score  Precursor m/z  MH<sup>+</sup>

Glyco Product Ions Score  Delta mass  Pep pl

Unmatched ions  # Sequences per MH<sup>+</sup>

Var mod sites  Protein MW  Prot pl

VML score s|t|y  Species

Solution charge  Accession #

Start AA position  Protein name

Proteogenomic features

Category:

Spectrum Grouping Options

Group missed cleavages containing VM site(s)

Show all grouped spectra

Combine repeat site-level observations with non-conflicting localizations

- Different precursor charge states
- Different sample-handling modifications
- Different trypsin missed cleavage states

Representative peptide: if confident localization, higher ID scores and longer peptide preferred, else best VML score.

# Collapsing PSM's to Phosphosite level – Show site-grouped PSMs - SSV

PSMexportVMsitesPolished.1.ssv

Log file created by  
Autovalidation  
VM site Polishing

Combine repeat site-level observations with non-conflicting localizations

- Different precursor charge states
- Different sample-handling modifications
- Different trypsin missed cleavage states

Representative peptide: if confident localization, higher ID scores and longer peptide preferred, else best VML score.

Sequence overlap complications: Ambiguous localizations grouped with confident ones preferentially over other overlapping ambiguous localizations that are inconsistent with the confident localization.

number	filename	parent_charge	score	deltaForwardReverseScore	backbone_cleavage_score	missedCleavage	numKorR	scoreVML	numPossibleMSites	numActualMSites	numLocalizedMSites	numAmbiguousMSites	sequence	sequenceVML	accession_number
6442	05CPTAC_LUAD_	4	5.51	5.51	4	1	3	0	4	1	0	1	KKEsSmLATVK	KKEs(0.33)S(0.33)S(0.33)M(1.0)LAT(0.0)VK	NP_001123917.1
6442	12CPTAC_LUAD_	3	5.27	3.72	5	1	2	0	4	1	0	1	KESsSmLATVK	KES(0.33)S(0.33)S(0.33)M(1.0)LAT(0.0)VK	NP_001123917.1
6442	09CPTAC_LUAD_	3	4.92	4.92	5	1	2	0	4	1	0	1	KESsSMLATVK	KES(0.33)S(0.33)S(0.33)MLAT(0.0)VK	NP_001123917.1
6442	04CPTAC_LUAD_	3	4.71	4.71	5	1	2	0	4	1	0	1	KESsSmLATVK	KES(0.33)S(0.33)S(0.33)M(1.0)LAT(0.0)VK	NP_001123917.1
6442	04CPTAC_LUAD_	2	4.27	4.27	6	0	1	0	4	1	0	1	ESSsmLATVK	ES(0.0)S(0.50)S(0.50)M(1.0)LAT(0.0)VK	NP_001123917.1
6443	20CPTAC_LUAD_	2	16.95	15.07	11	0	1	1.127	5	1	1	0	EEVSGsSAAVTENADSDR	EEVS(0.0)GS(0.99)S(0.0)AAVT(0.0)ENADS(0.0)DR	NP_001123917.1
6443	10CPTAC_LUAD_	2	9.55	9.55	9	0	1	0	5	1	0	1	EEVSGsSAAVTENADSDR	EEVS(0.33)GS(0.33)S(0.33)AAVT(0.0)ENADS(0.0)DR	NP_001123917.1
6444	10CPTAC_LUAD_	4	21.91	21.91	17	1	2	0	9	1	0	1	EEVSGSSAAVTENADSDRI	EEVS(0.0)GS(0.0)S(0.0)AAVT(0.0)ENADS(0.0)DRIS	NP_001123917.1
6444	24CPTAC_LUAD_	4	18.74	22.07	16	1	2	0	9	1	0	1	EEVSGSSAAVTENADSDRI	EEVS(0.0)GS(0.0)S(0.0)AAVT(0.0)ENADS(0.0)DRIS	NP_001123917.1
6444	25CPTAC_LUAD_	4	12.42	12.42	12	1	2	0	9	1	0	1	EEVSGSSAAVTENADSDRI	EEVS(0.0)GS(0.0)S(0.0)AAVT(0.0)ENADS(0.0)DRIS	NP_001123917.1
6445	02CPTAC_LUAD_	3	12.26	12.26	13	0	1	0	6	1	0	1	LDSDFNIshSELENSSELK	LDS(0.0)DFNIS(0.50)S(0.50)HS(0.0)ELENS(0.0)S(0.0)	NP_001123917.1
6445	10CPTAC_LUAD_	3	10.67	10.67	11	0	1	0	6	1	0	1	LDSDFNIshSELENSSELK	LDS(0.0)DFNIS(0.33)S(0.33)HS(0.33)ELENS(0.0)S(0.0)	NP_001123917.1
6445	09CPTAC_LUAD_	3	9.17	9.17	11	0	1	0	6	1	0	1	LDSDFNIshSELENSSELK	LDS(0.0)DFNIS(0.50)S(0.50)HS(0.0)ELENS(0.0)S(0.0)	NP_001123917.1
6445	05CPTAC_LUAD_	3	8.79	8.79	11	0	1	0	6	1	0	1	LDSDFNIshSELENSSELK	LDS(0.0)DFNIS(0.33)S(0.33)HS(0.33)ELENS(0.0)S(0.0)	NP_001123917.1
6445	14CPTAC_LUAD_	3	8.05	8.05	12	0	1	0	6	1	0	1	LDSDFNIshSELENSSELK	LDS(0.0)DFNIS(0.50)S(0.50)HS(0.0)ELENS(0.0)S(0.0)	NP_001123917.1
6445	09CPTAC_LUAD_	3	7.46	7.46	9	0	1	0	6	1	0	1	LDSDFNIshSELENSSELK	LDS(0.0)DFNIS(0.33)S(0.33)HS(0.33)ELENS(0.0)S(0.0)	NP_001123917.1
6446	15CPTAC_LUAD_	3	11.45	8.48	6	0	1	5.573	3	1	1	0	sVHISTPEK	S(0.99)VHIS(0.0)T(0.0)PEK	NP_001123917.1
6446	10CPTAC_LUAD_	3	10.19	10.19	6	0	1	4.53	3	1	1	0	sVHISTPEK	S(0.99)VHIS(0.0)T(0.0)PEK	NP_001123917.1
6446	07CPTAC_LUAD_	3	6.34	6.34	6	0	1	5.149	3	1	1	0	sVHISTPEK	S(0.99)VHIS(0.0)T(0.0)PEK	NP_001123917.1
6447	09CPTAC_LUAD_	4	20.97	20.97	15	1	2	1.15	5	1	1	0	SVHISTPEKEPCAPLTIPSIR	S(0.0)VHIS(0.0)T(0.99)PEKEPCAPLT(0.0)IPS(0.0)IR	NP_001123917.1
6447	12CPTAC_LUAD_	4	20.31	18.94	15	1	2	1.135	5	1	1	0	SVHISTPEKEPCAPLTIPSIR	S(0.0)VHIS(0.0)T(0.99)PEKEPCAPLT(0.0)IPS(0.0)IR	NP_001123917.1
6447	03CPTAC_LUAD_	4	20.2	20.2	16	1	2	1.157	5	1	1	0	SVHISTPEKEPCAPLTIPSIR	S(0.0)VHIS(0.0)T(0.99)PEKEPCAPLT(0.0)IPS(0.0)IR	NP_001123917.1
6447	08CPTAC_LUAD_	4	19.99	19.99	14	1	2	1.177	5	1	1	0	SVHISTPEKEPCAPLTIPSIR	S(0.0)VHIS(0.0)T(0.99)PEKEPCAPLT(0.0)IPS(0.0)IR	NP_001123917.1
6447	10CPTAC_LUAD_	4	19.43	19.43	15	1	2	1.184	5	1	1	0	SVHISTPEKEPCAPLTIPSIR	S(0.0)VHIS(0.0)T(0.99)PEKEPCAPLT(0.0)IPS(0.0)IR	NP_001123917.1
6447	16CPTAC_LUAD_	4	19.02	19.02	15	1	2	1.221	5	1	1	0	SVHISTPEKEPCAPLTIPSIR	S(0.0)VHIS(0.0)T(0.99)PEKEPCAPLT(0.0)IPS(0.0)IR	NP_001123917.1
6447	14CPTAC_LUAD_	3	13.06	14.67	8	0	1	1.341	3	1	1	0	SVHISTPEK	S(0.0)VHIS(0.0)T(0.99)PEK	NP_001123917.1
6447	10CPTAC_LUAD_	3	12.01	12	8	0	1	1.253	3	1	1	0	SVHISTPEK	S(0.0)VHIS(0.0)T(0.99)PEK	NP_001123917.1
6447	17CPTAC_LUAD_	3	11.53	12.31	7	0	1	1.406	3	1	1	0	SVHISTPEK	S(0.0)VHIS(0.0)T(0.99)PEK	NP_001123917.1
6447	14CPTAC_LUAD_	4	20.65	20.65	14	1	2	0	5	1	0	1	SVHISTPEKEPCAPLTIPSIR	S(0.0)VHIS(0.50)T(0.50)PEKEPCAPLT(0.0)IPS(0.0)IR	NP_001123917.1
6447	04CPTAC_LUAD_	4	19.99	19.99	15	1	2	0	5	1	0	1	SVHISTPEKEPCAPLTIPSIR	S(0.0)VHIS(0.50)T(0.50)PEKEPCAPLT(0.0)IPS(0.0)IR	NP_001123917.1
6447	15CPTAC_LUAD_	4	17.04	14.84	14	1	2	0	5	1	0	1	SVHISTPEKEPCAPLTIPSIR	S(0.0)VHIS(0.50)T(0.50)PEKEPCAPLT(0.0)IPS(0.0)IR	NP_001123917.1
6447	13CPTAC_LUAD_	4	15.89	15.89	12	1	2	0	5	1	0	1	SVHISTPEKEPCAPLTIPSIR	S(0.0)VHIS(0.50)T(0.50)PEKEPCAPLT(0.0)IPS(0.0)IR	NP_001123917.1
6447	05CPTAC_LUAD_	4	15	15	13	1	2	0	5	1	0	1	SVHISTPEKEPCAPLTIPSIR	S(0.0)VHIS(0.50)T(0.50)PEKEPCAPLT(0.0)IPS(0.0)IR	NP_001123917.1
6448	19CPTAC_LUAD_	4	26.17	24.25	15	0	1	11.191	6	1	1	0	DsFEMEEVQSTEGEAPHVF	DS(0.99)FEMEEVQS(0.0)T(0.0)EGEAPHVPAT(0.0)Y	NP_001123917.1
6448	19CPTAC_LUAD_	5	25.78	25.41	20	1	2	3.291	7	1	1	0	NIMTQQKDsFEMEEVQST	NIMT(0.0)QQKDS(0.99)FEMEEVQS(0.0)T(0.0)EGEAPHVPAT(0.0)Y	NP_001123917.1
6448	15CPTAC_LUAD_	4	25.72	23.78	15	0	1	10.052	6	1	1	0	DsFEMEEVQSTEGEAPHVF	DS(0.99)FEMEEVQS(0.0)T(0.0)EGEAPHVPAT(0.0)Y	NP_001123917.1
6448	12CPTAC_LUAD_	4	25.6	25.49	18	1	2	2.259	7	1	1	0	NIMTQQKDsFEMEEVQST	NIMT(0.0)QQKDS(0.99)FEMEEVQS(0.0)T(0.0)EGEAPHVPAT(0.0)Y	NP_001123917.1

# Protein / VM-site Summary Modes

**Summarize** Save As... Load...

Queue request  Excel export  
Mode: Protein - Protein Comparison  
Group columns by:  File  Directory  
Data directories: Select ...

17CPTAC\_BCProspective\_Proteome\_BI\_201706

Search result files: \*.spo  
Search result files exclude:

**Validation and Sorting**  
Filter results by: valid  
Protein grouping method: expand subgroups, top uses shared, SGT  
Sort proteins by: Score  
Filter by protein score: > 0.0  
Filter peptides by: Score: > 0.0 % SPI: > 0.0  
Required AAs: any  
Disallowed AAs: none  
Accession #'s:

**Review Fields**  
 Filename  Protein MW  Prot pl  
 Score  Species  
 Run specific  Accession #  
 % Coverage  Reporter Ratios: TMT10  
 Distinct peptides  
 Distinct peptide forms/mods  
 DEQ ratio  
 Reporter Ratios: TMT10  
Control: 130N 130C 131  
Correction Method: afRICA

**Protein Quantitation Options**  
 Exclude poor isotope quality Precursor XIC's: < 0.70 Chi<sup>2</sup> vs. Averagine  
 Exclude poor Precursor Isolation Purity: < 50 %

Spectrum Mill Extractor MS/MS Search Autovalidation Quality Metrics & FDR Workflows Tool Belt Process Report MRM Selector Spectrum Summary Help

**Summarize** Save As... Load...

Queue request  Excel export  
Mode: Protein - Var Mod Site Comparison  
Group columns by:  File  Directory  
Data directories: Select ...

20CPTAC\_LUAD\_Phosphoproteome\_BI\_201807

Search result files: \*.spo  
Search result files exclude:

**Validation and Sorting**  
Filter results by: valid  
Protein grouping method: unexpand subgroups  
Sort proteins by: Score  
Filter by protein score: > 0.0  
Filter peptides by: Score: > 0 % SPI: > 0.0  
Required AAs: sltly  
Disallowed AAs: none  
Accession #'s:

**Review Fields**  
 Filename  Sequence  b/y map  
 Score  Rev Sequence  Rank 2  
 Fwd-Rev score  VML sequence  
 Rank 1-2 score  Prec Av Chi<sup>2</sup>  Isol Pur  
 SPI (%)  Ret time, width  
 Ion mobility  
 # Backbone Cleavages  
 Unmatched ions  
 Var mod sites  
 Precursor m/z  MH+  
 Solution charge  Delta mass  Pep pl  
 VML score: sltly  # Sequences per MH+  
 Start AA position  Protein MW  Prot pl  
 Species  
 Accession #  
 Protein name

**Protein Quantitation Options**  
 Exclude poor isotope quality Precursor XIC's: < 0.70 Chi<sup>2</sup> vs. Averagine  
 Exclude poor Precursor Isolation Purity: < 50 %

**Spectrum Grouping Options**  
 Group missed cleavages containing VM site(s)

# PSM / Peptide Summary Modes

Spectrum Mill - Protein/Peptide Summary - Karl|PSMExport\_Enshg19-3nrB721glvarSmorfNuorf

Spectrum Mill | Extractor | MS/MS Search | Autovalidation | Quality Metrics & FDR | Workflows | Tool Belt | Process Report | MRM Selector | Spectrum Summary | Help

Summarize | Save As... | Load...

Queue request  Excel  AMRT export

Mode: Peptide - Spectrum Match

Filter by Proteogenomic Features: Off

Data directories: Select ...

Karl/B721GLvarSmorfNuORFri2/B5703\_2018082

Search result files: \*.spo

Search result files exclude:

**Validation and Sorting**

Filter results by: valid

Validation preset: none

Sort peptides by: Score

Filter peptides by: Score: > 0.0 % SPI: > 0.0

Required AAs: any

Disallowed AAs: none

Peptide pI: from 3.0 to 10.0  All

Accession #'s:

**Review Fields**

Filename  Score  Fwd-Rev score  Rank 1-2 score  SPI (%)  Backbone Cleavage Score  Glyco Product Ions Score  Unmatched ions  Var mod sites  VML score  Solution charge  Start AA position  Proteogenomic features

Sequence  Rev Sequence  VML sequence  Prec Av Chi<sup>2</sup>  Ret time, width  Ion mobility  Precursor m/z  Delta mass  Protein MW  Species  Accession #  Protein name

b/y map  Rank 2  Precursor XIC Intensity  Precursor S/N  DEQ ratios  Invert  Reporter Ratios: TMT 16

Control: 126C 127N 127C  Intensities

Correction Method: Static

Modification names  N-term  C-term  Cysteines  Fragmentation mode  Max tag length  Longest tag  # b/y pairs

Category: Select ..

Ensembl.human.hg19.clean3.fasta.categories.tsv:geneSymbol  
Ensembl.human.hg19.clean3.fasta.categories.tsv:chromosome  
Ensembl.human.hg19.clean3.fasta.categories.tsv:genomeCoordinates

Spectrum Mill - Protein/Peptide Summary - Karl|PeptideExport\_Enshg19-3nr

Spectrum Mill | Extractor | MS/MS Search | Autovalidation | Quality Metrics & FDR | Workflows | Tool Belt | Process Report | MRM Selector | Spectrum Summary | Help

Summarize | Save As... | Load...

Queue request  Excel  AMRT export

Mode: Peptide - Distinct

Filter by Proteogenomic Features: Off

Filter to distinct peptides: Case sensitive - CS

Data directories: Sel Off

Karl/B721GLvarSmorfNu

Search result files: \*.spo

Search result files exclude: File CS

**Validation and Sorting**

Filter results by: valid

Sort peptides by: Sequence file charge

Filter peptides by: Score: > 0.0 % SPI: > 0.0

Required AAs: any

Disallowed AAs: none

Peptide pI: from 3.0 to 10.0  All

Accession #'s:

**Review Fields**

Filename  Score  Fwd-Rev score  Rank 1-2 score  SPI (%)  Backbone Cleavage Score  Glyco Product Ions Score  Unmatched ions  Var mod sites  VML score  Solution charge  Start AA position  Proteogenomic features

Sequence  Rev Sequence  VML sequence  Prec Av Chi<sup>2</sup>  Ret time, width  Ion mobility  Precursor m/z  Delta mass  Protein MW  Species  Accession #  Protein name

b/y map  Rank 2  Precursor XIC Intensity  DEQ ratios  Invert  Reporter Ratios: TMT 16

Control: 126C 127N 127C  Intensities

Correction Method: afRIC

Modification names  N-term  C-term  Cysteines  Fragmentation mode

Category: Select ..

Ensembl.human.hg19.clean3.fasta.categories.tsv:transcriptID

**Distinct Peptide Quantitation Options**

Combine PSMs for Peptide Quantitation

## Refine notion of sameness for collapsing PSMs

**Case insensitive CI** - combine variable mods (lowercase AA's), no mods, diff precursor z, diff LC-MS/MS runs.

**Case sensitive CS** - separate mods & no mods

**Charge file CS** - separate diff precursor z, diff LC-MS/MS runs, diff mods.

**File CI** - separate diff LC-MS/MS runs

**File CS** - separate diff LC-MS/MS runs, diff mods

**Charge CS** - separate diff precursor z, diff mods

## **Autovalidation**

**Confident identification and False Discovery Rates (FDR) at the PSM, peptide, and protein levels.**

# Autovalidation – Peptide/PSM level

Spectrum Mill - MS/MS Autovalidation - Defaults\peptide\_auto\_z25\_ppm\_MSL6\_1\_0

Spectrum Mill | Extractor | MS/MS Search | Quality Metrics & FDR | Protein/Peptide Summary | Workflows | Tool Belt | Help

Validate Files  Queue request Undo Last Clear All Save As... Load...

Data Directories

Select ...

JurkatQC/Lacks/2021Q4/20210814

Fragmentation mode: All

Search result files: \*.spo

Validation Parameters

Strategy:  Fixed thresholds  Auto thresholds

Mode: Peptide Auto determine using score, delta R1-R2 thresholds to reach a target FDR

Optimize score & R1-R2 score thresholds with max FDR: 1.0 % across each:  LC run  Directory

Precursor charge range: 2 to 5 Min Sequence Length: 6 Max Sequence Length: 99 Min Backbone Cleavage Score: 0

Required AAs: any Disallowed AAs: none

Filtering

Filtering	Automatic variable range for each run	Fixed range for all runs
<input type="radio"/> None (ppm)	<input checked="" type="radio"/> Auto precursor mass error	<input type="radio"/> Fixed precursor mass error Low -1.0 High 30.0 ppm
<input checked="" type="radio"/> None (SC/pl)	<input type="radio"/> Auto SCX Solution Charge, pH3	<input type="radio"/> Fixed Solution Charge Low -2 High 6
	<input type="radio"/> Auto OGE/IEF peptide pl	<input type="radio"/> Fixed peptide pl Low 3.0 High 10.0

Results

Separate Score thresholds for each:

- LC-MS/MS run
- Precursor charge

Separate precursor mass error range for each LC-MS/MS run

- Median +/- 3 std dev

## Min Sequence length filter

- Short peptides often contribute to multiple proteins, inclusion may skew protein quantitation.
- Extractor MH+ lower limit: 750 for iTRAQ excludes short Arg peps
- Autovalidation MSL filter excludes short Lys & Arg peps

# Autovalidation – Peptide/PSM level – immunopeptidomics extras

Spectrum Mill - MS/MS Autovalidation - Defaults\HLA\_classI\_z1to4\_BCS5\_FDR1\_acrossDIR

Spectrum Mill | Extractor | MS/MS Search | Quality Metrics & FDR | Protein/Peptide Summary | Workflows | Tool Belt | Help

Validate Files  Queue request Undo Last Clear All Save As... Load...

Data Directories

Select ...

Fragmentation mode: All

Search result files: \*.spo

Susan/Covid\_lungs/SARS\_CoV2\_Lung01\_newprep\_2IPs\_ClassI\_20210714

Validation Parameters

Strategy:  Fixed thresholds  Auto thresholds

Mode: Peptide Auto determine using score, delta R1-R2 thresholds to reach a target FDR

Optimize score & R1-R2 score thresholds with max FDR: 1.0 % across each:  LC run  Directory

Precursor charge range: 1 to 4 Min Sequence Length: 7 Max Sequence Length: 99 Min Backbone Cleavage Score: 5

Required AAs: any Disallowed AAs: none

Filtering

Fixed range for all runs

None (ppm)

Fixed precursor mass error Low -1.0 High 30.0 ppm

None (SC/pl)

Fixed Solution Charge Low -2 High 6

Fixed peptide pl Low 3.0 High 10.0

Results

Separate Score thresholds for each:

- Precursor charge

## Optimize thresholds by directory

Useful for:

- Low abundance samples with <1,000 PSMs/LC run
- Low frequency matches in complex samples
  - PTM's in unenriched samples
  - Proteogenomic variants
  - High charge states (>4+)

## Backbone cleavage score

Enforces uniformly higher minimum sequence coverage for each PSM, at least 4 or 5 residues of unambiguous sequence.

**Built for immunopeptidomics.**

# Autovalidation - Protein Polishing SGT

Spectrum Mill - MS/MS Autovalidation - Defaults\proteinPolishing\_both\_minDir1\_minScore0\_protFDR0

Validate Files  Queue request Undo Last Clear All Save As... Load...

Data Directories

Select ...

CPTAC3pancancer/LSCC/Proteome/01CPTAC\_LSCC\_Proteome\_BI\_20190615

Validation Parameters

Strategy:  Fixed thresholds  Auto thresholds

Mode: Protein polishing Auto determine using score, delta R1-R2 thresholds to reach a target FDR

Group proteins across all directories

Protein grouping method: unexpand subgroups

Method for applying combined thresholds of Protein Score and Minimum number of directories:

Retain proteins above both thresholds (more strict)

Retain proteins above either threshold (less strict - retains recurrently observed proteins with scores below threshold)

Minimum number of directories protein group is observed in: 1

Minimum protein score: 0

Automatically raise minimum protein score to yield maximum protein FDR: 0.0 %

**Primary Purpose**  
Eliminate  
low scoring  
single-experiment  
single-peptide  
protein groups

Assembles protein groups from the autovalidated PSM's, determines the maximum protein level score of a protein group that consists entirely of distinct peptides estimated to be false-positive identifications (PSM's with negative delta forward-reverse scores). Then PSM's are unvalidated if they contribute to protein groups that **BOTH**:

- have protein scores at or below the larger of
  - the minimum protein score
  - the max false-positive protein score
- are observed in less than the minimum # of directories

Spectrum Mill - MS/MS Autovalidation - Defaults\proteinPolishing\_either\_2\_25\_SGT

Validate Files  Queue request Undo Last Clear All Save As... Load...

Data Directories

Select ...

CPTAC3pancancer/LSCC/Proteome/01CPTAC\_LSCC\_Proteome\_BI\_20190615

CPTAC3pancancer/LSCC/Proteome/02CPTAC\_LSCC\_Proteome\_BI\_20190618

CPTAC3pancancer/LSCC/Proteome/03CPTAC\_LSCC\_Proteome\_BI\_20190626

CPTAC3pancancer/LSCC/Proteome/04CPTAC\_LSCC\_Proteome\_BI\_20190629

CPTAC3pancancer/LSCC/Proteome/05CPTAC\_LSCC\_Proteome\_BI\_20190705

Validation Parameters

Strategy:  Fixed thresholds  Auto thresholds

Mode: Protein polishing Auto determine using score, delta R1-R2 thresholds to reach a target FDR

Group proteins across all directories

Protein grouping method: expand subgroups, top uses shared

Method for applying combined thresholds of Protein Score and Minimum number of directories:

Retain proteins above both thresholds (more strict)

Retain proteins above either threshold (less strict - retains recurrently observed proteins with scores below threshold)

Minimum number of directories protein group is observed in: 3

Minimum protein score: 25

**Primary Purpose**  
Eliminate  
low scoring  
single-experiment  
protein Subgroups

From all directories, assembles protein subgroups from the autovalidated PSM's. Then PSM's are unvalidated if they contribute to protein subgroups that **EITHER**:

- have protein subgroup-specific scores at or below the minimum protein score
- OR
- are observed in less than the minimum # of directories

# Autovalidation – VM site Polishing

Spectrum Mill - MS/MS Autovalidation - CPTAC3human\CPTAC\_VMsitePolish\_sty\_multiDir\_retainEither\_3dir\_8score

Spectrum Mill | Extractor | MS/MS Search | Quality Metrics & FDR | Protein/Peptide Summary | Workflows | Tool Belt | Help

Validate Files  Queue request Undo Last Clear All Save As... Load...

Data Directories

Select ...

Fragmentation mode: All

Search result files: \*.spo

Validation Parameters

Strategy:  Fixed thresholds  Auto thresholds

Mode: VM site polishing Auto determine using score, delta R1-R2 thresholds to reach a target FDR

Group proteins across all directories (unexpanded subgroups)

Group VM sites across all directories, and group missed cleavages containing the same VM site

Variable modification site localization s|t|y

Method for applying combined thresholds of VM site Score and Minimum number of directories:

Retain VM sites above both thresholds (more strict)

Retain VM sites above either threshold (less strict - retains recurrently observed VM sites with scores below threshold)

Minimum number of directories VM site is observed in: 3

Minimum VM site score: 8.0

Results

Primary Purpose  
Eliminate  
low-scoring,  
single-experiment  
phosphosites

From all directories, assembles VM sites from the autovalidated PSM's. Then PSM's are unvalidated if they contribute to VM sites that EITHER:

- have VM site scores at or below the minimum VM site score
- OR
- are observed in less than the minimum # of directories

## Quantitation – reporter ion ratios (TMT, iTRAQ)

# Key Features of Spectrum Mill Quantitation

- MS/MS reporter ion ratio based: iTRAQ, TMT
  - Median of all PSM ratios for each Protein, VM-site
    - Not sum of all PSM reporter ion intensities, then ratio for Protein, VM-site
  - Exclude ratios with Precursor Ion Purity, PIP < 50%
  - Exclude ½ of false positive ID's: Delta Fwd-Rev < 0
  - Exclude peptides with no label
- MS precursor ion ratio based: SILAC
  - Only 1 member of H/L pair H/M/L triplet needs to be triggered for MS/MS
  - Median of all H/L pair H/M/L triplet ratios for each Protein, VM-site
  - Exclude ratios from poor precursor ion isotope cluster shape
- Label – free
  - Only uses peptides identified by MS/MS
  - Sums up all precursor ion peak area for all PSM's for a protein
- Peptides shared between protein subgroups
  - Use shared in each subgroup
  - Use only unshared: Subgroup specific



# Providing Reporter Ion Correction Factors

Spectrum Mill - Tool Belt Utilities

Spectrum Mill | Extractor | MS/MS Search | Autovalidation | Quality Metrics & FDR | Protein/Peptide S

Select a Tool:

Stop process       List modifications details       File collector  
 Create saved results file       Create Reporter Ion correction factors       Export PepXML  
 Create MS/MS search summary file       Apply Reporter Ion correction factors       Convert spectra  
 View parameter table

**Create Reporter Ion correction factors**

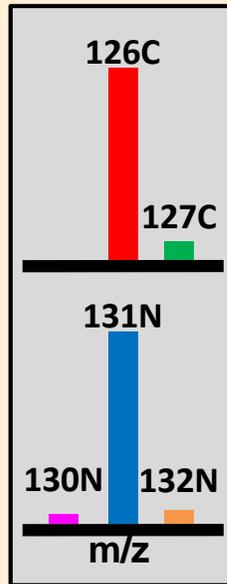
(iTRAQ 8-plex corrections do not vary by batch lot)

Create TMT 16-plex

Batch: VA299611

Mass	% of -2	% of -1	% of +1	% of +2
126C	0.0	0.0	9.31	0.32
127N	0.0	0.0	8.98	0.32
127C	0.0	0.93	8.63	0.27
128N	0.0	0.95	8.38	0.26
128C	0.0	1.47	6.91	0.15
129N	0.0	1.46	6.86	0.15
129C	0.51	2.74	6.15	0.11
130N	0.49	2.76	5.98	0.11
130C	0.04	3.10	4.82	0.06
131N	0.04	3.09	4.75	0.06
131C	0.08	3.81	3.29	0.03
132N	0.04	2.84	3.51	0.02
132C	0.10	4.14	1.80	0.0
133N	0.36	3.64	1.94	0.0
133C	0.88	4.70	1.01	0.0
134N	0.4	4.92	1.05	0.0

Enter the <sup>13</sup>C values from the Certificate of Analysis



Updates file:  
 \\spectrumMill\msparams\_mill\reporterlon.corrections.txt

Spectrum Mill - Tool Belt Utilities

Spectrum Mill | Extractor | MS/MS Search | Autovalidation | Quality Metrics & FDR | Protein/Peptid

Select a Tool:

Stop process       List modifications details       File collector  
 Create saved results file       Create Reporter Ion correction factors       Export PepXML  
 Create MS/MS search summary file       Apply Reporter Ion correction factors       Convert spectra  
 View parameter table

**Apply Reporter Ion correction factors**

Apply

Batch: TMT16\_VA299611

**Data Directory**

Select ... Hasmik/HIV\_RAPMS\_TMT16/HIV\_RAPMS\_TMT16\_fxns

Batch: TMT16\_VA299611

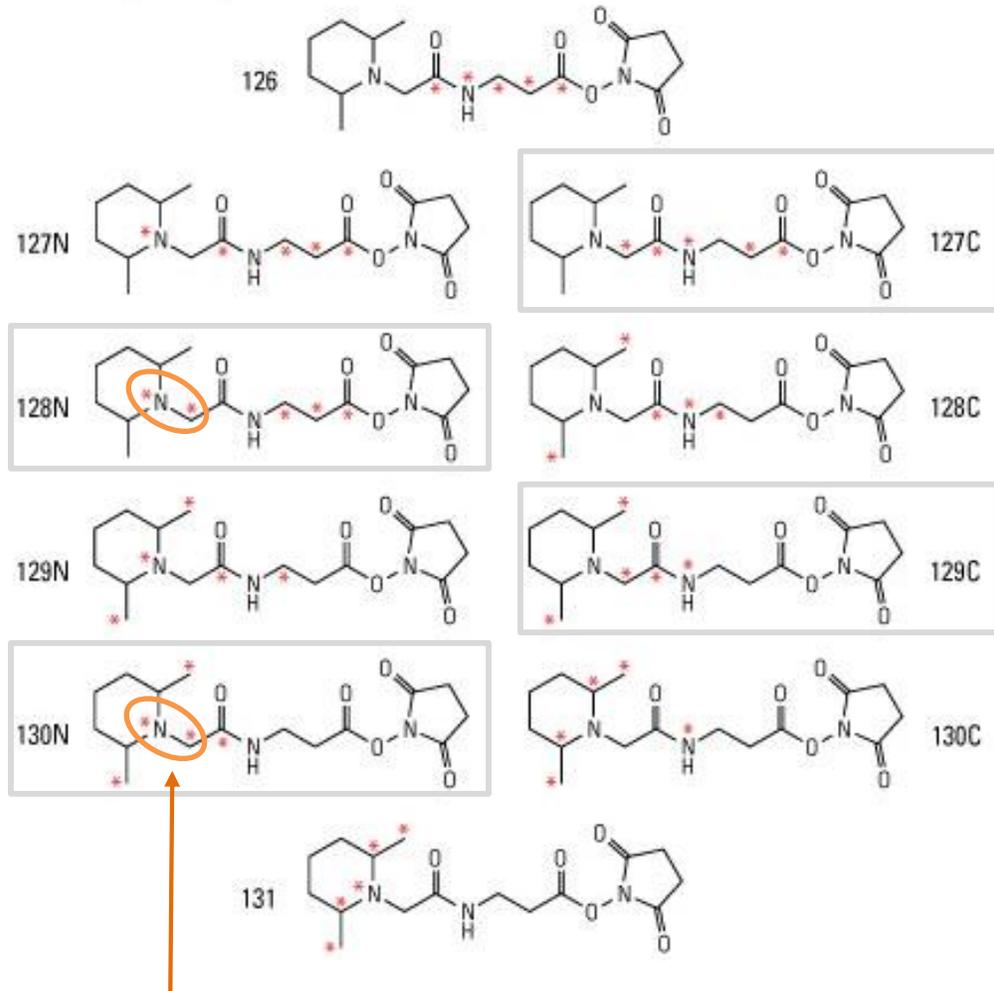
126C	0.00	0.00	9.31	0.32
127N	0.00	0.00	8.98	0.32
127C	0.00	0.93	8.63	0.27
128N	0.00	0.95	8.38	0.26
128C	0.00	1.47	6.91	0.15
129N	0.00	1.46	6.86	0.15
129C	0.51	2.74	6.15	0.11
130N	0.49	2.76	5.98	0.11
130C	0.04	3.10	4.82	0.06
131N	0.04	3.09	4.75	0.06
131C	0.08	3.81	3.29	0.03
132N	0.04	2.84	3.51	0.02
132C	0.10	4.14	1.80	0.00
133N	0.36	3.64	1.94	0.00
133C	0.88	4.70	1.01	0.00
134N	0.40	4.92	1.05	0.00

Creates file:  
 \\spectrumMill\msdataSM\Hasmik\HIV\_RAPMS\_TMT16\HIV\_RAPMS\_TMT16\_Allreporterlon.correction.txt

Used in SM modules:  
 Quality Metrics  
 Protein/Peptide Summary

# Source of Isotope Impurities

## B. TMT10plex Reagents (TMT<sup>10</sup>)



Probably most difficult combination to synthesize with high incorporation  
 >1% signal for 127C (in 128N) and for 129C (in 130N).

- +1 impurities originate from natural <sup>13</sup>C levels in the unlabeled C positions
  - C-reagents have only C-to-C +1.003 impurities
  - N-reagents have only N-to-N +1.003 impurities
- 1 impurities originate from incomplete incorporation of <sup>15</sup>N or <sup>13</sup>C in the synthetic precursors
  - C-reagents may have only C-to-C -1.003 impurities
  - N-reagents may have both:
    - N-to-C -0.997 impurities (<sup>15</sup>N-<sup>14</sup>N)
    - N-to-N -1.003 impurities (<sup>13</sup>C-<sup>12</sup>C)

### TMT6/TMT10 Reporters

126	+0		
127N	+1		<sup>15</sup> N
127C	+1	<sup>13</sup> C	
128N	+2	<sup>13</sup> C	<sup>15</sup> N
128C	+2	<sup>13</sup> C <sub>2</sub>	
129N	+3	<sup>13</sup> C <sub>2</sub>	<sup>15</sup> N
129C	+3	<sup>13</sup> C <sub>3</sub>	
130N	+4	<sup>13</sup> C <sub>3</sub>	<sup>15</sup> N
130C	+4	<sup>13</sup> C <sub>4</sub>	
131	+5	<sup>13</sup> C <sub>4</sub>	<sup>15</sup> N

# afRICA Reporter Ion Intensity Correction Method in SM

Shadforth IP, Dunkley TPJ, Lilley KS, and Bessant C  
 BMC Genomics 2005, 6:145 doi:10.1186/1471-2164-6-145

**BMC Genomics**



Software

Open Access

## i-Tracker: For quantitative proteomics using iTRAQ™

Ian P Shadforth\*<sup>1</sup>, Tom PJ Dunkley<sup>2</sup>, Kathryn S Lilley<sup>2</sup> and Conrad Bessant<sup>1</sup>

Address: <sup>1</sup>Department of Analytical Science and Informatics, Cranfield University at Silsoe, Silsoe, Bedfordshire, UK and <sup>2</sup>Cambridge Centre for Proteomics, Biochemistry Department, Cambridge University, Cambridgeshire, UK

Email: Ian P Shadforth\* - i.p.shadforth.s01@cranfield.ac.uk; Tom PJ Dunkley - tpjd2@cam.ac.uk; Kathryn S Lilley - k.s.lilley@bioc.cam.ac.uk; Conrad Bessant - c.bessant@cranfield.ac.uk

\* Corresponding author

### Purity correction

Each batch of iTRAQ reagents supplied by ABI is labelled with sixteen purity values indicating the percentages of each reporter ion that have masses differing by -2, -1, +1 and +2 Da from the nominal reporter ion mass due to isotopic variants. This information can be used to correct the peak areas calculated for each reporter ion to account for the losses to, and gains from, other reporter ions. Losses to ion peaks not in the reporter ion range are also accounted for in this method.

w,x,y,z represent the percentage of each peak expected to be present at the mass of the reporter ion associated with that peak. Here, w is for 114.1, x for 115.1 etc.:

$$w = (100 - (a + e + i + m))$$

$$x = (100 - (b + f + j + n))$$

$$y = (100 - (c + g + k + o))$$

$$z = (100 - (d + h + l + p))$$

The area ( $A_r$ ) of each reporter ion peak ( $r$ ), as calculated above, can now be written in terms of the true areas of peaks ( $T_r$ ):

$$\begin{aligned} A_{114.1} &= (w * T_{114.1}) & +(f * T_{115.1}) & +(c * T_{116.1}) \\ A_{115.1} &= (i * T_{114.1}) & +(x * T_{115.1}) & +(g * T_{116.1}) + (d * T_{117.1}) \\ A_{116.1} &= (m * T_{114.1}) & +(j * T_{115.1}) & +(y * T_{116.1}) + (h * T_{117.1}) \\ A_{117.1} &= & (n * T_{115.1}) & +(k * T_{116.1}) + (z * T_{117.1}) \end{aligned}$$

The task is now to calculate each  $T_r$  according to these equations.

The simultaneous equations needed to solve this problem are fairly complicated, but can be framed such that Cramer's rule may be applied. A detailed explanation of

a,b,c,d,e,f,g,h,i,j,k,l,m,n,o,p are the sixteen purity correction values (as percentages) in the order:

114.1 - 2 Da, 115.1 - 2 Da, 116.1 - 2 Da, 117.1 - 2 Da, 114.1 - 1 Da, etc...

### Algorithm for Reporter Ion Correction and Adjustment (afRICA)

afRICA extends the i-TRACKER algorithm from iTRAQ to TMT

- 4x4 = 16 correction factors for iTRAQ4
- 6x4 = 24 correction factors for TMT6
- 10x4 = 40 correction factors for TMT10

### Matrix of Coefficients

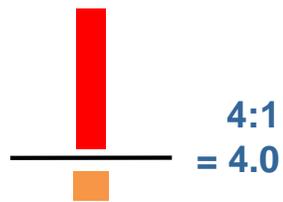
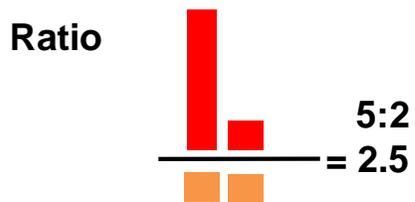
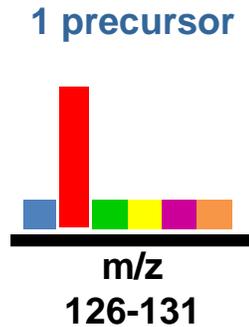
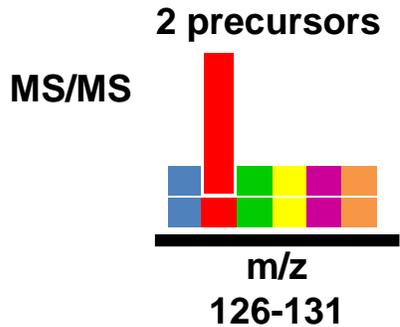
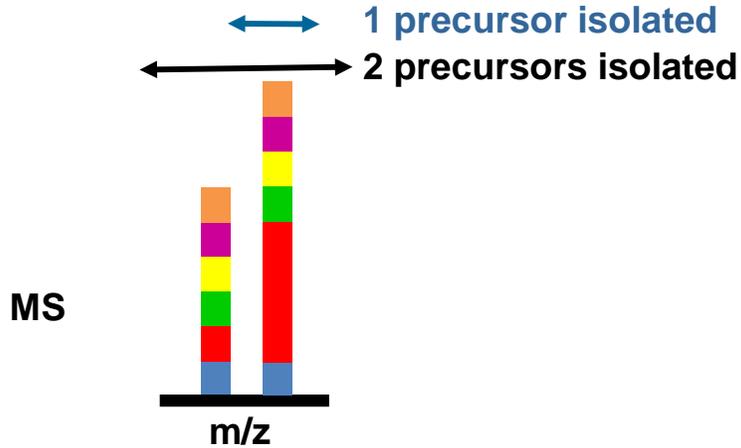
- iTRAQ4 4X4
- TMT6 6x6
- TMT10 5x5 (N), 5x5 (C)
- TMT11 5x5 (N), 6x6 (C)
- TMT16 8x8 (N), 8x8 (C)
- TMT18 9x9 (N), 9x9 (C)

4 equations and 4 unknown  
 -> Solve for unknowns using **Cramer's rule**

The determinant of the matrix of coefficients can be found:

$$|C| = \begin{vmatrix} w, f, c, 0 \\ i, x, g, d \\ m, j, y, h \\ 0, n, k, z \end{vmatrix}$$

# Interference and Ratio Compression



## Minimizing Interference

### Data generation

- Use narrow MS1 precursor isolation
  - 0.7 m/z instead of 2.5

### Data analysis

- Exclude MS/MS with low precursor ion purity (PIP) i.e. <50%
  - $PIP = \frac{100\% \times \text{intensity precursor 1}}{\text{intensity all isolation window}}$

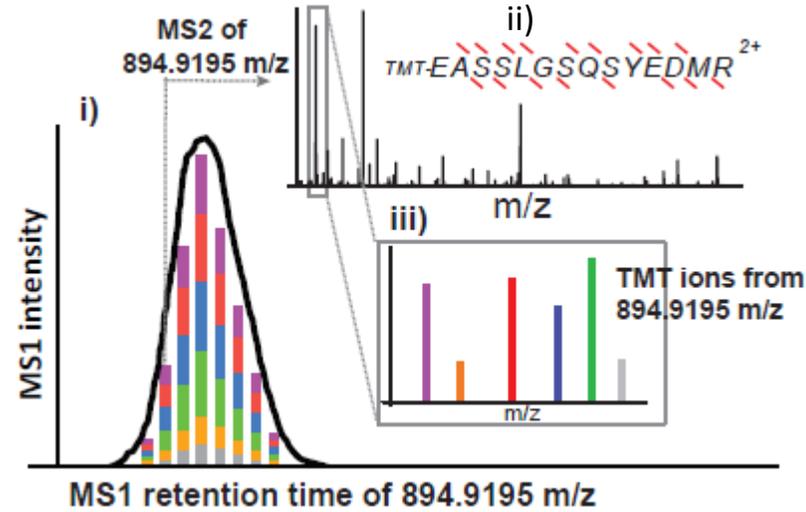
For 2 peptides, relative precursor intensities not proportional to reporter ion intensities

Factors:

- Precursor charge
- MS/MS Collision energy
- # TMT labels (Lys vs Arg)

# De-multiplexing iTRAQ and TMT Experiments

Method is used when  
instead of  
quantifying across samples ↔  
one seeks to  
quantify along proteins ⇕



$$\text{Peptide}_{126} \text{ abundance} = \frac{\text{MS2 reporter}_{126} \text{ abundance}}{\sum_{k=126}^{131} \text{MS2 reporter}_k \text{ abundance}} \times \text{MS1 precursor abundance}$$

$$\text{Protein}_{126} \text{ abundance} = \sum_{j=0}^n \text{Peptide}_{126} \text{ abundance}$$

n is # peptides in the protein

Repeat for each  
reporter ion

- 126
- 127
- 128
- 129
- 130
- 131

Not yet done directly in SM – done manually using values exported in SM reports, aka fractional intensity method



# Process Report

# Process Report

Spectrum Mill - Process Report - CPTAC3human\VMsite\_TMT10\_medMAD\_Gencode-smORFs-nuORFs

Spectrum Mill | Extractor | MS/MS Search | Autovalidation | Quality Metrics & FDR | Protein/Peptide Summary | Workflows | Tool Belt | Protigy | Help

## Process Report

Process

Save As...

Load...

## Data Directories

Select ...

CPTAC3pancancer/MEDUL/Phosphoproteome\_GencodeGS/TMTplex01

Protein - Protein Comparison, all use shared  
Protein - Protein Comparison SGT, top uses shared  
Protein - Protein Comparison SGS, ignore shared  
**Protein - Var Mod Site Comparison**  
Protein - Peptide Comparison  
Protein - Prot Genom Site Comparison Variants  
Protein - Prot Genom Site Comparison Spliceforms

## Processing Parameters

Project: Gencode-smORFs-nuORFs

SM report type:

Parse Report

Reporter ion label type:

iTRAQ 4  
iTRAQ 8  
TMT 0  
TMT0\_16  
TMT 2  
TMT 6  
**TMT 10**  
TMT 11  
TMT 16  
iTRAQ 4 mean/median multi  
TMT 6 mean/median multi  
TMT 10 mean/median multi  
TMT 11 mean/median multi  
TMT 16 mean/median multi  
Label free  
SILAC 2

Plot Ratio Distributions

Max plots per row:

3

Normalize Reporter Ratios

Method:

Median / MAD

Plot Ratio Distributions 2

Max plots per row:

3

## Results

- Keep most useful columns from \*.ssv
- Revise headers
- Map samples to reporter ions
- Filter species, exclude contaminants
- Normalize TMT ratios
- Reformats .TXT, .GCT

Launches Perl script processReport.pl which runs consecutive R scripts

- parseSMreport.r
- plotRatioDistributions.r
- normalizeReporterRatios.r
- plotRatioDistributions.r

Project-specific features in R scripts

Generates GCT format, if \*-sample-annotation.csv present



# Create vertical layout file: \*-sample-annotation.csv

## Create a sample-annotation.csv file (vertical layout)

The requirements described below are, by design, harmonized between tools developed in the Broad Institute Proteomics Platform group including: Spectrum Mill and the downstream tools **Protigy** and **Panoply**.

In order for the parseSMreport script to correlate each reporter ion label with a sample name and metadata to produce a **GCT v1.3 format** output file you must provide a **comma-delimited .csv file** in the SM directory that contains the report that is parsed. While the filename may be prefixed at the user's discretion, the text file must be named with the suffix

sample-annotation.csv

And be organized like the example below:

The following considerations apply:

- Required columns are only: **Sample.ID**, **Experiment**, and **Channel**.
- Though not required for SM, the column: **Type** (Tumor, NAT, etc) is required downstream for Panoply use.
- Note that the sample-annotation file **must not contain a row for the sample(s) used as a control ion (denominator in ratios)**. However, users are encouraged to include the isobaric label name and control ion channel in the filename of the sample-annotation file.
- The headers of metadata columns need to all be R-compatible. No spaces, no control characters like (), and can not begin with a number.
- Replicates: Sample.IDs must be unique and cannot have duplicates. If replicate samples are present, they must have unique Sample.IDs, but can be identified as replicates by using an additional column (named "Participant" for downstream Panoply use) with identical ids.
- The column **Experiment** (integer values) **must be numbered in the exact same order** as the data directories appear in the report (or Sample.ID and metadata will be mis-applied), which is the same order as they were selected when creating the report in P/P Summary. Note that the SM data directory need not, but could be a column in the sample-annotation file.
- The column **Channel** should include the reporter ion label used in SM reports. Please note the following nuances: For TMT6, 10 SM reports: 126, 131, while for TMT11, 16, 18 SM reports: 126C, 131N  
. The use of and N and C on every channel was implemented for clarity, and I regret not implementing the convention for the earliest reagents, in deference to the manufacturer's nomenclature, though considered it at the time. Now to ensure backwards compatibility the difference is what it is.

### Sorting of the sample-annotation.csv file:

Sort as you wish. It has been expected that when some users prepare a sample-annotation file they might frequently end up with rows sorted in a different order than the data columns in the SM report to which they will be mated, for reasons like sorted by Sample.IDs, or the Channel column sorted alphabetically instead of by mass. So to reduce excess "quick questions" back to the developer, the script has been written to tolerate differences in sort. However, keep in mind that **the output .GCT file will maintain the order of the columns in the input .ssv file**. It is expected that the metadata columns in the sample-annotation file which become rows in the .GCT file will enable facile re-sorting of the data columns during downstream analyses.

### Example sample-annotation.csv file

Download a template suited to your labeling reagent: **TMT6 TMT10 TMT11**

- Manually create file, place in SM data directory
- Automatically read by script: parseSMreport.r



Sample.ID	Participant	Experiment	Channel	Type	Aliquot	Smoking.History	Stage	Country.of.Origin	Age	Gender	Ethnicity	Cigarettes.per.Day	Pack.Years.Smoked	Secondhand.Smoke	QC.status
C3L.02665	C3L-02665	1	126C	Tumor	CPT0192540004	Smoking history not available	IB	bulgaria	51	male	caucasian			Exposure t	QC.pass
C3L.02665.N	C3L-02665	1	127N	NAT	CPT0192590004	Smoking history not available		bulgaria	51	male	caucasian			Exposure t	QC.pass
C3L.01663	C3L-01663	1	127C	Tumor	CPT0091610003	Smoking history not available	IIIA	usa	73	female	caucasian			Exposure t	QC.pass
C3L.01663.N	C3L-01663	1	128N	NAT	CPT0091620003	Smoking history not available		usa	73	female	caucasian			Exposure t	QC.pass
C3N.02575	C3N-02575	1	128C	Tumor	CPT0147680004	Current reformed smoker	IA	china	69	male	han	40	100	Exposure t	QC.pass
C3N.02575.N	C3N-02575	1	129N	NAT	CPT0147700003	Current reformed smoker within		china	69	male	han	40	100	Exposure t	QC.pass
C3L.02546	C3L-02546	1	129C	Tumor	CPT0112650003	Current reformed smoker, IB		usa	74	male	caucasian	20	40	Exposure t	QC.pass
C3L.02546.N	C3L-02546	1	130N	NAT	CPT0112660003	Current reformed smoker, more		usa	74	male	caucasian	20	40	Exposure t	QC.pass
C3L.00965	C3L-00965	1	130C	Tumor	CPT0101130003	Current reformed smoker	IIA	russia	63	male	white(cauc	20	42	Yes	QC.pass
LUAD.Global.CR..pool.1	LUAD Globa	1	131N		LUAD Global CR (pool 1)										QC.fail
C3L.02963	C3L-02963	2	126C	Tumor	CPT0197090003	Current smoker: Includes (I		usa	71	male	caucasian	20	56	Exposure t	QC.pass
C3L.02963.N	C3L-02963	2	127N	NAT	CPT0197100003	Current smoker: Includes daily a		usa	71	male	caucasian	20	56	Exposure t	QC.pass
C3N.04162	C3N-04162	2	127C	Tumor	CPT0222730003	Lifelong non-smoker: Less IIA		ukraine	46	male	slavic			Exposure t	QC.pass
C3N.04162.N	C3N-04162	2	128N	NAT	CPT0222740003	Lifelong non-smoker: Less than 1		ukraine	46	male	slavic			Exposure t	QC.pass
C3L.02646	C3L-02646	2	128C	Tumor	CPT0202530004	Current smoker: Includes (IIIA		bulgaria	74	male	caucasian	50	5	Exposure t	QC.pass
C3L.02646.N	C3L-02646	2	129N	NAT	CPT0202580004	Current smoker: Includes daily a		bulgaria	74	male	caucasian	50	5	Exposure t	QC.pass
C3N.02285	C3N-02285	2	129C	Tumor	CPT0126740003	Current reformed smoker	IIB	poland	72	female	caucasian	15	35.3	Exposure t	QC.pass
C3N.02285.N	C3N-02285	2	130N	NAT	CPT0126760003	Current reformed smoker within		poland	72	female	caucasian	15	35.3	Exposure t	QC.pass
C3N.03875	C3N-03875	2	130C	Tumor	CPT0237410004	Current reformed smoker	IB	china	68	male	han	20		Exposure t	QC.pass
C3N.03875.N	C3N-03875	2	131N	NAT	CPT0237430003	Current reformed smoker within		china	68	male	han	20		Exposure t	QC.pass
C3N.03424	C3N-03424	3	126C	Tumor	CPT0245780003	Current reformed smoker, IA		poland	74	male	caucasian	Unknown		Exposure t	QC.pass
C3N.03424.N	C3N-03424	3	127N	NAT	CPT0245790003	Current reformed smoker, more		poland	74	male	caucasian	Unknown		Exposure t	QC.pass
C3N.01017	C3N-01017	3	127C	Tumor	CPT0069890003	Current smoker: Includes (IIB		vietnam	56	male	asian	5	5	No or mini	QC.pass
C3N.01017.N	C3N-01017	3	128N	NAT	CPT0069900003	Current smoker: Includes daily a		vietnam	56	male	asian	5	5	No or mini	QC.pass
C3L.00081	C3L-00081	3	128C	Tumor	CPT0001020003	Lifelong non-smoker: Less IIA		usa	61	female	caucasian			Exposure t	QC.pass
C3L.00081.N	C3L-00081	3	129N	NAT	CPT0001030003	Lifelong non-smoker: Less than 1		usa	61	female	caucasian			Exposure t	QC.pass
C3L.02649	C3L-02649	3	129C	Tumor	CPT0204800003	Current reformed smoker	IIIA	bulgaria	62	male	caucasian	40		Yes	QC.pass
C3L.02649.N	C3L-02649	3	130N	NAT	CPT0204820003	Current reformed smoker within		bulgaria	62	male	caucasian	40		Yes	QC.pass
C3L.01838	C3L-01838	3	130C	Tumor	CPT0091820003	Current reformed smoker, IB		usa	70	male	caucasian	20	33	Exposure t	QC.pass
LSCC.Tumor.ONLY.CR	LSCC Tumor	3	131N		LSCC Tumor ONLY CR										QC.fail
C3N.02252	C3N-02252	4	126C	Tumor	CPT0128890003			other	68	female					QC.pass
C3N.02252.N	C3N-02252	4	127N	NAT	CPT0128900003			other	68	female					QC.pass
C3L.03963	C3L-03963	4	127C	Tumor	CPT0227400003	Current smoker: Includes (IIB		bulgaria	73	male	caucasian	20	60	Yes	QC.pass
C3L.03963.N	C3L-03963	4	128N	NAT	CPT0227420003	Current smoker: Includes daily a		bulgaria	73	male	caucasian	20	60	Yes	QC.pass
C3N.03072	C3N-03072	4	128C	Tumor	CPT0209020004	Current reformed smoker	IIIA	china	62	male	han	20		Exposure t	QC.pass
C3N.03072.N	C3N-03072	4	129N	NAT	CPT0209040003	Current reformed smoker within		china	62	male	han	20		Exposure t	QC.pass
C3L.02969	C3L-02969	4	129C	Tumor	CPT0217420003	Current reformed smoker	I	usa	76	male	caucasian	20	45	Exposure t	QC.pass
C3L.02969.N	C3L-02969	4	130N	NAT	CPT0217490003	Current reformed smoker within		usa	76	male	caucasian	20	45	Exposure t	QC.pass
C3N.03662	C3N-03662	4	130C	Tumor	CPT0219360003	Lifelong non-smoker: Less IIA		ukraine	64	male	slavic			Exposure t	QC.pass
C3N.03662.N	C3N-03662	4	131N	NAT	CPT0219390003	Lifelong non-smoker: Less than 1		ukraine	64	male	slavic			Exposure t	QC.pass
C3N.03051	C3N-03051	5	126C	Tumor	CPT0208770003	Current smoker: Includes (IIIA		china	65	male	han	20		Exposure t	QC.pass



# Parsed Report - Normalized, Formatted - TXT

- Revise headers based on reporter\_sample\_template.txt
- Keep only useful columns
- Normalize TMT ratios

	A	B	C	D	E	F	G	H	I	id	accession_number	GeneSymbol	Division	Category	proteinObserved	proteinExpected	numPeptides	accession_number	species	protein_mw	subgroup	score	Unique	entry_name
1	bRPfxns_0813, 126:127C, Control_205:LM2.3(9.1)	bRPfxns_0813, 127N:127C, Control_206:LM2.3(9.1)	bRPfxns_0813, 128N:127C, empty:LM2.3(9.1)	bRPfxns_0813, 128C:127C, Control_207:LM2.3(9.1)	bRPfxns_0813, 129N:127C, SNED1_KD_208:LM2.3(9.1)	bRPfxns_0813, 129C:127C, 231.1:LM2.3(9.1)	bRPfxns_0813, 130N:127C, 231.2:LM2.3(9.1)	bRPfxns_0813, 130C:127C, SNED1_KD_210:LM2.3(9.1)	bRPfxns_0813, 131:127C, SNED1_KD_211:LM2.3(9.1)	PLEC	Q15149				1	1147	416	Q15149 Q15149-4	HUMAN	533,778	1.1	5897.5	Plectin	
2	-1.58	-1.44	-1.25	-1.17	-0.77	0.05	-0.51	-0.96	-1.51	PLEC	Q15149-4			1	1135	410	Q15149-4	HUMAN	518,129	1.2	5837.2	Isoform 4		
3	-1.60	-1.47	-1.30	-1.18	-0.79	0.04	-0.53	-0.97	-1.52	PLEC	Q15149-9			1	1132	409	Q15149-9	HUMAN	516,705	1.3	5830.8	Isoform 9		
4	-1.60	-1.46	-1.28	-1.18	-0.79	0.04	-0.53	-0.96	-1.52	PLEC	Q15149-5			1	1132	409	Q15149-5	HUMAN	518,207	1.4	5828.1	Isoform 5		
5	-1.61	-1.45	-1.31	-1.18	-0.80	0.01	-0.54	-0.97	-1.52	PLEC	Q15149-8			1	1133	409	Q15149-8	HUMAN	515,637	1.5	5818.7	Isoform 8		
6	-1.61	-1.47	-1.31	-1.18	-0.81	0.00	-0.54	-0.97	-1.53	PLEC	Q15149-3			1	1128	407	Q15149-3	HUMAN	519,963	1.6	5811.7	Isoform 3		
7	-1.27	-1.15	-1.03	-1.04	-0.54	0.66	-0.29	-0.58	-1.26	Plec	Q9QXS1			1	647	276	Q9QXS1 Q9QXS1-1	MOUSE	536,118	1.7	3648.7	Plectin		
8	-2.18	-0.93	-0.73	-1.10	0.47	-1.15	-1.10	-2.67	-2.25	AHNAK	Q09666			1	726	289	Q09666 Q09666-1	HUMAN	629,604	2.1	4100.7	Neuroblas		
9	-1.48	-0.21	-0.35	-0.66	0.35	-0.21	-0.72	-1.69	-1.54	Ahnak	E9Q616			1	268	125	E9Q616 G9Q616	MOUSE	605,041	2.2	1520.7	Protein Ah		
10	0.62	0.48	0.41	0.53	1.07	0.22	0.02	0.09	0.48	DYNCLH	Q14204			1	867	270	Q14204	HUMAN	535,137	3.1	3800.6	Cytoplasm		
11	0.60	0.48	0.37	0.50	1.05	0.16	0.01	0.07	0.53	Dynd1h	Q9JHU4			1	802	260	Q9JHU4	MOUSE	534,774	3.2	3614.6	Cytoplasm		
12	5.47	3.10	2.74	0.96	3.99	3.87	4.69	2.04	4.99	Col1a1	P11087	Col1a1	Core matr	Collagens	1	1856	159	P11087 P11087-1	MOUSE	139,057	4.1	2828.7	Collagen a	
13	5.96	3.66	3.06	1.80	3.79	3.22	3.90	2.41	5.18	COL1A1	P02452	COL1A1	Core matr	Collagens	1	1222	119	P02452	HUMAN	139,966	4.2	2046.2	Collagen a	
14	4.58	2.23	1.56	0.52	3.29	2.85	4.26	1.28	4.03	COL2A1	P02458	COL2A1	Core matr	Collagens	1	158	23	P02458 P02458-1	HUMAN	142,867	4.3	236.4	Collagen a	
15	4.44	2.21	1.46	0.40	3.11	2.72	4.17	1.15	3.73	Col2a1	P28481-5	Col2a1	Core matr	Collagens	1	170	19	P28481-5 P28481-5-1	MOUSE	138,950	4.4	218.4	Isoform 5	
16	2.46	0.07	0.35	-0.07	-0.37	2.76	0.59	0.40	1.20	Col6a3	E9PWQ3	Col6a3	Core matr	Collagens	1	588	173	E9PWQ3	MOUSE	355,523	5.1	2757.7	Protein Co	
17	2.65	0.13	0.32	-0.04	-0.29	2.76	0.71	0.40	1.33	Col6a3	J3QQ16	Col6a3	Core matr	Collagens	1	539	157	J3QQ16	MOUSE	290,174	5.2	2450.5	Protein Co	
18	2.44	1.06	0.52	0.44	-0.52	1.24	0.07	-0.17	0.64	COL6A3	P12111	COL6A3	Core matr	Collagens	1	303	125	P12111 P12111-1	HUMAN	345,375	5.3	1808.9	Collagen a	
19	2.56	1.14	0.76	0.54	-0.34	1.47	0.33	-0.04	0.90	COL6A3	P12111-4	COL6A3	Core matr	Collagens	1	270	112	P12111-4	HUMAN	279,796	5.4	1603.3	Isoform 4	
20	-0.30	-0.28	-0.32	-0.46	-0.50	-0.40	-0.48	-0.11	-0.36	PRKDC	P78527			1	511	200	P78527 P78527-1	HUMAN	474,044	6.1	2604.0	DNA-dene		

# Auxiliary Files

reporter\_sample\_template.txt

- Manually create file
  - ✓ Sort columns into order desired in output
  - ✓ Place in SM data directory
- Automatically read by script: parseSMreport.r

	A	B	C	D	E	F	G	H	I	J
1	129C	130N	126	127N	128C	128N	127C	129N	130C	131
2	231.1	231.2	Control_205	Control_206	Control_207	empty	LM2.3(9.1)	SNED1_KD_208	SNED1_KD_210	SNED1_KD_211
3	bRPfxns_0813									

\*norm-stats.txt

This file is automatically generated by the script:  
normalizeReporterRatios.r

	A	B	C	D	E	F
1	norm.params.sample	norm.params.mean	norm.params.median	norm.params.sd	norm.params.mad	
2	bRPfxns_0813, 126:127C, Control_205:LM2.3(9.1)	-0.187	-0.233	0.616	0.385	
3	bRPfxns_0813, 127N:127C, Control_206:LM2.3(9.1)	0.135	0.090	0.637	0.420	
4	bRPfxns_0813, 128N:127C, empty:LM2.3(9.1)	-4.958	-4.947	0.711	0.471	
5	bRPfxns_0813, 128C:127C, Control_207:LM2.3(9.1)	-0.040	-0.077	0.554	0.387	
6	bRPfxns_0813, 129N:127C, SNED1_KD_208:LM2.3(9.1)	-0.488	-0.528	0.729	0.513	
7	bRPfxns_0813, 129C:127C, 231.1:LM2.3(9.1)	-0.043	-0.140	0.562	0.350	
8	bRPfxns_0813, 130N:127C, 231.2:LM2.3(9.1)	-0.825	-0.912	0.800	0.593	
9	bRPfxns_0813, 130C:127C, SNED1_KD_210:LM2.3(9.1)	0.090	0.053	0.553	0.353	
10	bRPfxns_0813, 131:127C, SNED1_KD_211:LM2.3(9.1)	-0.186	-0.219	0.628	0.359	
11						

# Create horizontal layout file: reporter\_sample\_template.txt

## Create a reporter\_sample\_template.txt file (horizontal layout)

When the parseSMreport script is run to process a P/P Summary report it can update the column headers in the report to include meta information, like more specific sample names attached to each reporter ion. These names will also be propagated through to the ratio distribution plots, and reports following normalization.

In order to provide the script with information on a sample name for each reporter ion you must create a **tab-delimited text file** in the SM directory that contains the initial report that is parsed. The text file must be named:

reporter\_sample\_template.txt

And be organized like the following example:

This example is for an experiment involving 18 samples that were run in 2 separate TMT-10 plexes, with 9 samples and a common control (mix of all 18).

1. line 1: reporter ion. The parser looks only for a substring in each cells text for the reporter ion (126, 127N, etc). The A and B are optional, and added only for user convenience.
2. line 2: sample name, free text.
3. line 3: directory name, directory the data was present in when the report was generated. This is necessary to match up the name and reporter ion, when more than 1 directory contributed to the report.

### Special considerations;

- Do not use a colon character : in any cells of the report, because colon will later be inserted by the parser to denote the numerator : denominator involved in a ratio.
- The examples have the columns in mass order. However, you can put the columns in any order. The parsed report will come out with the same column order as the reporter\_sample\_template.txt.
- When saving reporter\_sample\_template.txt in Excel  
Save as type:  
Text (Tab delimited) (\*.txt)  
not  
Unicode Text (\*.txt)  
You can diagnose problems caused by saving as Unicode Text by the funky characters in the script output pane in the browser in the reporterIons field.  
numerator: 126 denominator: 127C reporterIons: ÿþ1

A126	A127N	A127C	A128N	A128C	A129N	A129C	A130N	A130C	A131	B126	B127N	B127C	B128N	B128C	B129N
MD-8214C long survival, no chemo	MD-8226C long survival, unknown chemo	NP-8932N	PC-8592T	PI-8592N2	PI-8762N	PD-8832T long survival, palliative chemo	PD-8715 short survival, adj chemo	WD-8996C long survival, adj chemo	CC-all18	MD-7800T, long survival, unknown chemo	MD-1044C short survival, adj chemo	NP- 8926N0308	PC-1254T	PI-8926C	PD-8216C long survival, adj chemo
TMT10A_bRP	TMT10A_bRP	TMT10A_bRP	TMT10A_bR	TMT10A_bRP	TMT10A_bRP	TMT10A_bRP	TMT10A_bRP	TMT10A_bRP	TMT10A_bRP	TMT10B_bRP	TMT10B_bRP	TMT10B_bRP	TMT10B_bRP	TMT10B_bRP	TMT10B_bR

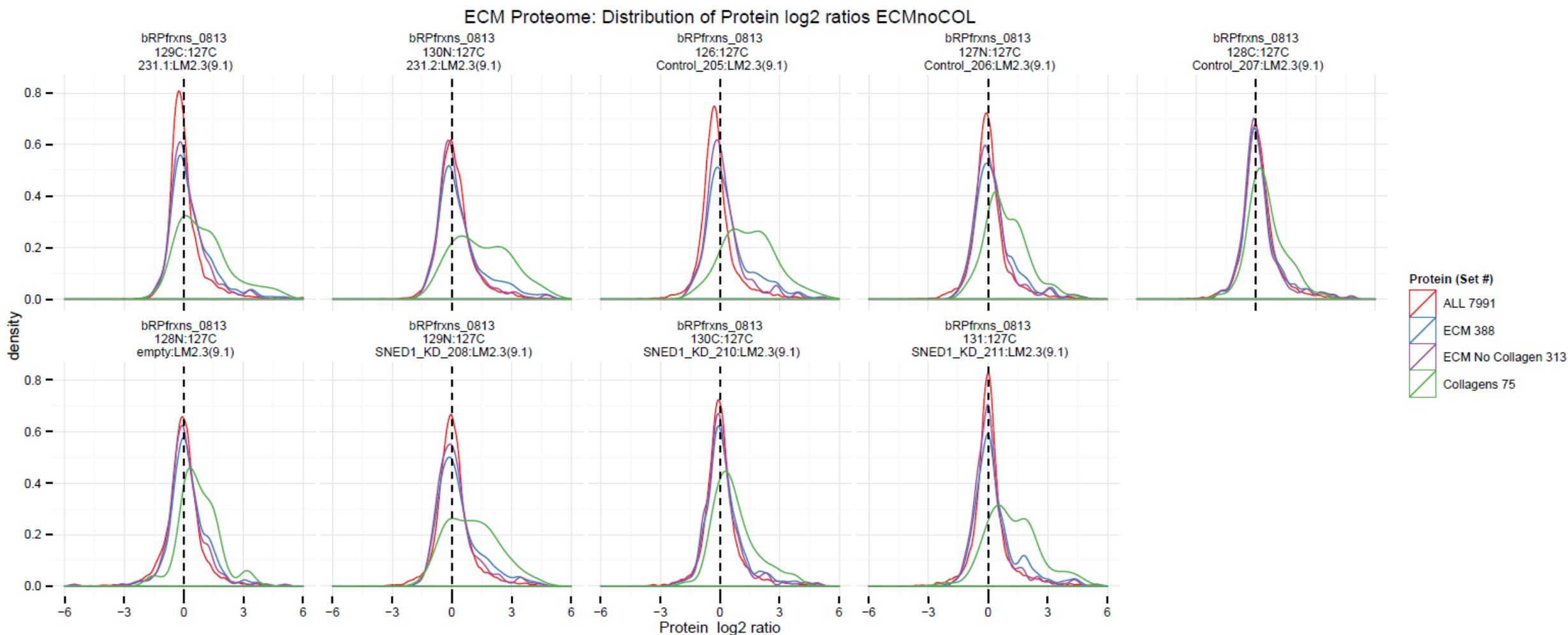
### MeanMulti

If the control ion was designated as MeanMulti, when generating the SM report, this leads to the mean intensity of multiple reporter ions being used as the denominator in ratios for each of the TMT10 mass labels. Consequently, there will be 1 additional ratio for each directory in the report. In the corresponding reporter\_sample\_template, indicate which ions were used as the denominator by joining the ions with a dot character (and keep them in alphabetic order, to match the column headers in the .ssv file). The dot delimiter enables the parser to prevent the parser from expecting this channel to also be used as a numerator.

126	127N	127C	128N	128C	129N	129C	130N	130C	131	126.127N.128C
Control_205	Control_206	LM2.3(9.1)	empty	Control_207	SNED1_KD_208	231.1	231.2	SNED1_KD_210	SNED1_KD_211	Control_3_mean
bRPfrxns_0813										

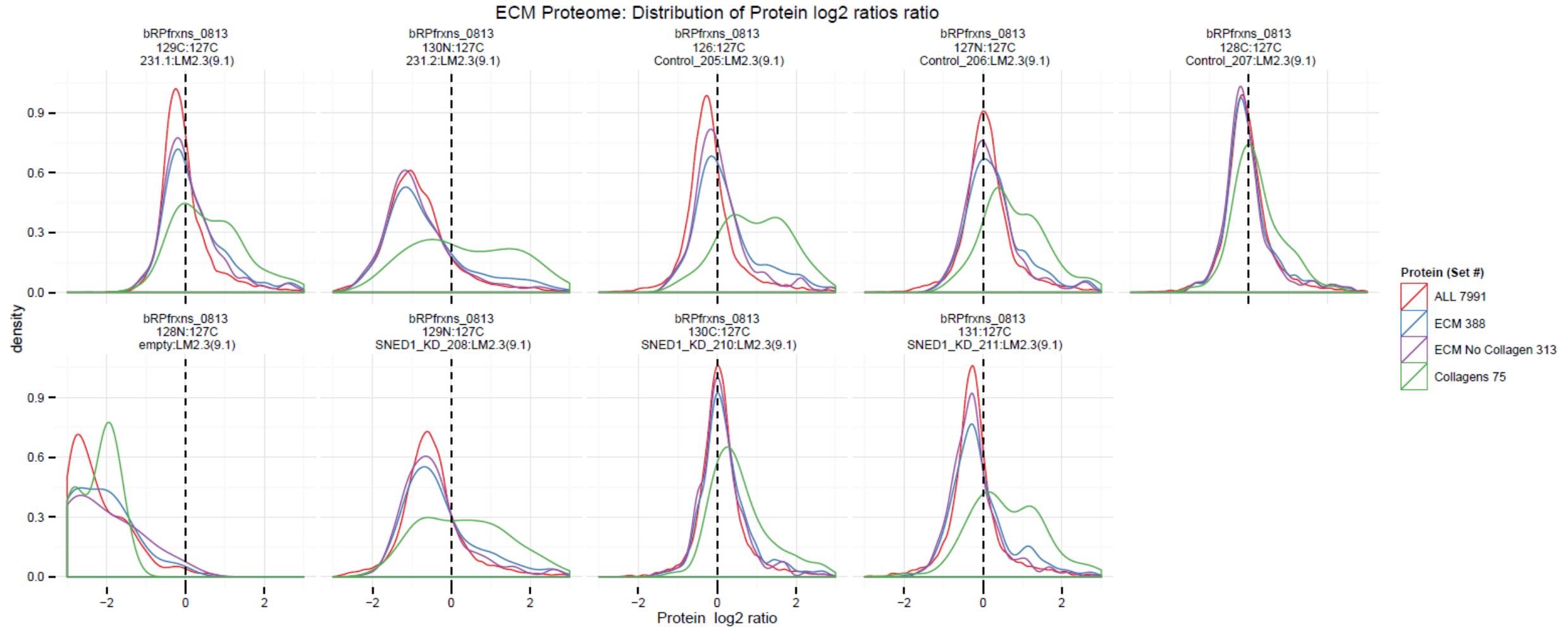
# TMT ratio Distributions – Normalized by median/SD of ECM(no Coll) - SGS

proteinProteinCentricColumnsExport.SGS.10-ratio-normmedianSD-ECMnoCOL-ECM-Protein-distributions4EnC.pdf



# TMT ratio Distributions – Unnormalized LM2 denom - SGS

proteinProteinCentricColumnsExport.SGS.10-ratio-ECM-Protein-distributions4EnC.pdf



# Automation

Service Request Manager (SRM)  
and Workflows

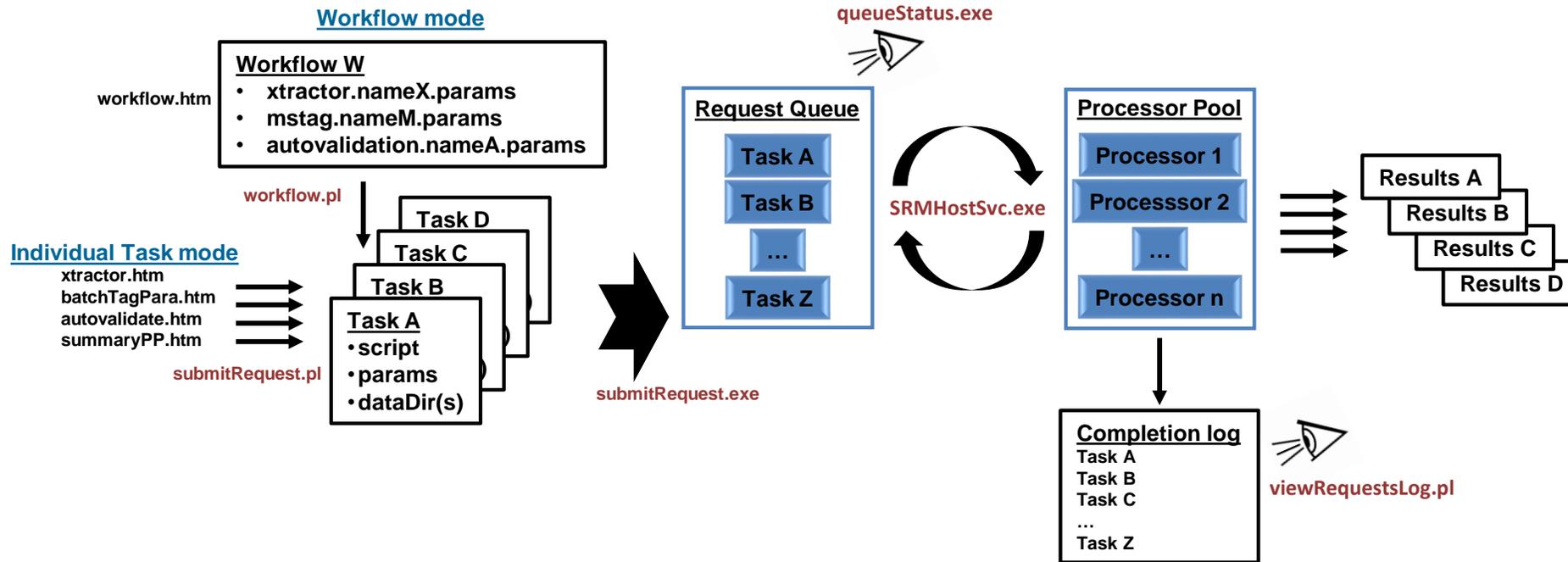
# Simple Workflow – Broad Institute QC sample

The screenshot displays the Spectrum Mill web interface. At the top, there's a navigation menu with tabs like 'Spectrum Mill', 'Request Queue', 'Completion Log', 'Extractor', 'MS/MS Search', 'Autovalidation', 'Quality Metrics & FDR', 'Protein/Peptide Summary', 'Tool Belt', 'Archive Data', and 'Help'. Below this, the 'Data Directories' section shows a selected directory: 'JurkatQC/Lacks/2021Q4/20210814'. The 'Workflow' section contains an 'Execute' button, 'Maximize CPUs' options, and a list of tasks. A red arrow points to the 'Tasks' dropdown menu, which lists various tasks such as 'JurkatJurkat\_ETD-QC\_QE\_EMAQ', 'JurkatJurkat\_ETHcD-QC\_QE\_EMAQ', 'JurkatJurkat\_QC\_QE\_EMAQ\_2020', 'JurkatJurkat\_QC\_QE\_EMAQ\_2020v430', 'JurkatJurkatTMT10\_QC\_QE\_EMAQ', 'JurkatJurkatTMTpro\_QC\_QE\_EMAQ', 'KarlHLA\_I\_IAA\_V3\_Enshg19B721glvarSmorfNuorfRI\_EMA', and 'KarlHLA\_I\_noenzyme\_IAA\_V3\_Enshg19B721glvarSmorfNuorf\_EMA'. Below the workflow section, a 'Read-only preview of parameters for selected task' is shown for the task 'Spectrum Mill - MS/MS Search - JurkatUniprotHuman20\_20\_20\_Amnqc\_SPI40\_QE\_HCDv430'. This preview includes search parameters like 'Validation filter', 'Batch size', 'Search previous hits', 'Database', 'Species', 'Digest', 'Maximum # missed cleavages', and 'Modifications'. At the bottom, there are sections for 'Search Criteria' (Matching Tolerances, Search Mode, Spectral Quality Filtering) and 'Data Files'.

List of tasks in workflow

Task parameter preview in normal layout (inactive)

# Mechanisms for Controlling Processes - Service Request Manager (SRM)



SRM coordinates task execution on the queue to maximize CPU usage while maintaining workflow order dependencies

Tasks in each user's workflow execute in serial

- `runXtractor.pl` → `*.mzXML, specFeatures.1.tsv`
- `batchTagPara.pl` → `spo.zip, tagSummary.1.tsv`
- `validateTable.pl` → `hitTable.1.tsv, spectrumTable.1.tsv`
- `summaryPP.pl` → `reportMode.1.ssv`

*eyes peeking, not rocket ships blasting off*

# SRM Request Queue Status

Workflow Requests - Spectrum Mill - Google Chrome

Not secure | manzano/millhtml/requests.htm

Spectrum Mill - Request Queue **manzano**

**Request Queue** Completion Log Help

There are 8 requests in the queue. Available memory: 25.8 Gb of 51.5 Gb

Remove

#	Task Id	Status	Task Type	Data Directory	Dependencies	User
<input type="checkbox"/>	1 210817081507.24689	Queued	Autovalidation	CPTAC3/HNSCC/Phosphoproteome/17CPTAC_HNSCC_Phosphoproteome_JHU_20190922	210817081506.24673P	Anonymous gp174-e69
<input type="checkbox"/>	2 210817081508.24690	Queued	Autovalidation	CPTAC3/HNSCC/Phosphoproteome/17CPTAC_HNSCC_Phosphoproteome_JHU_20190922	210817081506.24673P	Anonymous gp174-e69
<input type="checkbox"/>	3 210817160303.27503	Queued	Autovalidation	Meg/SoukasNT/Nucleus/Nucleus_Proteome	210817160126.27501P	Anonymous wp707-416
<input type="checkbox"/>	4 210817160324.27504	Queued	Autovalidation	Meg/SoukasNT/Nucleus/Nucleus_Proteome	210817160126.27501P	Anonymous wp707-416
<input type="checkbox"/>	5 210817160345.27505	Queued	Autovalidation	Meg/SoukasNT/Nucleus/Nucleus_Proteome	210817160126.27501P	Anonymous wp707-416
<input type="checkbox"/>	6 210817160449.27506	Queued	Quality Metrics	Meg/SoukasNT/Nucleus/Nucleus_Proteome	210817160126.27501P	Anonymous wp707-416
<input type="checkbox"/>	7 <a href="#">210817081506.24673P Monitor</a>	Running (929:1014)	MS/MS Search	CPTAC3/HNSCC/Phosphoproteome/17CPTAC_HNSCC_Phosphoproteome_JHU_20190922	210817081504.24671P	Anonymous gp174-e69
<input type="checkbox"/>	8 <a href="#">210817160126.27501P Monitor</a>	Running (986:1583)	MS/MS Search	Meg/SoukasNT/Nucleus/Nucleus_Proteome	none	Anonymous wp707-416

User A's workflow executes in parallel with user B's workflow

Any user can add/remove tasks from queue

# SRM Completion Log

ID	Task Name	Status	Date	User
210806213119.9974	MS/MS Search CPTAC3pancancer/BRCA/Proteome_GencodeGS/04CPTAC_BCProspective_Proteome_BI_20160930	NoError	Fri Aug 06 21:53:04 2021	Anonymous shiva
210806213118.9971P	Extraction CPTAC3pancancer/BRCA/Proteome_GencodeGS/04CPTAC_BCProspective_Proteome_BI_20160930	NoError	Fri Aug 06 21:51:17 2021	Anonymous shiva
210806213118.9972	Extraction CPTAC3pancancer/BRCA/Proteome_GencodeGS/04CPTAC_BCProspective_Proteome_BI_20160930	NoError	Fri Aug 06 21:36:44 2021	Anonymous shiva
210806211903.9941	MS/MS Search CPTAC3pancancer/BRCA/Proteome_GencodeGS/03CPTAC_BCProspective_Proteome_BI_20160916	NoError	Fri Aug 06 21:35:36 2021	Anonymous shiva
210806211901.9938P	Extraction CPTAC3pancancer/BRCA/Proteome_GencodeGS/03CPTAC_BCProspective_Proteome_BI_20160916	NoError	Fri Aug 06 21:33:55 2021	Anonymous shiva
210806211901.9939	Extraction CPTAC3pancancer/BRCA/Proteome_GencodeGS/03CPTAC_BCProspective_Proteome_BI_20160916	NoError	Fri Aug 06 21:22:02 2021	Anonymous shiva
210806205252.8006	MS/MS Search CPTAC3pancancer/BRCA/Proteome_GencodeGS/02CPTAC_BCProspective_Proteome_BI_20160913	NoError	Fri Aug 06 21:02:42 2021	Anonymous shiva
210806205251.8003P	Extraction CPTAC3pancancer/BRCA/Proteome_GencodeGS/02CPTAC_BCProspective_Proteome_BI_20160913	NoError	Fri Aug 06 21:00:57 2021	Anonymous shiva
210806205251.8004	Extraction CPTAC3pancancer/BRCA/Proteome_GencodeGS/02CPTAC_BCProspective_Proteome_BI_20160913	NoError	Fri Aug 06 20:53:02 2021	Anonymous shiva
210806184914.6037	P/P Summary CPTAC3pancancer/GBM/Phosphoproteome_GencodeGS/01CPTAC_GBM_Phosphoproteome_PNNL_20190123	NoError	Fri Aug 06 20:21:47 2021	Anonymous shiva
210806185535.6048	MS/MS Search CPTAC3pancancer/BRCA/Proteome_GencodeGS/01CPTAC_BCprospective_Proteome_BI_20160911	NoError	Fri Aug 06 18:57:19 2021	Anonymous shiva
210806185534.6046	Extraction CPTAC3pancancer/BRCA/Proteome_GencodeGS/01CPTAC_BCprospective_Proteome_BI_20160911	Error	Fri Aug 06 18:55:34 2021	Anonymous shiva
210806185101.6039	Extraction CPTAC3pancancer/BRCA/Proteome_GencodeGS/01CPTAC_BCprospective_Proteome_BI_20160911	Aborted	Fri Aug 06 18:55:09 2021	Anonymous shiva
210806185103.6041	MS/MS Search CPTAC3pancancer/BRCA/Proteome_GencodeGS/01CPTAC_BCprospective_Proteome_BI_20160911	Aborted	Fri Aug 06 18:55:09 2021	Anonymous shiva
210806185106.6044	Autovalidation CPTAC3pancancer/BRCA/Proteome_GencodeGS/01CPTAC_BCprospective_Proteome_BI_20160911	Aborted	Fri Aug 06 18:55:09 2021	Anonymous shiva
210806185105.6043	Autovalidation CPTAC3pancancer/BRCA/Proteome_GencodeGS/01CPTAC_BCprospective_Proteome_BI_20160911	Aborted	Fri Aug 06 18:55:09 2021	Anonymous shiva
210806185104.6042	Autovalidation CPTAC3pancancer/BRCA/Proteome_GencodeGS/01CPTAC_BCprospective_Proteome_BI_20160911	Aborted	Fri Aug 06 18:55:09 2021	Anonymous shiva
210806123207.4059	Autovalidation CPTAC3pancancer/BRCA/Proteome_GencodeGS/01CPTAC_BCprospective_Proteome_BI_20160911	Error	Fri Aug 06 18:55:09 2021	Anonymous shiva
210806123206.4058	Autovalidation CPTAC3pancancer/BRCA/Proteome_GencodeGS/01CPTAC_BCprospective_Proteome_BI_20160911	NoError	Fri Aug 06 18:42:50 2021	Anonymous shiva
210806123204.4057	Autovalidation CPTAC3pancancer/BRCA/Proteome_GencodeGS/01CPTAC_BCprospective_Proteome_BI_20160911	NoError	Fri Aug 06 18:42:26 2021	Anonymous shiva

If a task fails, the SRM aborts all subsequent tasks in the workflow





# Workflow and Saved Parameter files

## workflow.LSCC\_phospho\_TMT11\_v4\_Gencode\_nuORFs\_EMAA.tsv

requestScript	paramsFile
runXtractor.pl	CPTAC3human\xtractor.CPTAC-TMT11fullLysonly_HCD_v4_35_Xcent_800_6000_45sec_z6.params
batchTagPara.pl	CPTAC3human\mstag.LSCC_Phospho_TMT11FullLysOnly_Phospho_v4_CU_AmqcstynG_GencodeNuORFs.params
validateTable.pl	CPTAC3human\autovalidation.CPTAC_peptidevalidation_z24_acrossRun_MSL7_0_8_BCS3.params
validateTable.pl	CPTAC3human\autovalidation.CPTAC_peptidevalidation_z56_acrossDir_MSL7_0_4_BCS3.params

Workflow file lists each task and its parameter file name

Workflows and saved parameters files located at:  
SpectrumMill/millauto

## mstag.LSCC\_Phospho\_TMT11FullLysOnly\_Phospho\_v4\_CU\_AmqcstynG\_GencodeNuORFs.params

```
mparams_dir=mparams_mill%2F
seqdb_dir=D%3A%5CSeqDB%5C
requestScript=batchTagPara.pl
input_program_name=mstag
removeResults=1
hide_html_links=1
validationState=spectrum-not-marked-sequence-not-validated
batchSize=500
max_reported_hits=5
database=Gencode_v34_3nr.602contams.2043smorfs.nuORFv1.
110LSCCs.fasta
full_mw_range=1
full_pi_range=1
results_to_file=1
use_instrument_ion_types=1
species=All
enzyme=Trypsin%20allow%20P
missed_cleavages=4
fixedMods=carbamidomethylationCU
fixedMods=TMT11-Full-Lys
varMods=Acetyl
varMods=Oxidized-Methionine
varMods=pyroGlu
varMods=Deamidated-NG
varMods=Phosphorylated-S
varMods=Phosphorylated-T
varMods=Phosphorylated-Y
varMods=pyroCarbamidomethylCys
instrument_nameMSMS=3
instrument_name=ESI-QEXACTIVE-HCD-v4-35-20
minMatchedPercent=30
parent_mass_convert=monoisotopic
parent_mass_tolerance=20
tolerance_units_precursor=ppm
fragment_mass_tolerance=20
tolerance_units=ppm
max_ms_prod_charge=3
search_type=variable
enable_multiple_mods=1
user_min_parent_shift=-18.0
user_max_parent_shift=272
parent_shift_type=%2B%2F-
parent_shift=130.0
mutation_matrix_off=0
star_ions_off=1
homology_gap_mode=0
unknome=1
dissociationMethod=ALL
spectrumFileNames=%2A.pkl
contaminantProductionsScoreFilter=4.5
=
```

Parameter file lists:  
name=value  
pairs  
Perl-CGI style

Same as .params files  
written in each SM  
data directory



# Quality Metrics

## Troubleshooting

- LC Gradient and column performance
- Measuring data acquisition strategy changes
- Mass calibration drift
- MS ion optics performance
- Digestion efficiency
- IMAC phospho enrichment
- TMT/iTRAQ labeling efficiency
- Cysteine reduction/alkylation failure



# Quality Metrics & FDR

Spectrum Mill - Quality Metrics and False Discovery Rates (FDR) - KarlIMzC\_PIP\_PAU\_Chrom\_Mods\_TMT10\_131

Spectrum Mill | Extractor | MS/MS Search | Autovalidation | Protein/Peptide Summary | Workflows | Tool Belt | Spectrum Summary | Help

Report  Queue request  Excel Export Save As... Load...

<b>FDR Metrics (spectra, peptide, protein)</b> <input checked="" type="checkbox"/> FDR at the peptide & spectra level (from valid hits) <input type="checkbox"/> FDR at the protein level <input type="checkbox"/> Group proteins across all directories Grouping method: 1 shared, expand subgroups ▾	<b>MS/MS Interpretation Metrics</b> <input checked="" type="checkbox"/> Identification Scores <input type="checkbox"/> Fragmentation Metrics <input type="checkbox"/> Fragmentation Mode <input type="checkbox"/> Variable modification site localization sitly ▾	<b>Peptide Separation Metrics</b> <input checked="" type="checkbox"/> Chromatography metrics for each run Middle RT portion of validated spectra in LC run: 90 ▾ (%) <input type="checkbox"/> Peptide pI median for each run <input type="checkbox"/> Peptide subset reports for each data directory (seqdb/peptideQMLists/*.txt)
<b>Precursor Ion Metrics</b> <input checked="" type="checkbox"/> Precursor mass error mean (ppm) <input checked="" type="checkbox"/> Precursor charge count (from valid spectra) <input checked="" type="checkbox"/> Precursor Isolation Purity & Averagine Chi <sup>2</sup> <input checked="" type="checkbox"/> Precursor Acquisition Uncertainty: m/z and z <input type="checkbox"/> MS1 intensity stats for each run. <input type="checkbox"/> Precursor Ion Fragmentation Table	<b>MS/MS Spectral Identifiability Metrics (define thresholds below)</b> <input checked="" type="checkbox"/> Sequence tag length: > 3 <input checked="" type="checkbox"/> Precursor isotope quality XIC's (Chi <sup>2</sup> vs. Averagine): > 0.7 <input checked="" type="checkbox"/> Precursor Isolation Purity: > 70 %	<b>Sample Handling Metrics</b> <input checked="" type="checkbox"/> Isobaric label incorporation for each run TMT 10 ▾ Control: 126, 127N, 127C <input type="checkbox"/> Digestion completeness <input checked="" type="checkbox"/> Observed modifications by: Peptide Spectrum Matches ▾
<b>Peptide Fraction Overlap</b> <input type="checkbox"/> Distinct Peptide Fraction Overlap Table	Distinct peptide comparison method: Case Sensitive(CS) ▾ For fraction overlap, sample handling, or pI choose CS or CI. FDR calculations always use CI. Filtering to distinct peptides retains each highest scoring representative after CS or CI string comparison of sequences. Variable mods are lowercase.	

**Data Directories**  
Select ...  Development/ECM/PDAC\_human2019

Concepts for some of these metrics were developed and described in: Rudnick PA, Clauser KR, Kilpatrick LE, et. al., "Performance metrics for liquid chromatography-tandem mass spectrometry systems in proteomics analyses", Mol Cell Proteomics. 2010 Feb;9(2):225-41 (<http://www.ncbi.nlm.nih.gov/pubmed/19837981>).



# MS and Chromatography Metrics

Measure effect of changes in acquisition parameters and chromatography on newly installed QExactive Plus at Broad Institute

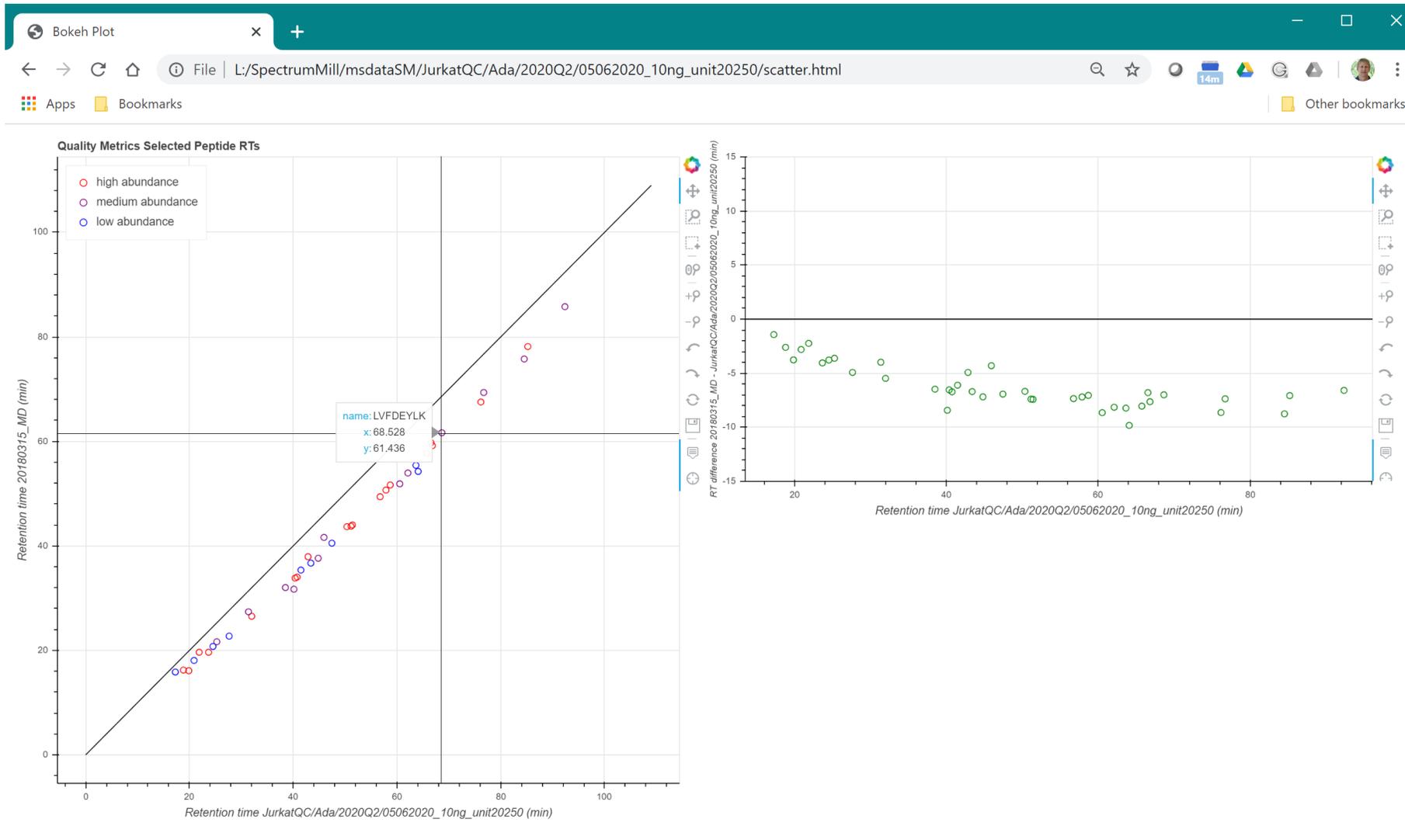
Comment	Pep Match	MS/MS spectra collected C	MS/MS spectra collected MI m/z = 0.0	MS/MS spectra filtered F	MS/MS spectra valid V	MS/MS spectra valid MI m/z = 0.0	Median ID Score	Mean Precursor Mass Error (ppm)	Start time mid 90 % run (min)	End time mid 90 % run (min)	Time span mid 90 % run (min)	Gradient Shape mid 90 % filtered spectra in run	Median MS1 peak width mid 90 % run (sec)	Total precursor XIC mid 90 % matched spectra in run	PIP bin 1 100-90 Spec tra (%)	PIP bin 6 <50 Spec tra (%)	Distinct Peps CS Total (#)	FDR Spec tra (%)
	on	52,270	-	47,892	34,811	-	12.07	0.3	21	98	78	9997754	12	6.15E+13	22.1	9.7	27,929	0.72
Isolation width decreased 2.5 to 2.0	on	51,827	-	49,963	36,876	-	12.68	1.1	17	98	81	9987753	12	6.31E+13	29.7	6.4	29,598	0.73
LC-Klaus	on	49,015	-	47,140	32,487	-	11.73	-0.2	18	98	80	9999876	15	5.35E+13	26.1	9.4	26,176	0.78
	pref	60,705	19,121	57,497	38,127	8,198	11.79	2.2	21	100	79	9897654	11	4.88E+13	27.0	7.6	30,759	0.87
16 min longer gradient, isolation width decreased 2.0 to 1.6	pref	68,003	23,946	63,633	42,998	10,670	12.36	-0.4	20	115	94	8898754	13	6.13E+13	35.0	4.8	34,719	0.96

Pep Match (on/pref)  
Precursor isotope cluster

Lengthened  
Gradient

Narrowed  
isolation width

# Quality Metrics - Selected Peptide RTs vs gold standard run



Gold Standard

Current Experiment

Interactive graphics  
Python script feeds  
Bokeh  
javascript library

Namrata Udeshi

SeqDB\peptideQMLists\\*.txt  
Lists of readily observable  
human peptides in cell  
lysates.



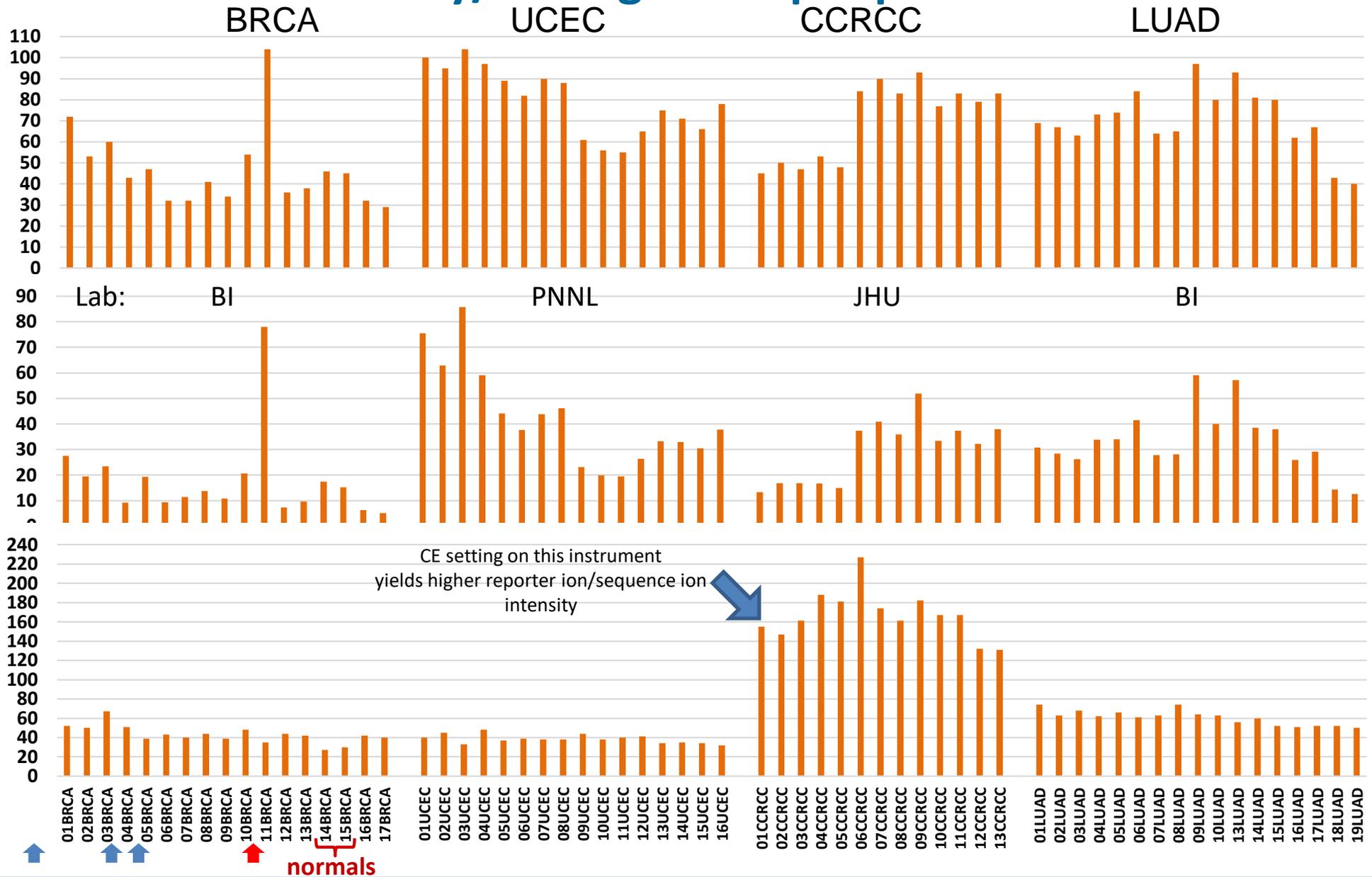
All are Thermo Fusion Lumos

# Instrument Sensitivity/Tuning - Phosphoproteome

Median MS2 fill time mid 90% matched spectra (msec)

Spectra Reaching max MS2 fill time mid 90% matched spectra (%)

Median S/N All Reporters

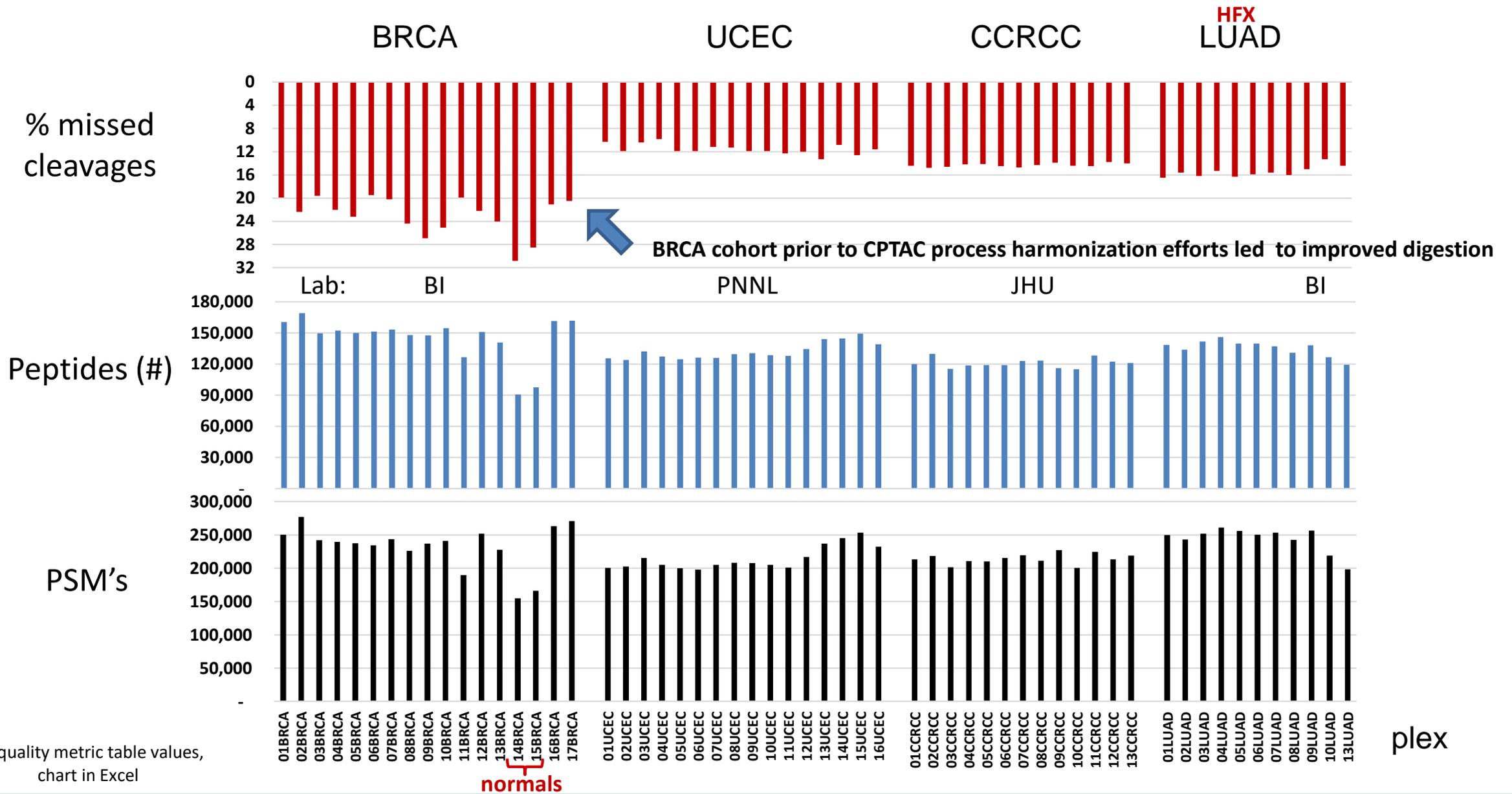


Instrument Maintenance

- Ion optics cleaning
- Replaced turbo & ion funnel

plex

# Tracking Digestion Efficiency



SM quality metric table values, chart in Excel

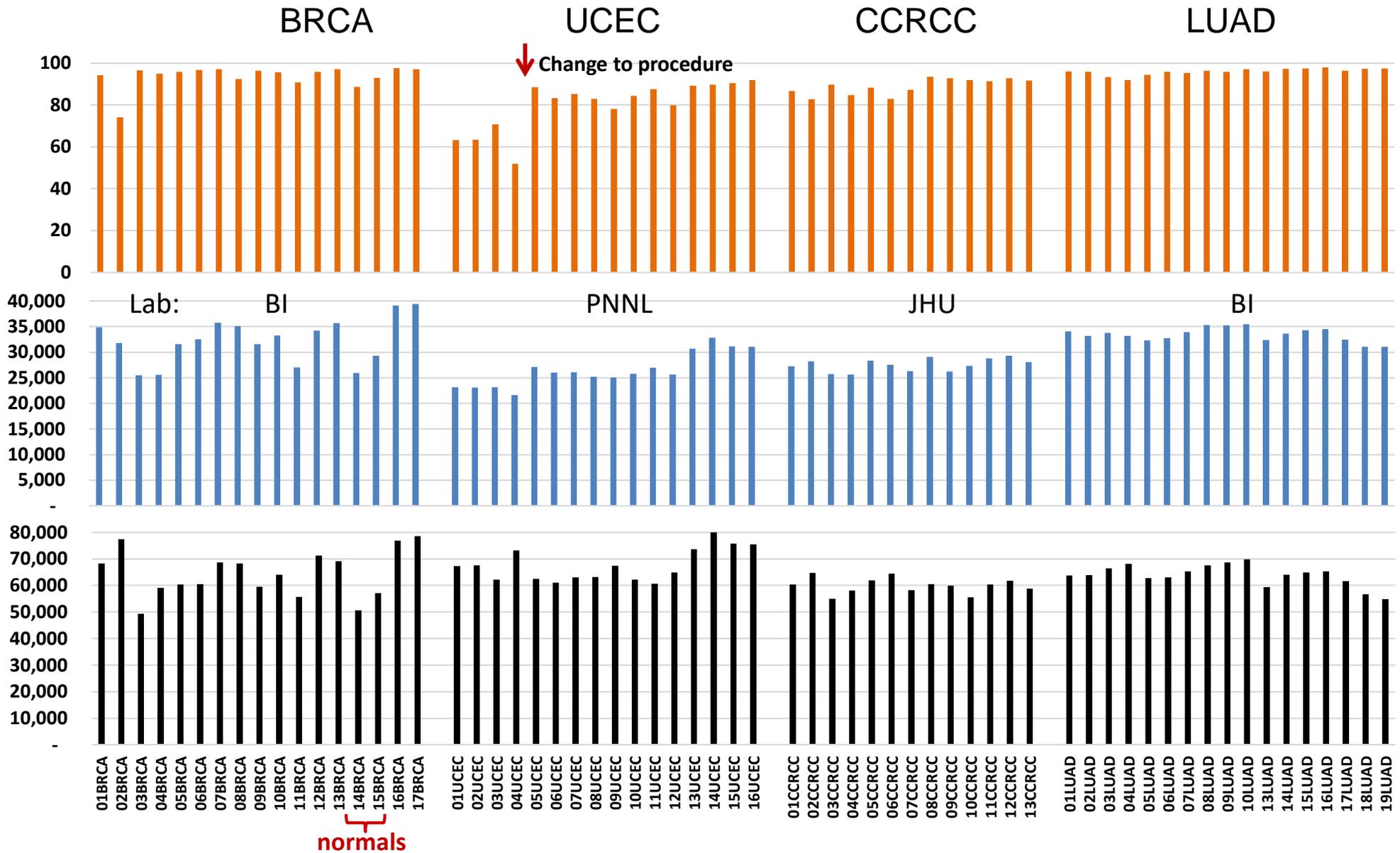
plex

# Tracking IMAC Phospho -enrichment

% enrichment  
sty PSM's  
all PSM's

Phospho  
Peptides (#)

PSM's



plex

SM quality metric table values,  
chart in Excel



# Phosphorylation Quality Metrics

Directory	Raw Files	MS/MS spectra valid V	s t y Sites spectra (%)	s t y Sites spectra (#)	s t y Sites Localized spectra (%)	s t y Sites Distinct Peptides (CI#)	s t y Sites Distinct Peptides (CI%)	s t y Sites FDR Spectra (%)	s t y Sites FDR Distinct Peptide (%)	s t y Sites Distinct Peptides (CI#)	s t y Sites FDR Spectra (%)	s t y Sites FDR Distinct Peptide (%)	s t y Sites Distinct Peptides (CI#)
ome2/01TCGA	13	77,152	94.5	72,930	56.4	32,179	92.4	0.47	0.92	32,179	0.51	0.98	34,830
ome2/02TCGA	13	82,578	90.0	74,288	55.7	31,843	85.9	0.51	1.00	31,843	0.52	0.98	37,072
ome2/03TCGA	13	91,894	77.8	71,460	57.0	31,193	75.6	0.48	0.95	31,193	0.50	0.96	41,256
ome2/04TCGA	13	98,744	72.6	71,690	57.9	31,921	68.1	0.57	1.12	31,921	0.51	0.96	46,887
ome2/05TCGA	13	90,028	83.4	75,068	57.6	31,748	78.2	0.49	1.05	31,748	0.50	0.98	40,620
ome2/06TCGA	13	96,310	53.2	51,264	58.7	22,933	51.4	0.64	1.18	22,933	0.51	0.93	44,613
ome2/07TCGA	13	83,351	73.0	60,858	58.4	26,668	68.9	0.54	1.03	26,668	0.47	0.86	38,695
ome2/08TCGA	13	102,722	66.7	68,521	55.1	30,355	66.3	0.65	1.23	30,355	0.49	0.92	45,811
ome2/09TCGA	13	87,016	67.1	58,360	55.5	27,513	65.1	0.62	1.19	27,513	0.51	0.91	42,250
ome2/10TCGA	13	88,406	78.5	69,368	55.7	29,632	74.7	0.58	1.15	29,632	0.52	1.00	39,650
ome2/11TCGA	13	60,593	92.5	56,029	49.4	25,075	90.0	0.52	0.95	25,075	0.50	0.89	27,864
ome2/12TCGA	13	63,858	94.1	60,117	50.5	26,540	92.5	0.45	0.81	26,540	0.45	0.80	28,689
ome2/13TCGA	13	81,810	85.2	69,701	49.5	28,186	79.7	0.51	0.93	28,186	0.49	0.86	35,384
ome2/14TCGA	13	68,580	90.7	62,196	48.5	26,335	87.6	0.49	0.92	26,335	0.49	0.91	30,077
ome2/15TCGA	13	80,970	95.0	76,933	57.2	31,203	92.8	0.47	0.97	31,203	0.50	1.02	33,629
ome2/16TCGA	13	83,601	93.1	77,874	56.2	32,908	90.0	0.46	0.91	32,908	0.49	0.91	36,573
ome2/17TCGA	13	71,216	94.5	67,305	56.7	29,337	91.9	0.48	0.95	29,337	0.49	0.93	31,907
ome2/18TCGA	13	71,290	93.9	66,930	53.0	27,688	91.6	0.47	0.91	27,688	0.49	0.93	30,234
ome2/19TCGA	13	85,670	92.5	79,252	55.9	35,012	90.2	0.5	0.98	35,012	0.51	0.97	38,814
ome2/20TCGA	13	85,860	94.7	81,314	55.8	35,404	93.1	0.48	0.96	35,404	0.51	1.00	38,046
ome2/21TCGA	13	90,324	96.6	87,277	55.3	35,215	95.1	0.47	0.98	35,215	0.50	1.04	37,043
ome2/22TCGA	13	83,748	97.5	81,684	56.0	34,594	96.3	0.46	0.97	34,594	0.49	1.01	35,935
ome2/23TCGA	13	84,628	94.6	80,032	55.2	31,800	92.1	0.43	0.94	31,800	0.46	0.98	34,537
ome2/24TCGA	13	75,699	95.7	72,416	55.4	28,790	93.8	0.46	0.97	28,790	0.48	0.99	30,697

- Track level of phospho enrichment across IMAC experiments

- Separate FDR calculation for phosphopeptides and all peptides



# Determining if you should relabel your samples

## Zecha MCP 2019

		Possibilities		
N-term		K		
N-term		R		
<b>Fully</b>				
	N-term		K	1
	N-term		R	2
pyro, ac	N-term		K	3
<b>Partial</b>				
	N-term		K	1
	N-term		K	2
<b>No labeling</b>				
	N-term		K	1
	N-term		R	2
pyro, ac	N-term		K	3
<b>Can't be labeled</b>				
pyro, ac	N-term		R	1

## SM terminology conflicts with Zecha 2019

Labeled – PSMs from A+B

**Fully** labeled – PSMs from A excluding A3

No label – PSMs from C and D

Partially labeled – PSMs from B

**Completely** labeled: 100- (Partially labeled + No label)

### **Metrics to determine if the samples should be relabeled**

**Underlabeled** = B + C = Partially labeled + No label

**Completely labeled** = 100- (B + C) = 100 -(Partially labeled + No label)

# SM Quality Metrics – TMT labeling

	MS/MS spectra valid V	Isobaric Labeled Spectra (%)	Fully Labeled Spectra (%)	No Label Spectra (%)	Isobaric Under-labeled both Spectra (%)	Isobaric Completely Labeled Spectra (%)	Isobaric Only Lys Labeled Spectra (%)	Isobaric Only Nterm Labeled Spectra (%)	Isobaric Only Lys N-term Labeled Spectra (%)	Isobaric No Label Labeled Spectra (%)	Isobaric Under-labeled both total Spectra (%)	PSMs Containing s t y (#)													
													ProteomeDirectory	305,976	97.2	93.6	2.8	2.6	94.6	3.1	0.5	2.1	0.3	2.9	34855
													Karl/Zecha_MCP_2019/ProteomeOverlabel/01Pi	296,635	96.0	91.6	4.0	3.5	92.6	3.8	0.6	2.8	0.5	3.9	34391
													Karl/Zecha_MCP_2019/ProteomeOverlabel/Corr												
Proteome	PSMs	F	P	U	O	DP	FDR																		
Mertins	305,976	94.6	2.6	0.3	11.4	163,322	0.6																		
Zecha	296,635	92.6	3.5	0.5	11.6	168,606	0.6																		

	MS/MS spectra valid V	Isobaric Labeled Spectra (%)	Fully Labeled Spectra (%)	No Label Spectra (%)	Isobaric Under-labeled both Spectra (%)	Isobaric Completely Labeled Spectra (%)	Isobaric Only Lys Labeled Spectra (%)	Isobaric Only Nterm Labeled Spectra (%)	Isobaric Only Lys N-term Labeled Spectra (%)	Isobaric No Label Labeled Spectra (%)	Isobaric Under-labeled both total Spectra (%)	PSMs Containing s t y (#)												
													PhosphoDirectory	88,921	98.1	92.5	1.9	3.9	94.2	4.8	0.7	3.2	0.3	4.1
													Karl/Zecha_MCP_2019/Phosphoproteome/01Ph	90,734	97.2	89.3	2.9	6.2	91.0	7.0	0.8	5.4	0.3	6.4
													Karl/Zecha_MCP_2019/Phosphoproteome/Com											
Phospho	PSMs	F	P	U	O	DP	FDR																	
Mertins	88,921	94.2	3.9	0.3	-	42,124	0.9																	
Zecha	90,734	91.0	6.2	0.3	-	42,785	0.8																	

F	PSMs with all available amine groups are labeled. (includes fully labeled, blocked N-terms)
P	PSMs with a label, but also unlabeled amine on N-term or Lysine.
U	PSMs with no label, but have unlabeled amine on N-term or Lysine.
O	PSMs with a label on Ser,Thr, Tyr; in peptides that contain His

## Fully labeled - old

PSMs with an N-terminal label and a label on Lys (if present).

## Excludes blocked N-termini

(acetyl or pyro) on peptides with labeled lysines.

## Completely labeled - new

100 - No label - Partially labeled

## Includes blocked N-termini

(acetyl or pyro) on peptides with labeled lysines.

Too confusing!  
should eliminate column: Fully

# Tumor/NAT – Evaluate Relative Peptide Load Across Samples

	126C_	127N_	127C_	128N_	128C_	129N_	129C_	130N_	130C_	131N_
plex	131C									
1	1.38	0.97	1.07	0.89	1.09	1.03	0.97	0.99	1.25	0.94
2	1.14	0.99	1.35	1.07	1.11	0.84	0.94	1.02	1.25	0.94
3	1.17	1.14	1.12	0.92	1.04	0.98	1.05	1.00	0.95	1.04
4	1.16	1.18	1.09	0.99	1.07	0.97	1.05	1.00	1.03	1.05
5	1.13	1.17	1.05	0.86	0.97	1.00	0.90	0.99	1.05	1.07
6	1.17	0.87	1.07	0.85	1.05	0.87	1.10	0.93	1.22	0.89
7	1.10	1.04	1.14	1.12	1.19	0.96	1.01	1.00	1.08	1.10
8	1.00	0.97	0.91	0.91	1.00	0.86	0.96	0.89	0.97	0.92
9	1.17	1.10	1.24	0.99	1.02	0.99	1.08	1.02	1.08	0.94
10	0.95	1.02	1.15	1.01	0.94	1.05	1.05	1.02	1.01	0.96
11	1.17	0.91	1.13	1.04	1.06	1.09	1.10	0.82	1.07	0.95
12	1.09	1.05	0.96	0.98	1.09	1.03	0.98	1.00	1.10	1.03
13	1.08	1.01	1.13	1.04	1.19	1.09	1.12	1.02	1.12	1.09
14	1.13	1.00	1.14	0.97	0.99	0.93	0.92	0.90	0.92	0.85
15	1.10	1.09	1.05	0.95	1.09	1.02	1.15	0.99	1.32	1.06
16	1.07	0.94	1.02	0.95	1.14	0.96	1.08	0.92	1.03	0.90
17	1.21	1.10	1.09	0.98	1.04	1.00	1.07	0.93	0.93	0.93
18	1.31	1.08	1.27	1.10	1.09	0.94	1.15	0.89	1.15	0.97
19	1.14	1.07	1.15	1.03	1.05	1.01	1.05	1.01	1.04	0.84
20	1.07	0.99	1.06	0.97	1.21	0.96	1.09	0.99	1.07	0.91
21	1.17	1.03	1.16	1.13	1.10	1.00	1.07	0.97	1.08	1.09
22	1.16	1.04	1.20	0.81	1.09	1.01	1.25	0.98	1.08	0.98

LSCC

LUAD

Proteome

LSCC tighter thresholds applied for mixing controls  
+/- 15% (LSCC) from CR instead of +/-25%(LUAD)

	126C	127N	127C	128N	128C	129N	129C	130N	130C	131N
1	1.83	0.63	1.27	0.72	1.11	0.83	1.08	0.90	1.64	0.88
2	0.89	0.48	1.58	0.96	1.26	0.49	1.19	0.91	1.75	0.92
3	1.13	0.99	1.00	0.47	1.15	1.02	1.02	0.48	0.94	1.18
4	1.37	1.03	0.91	0.71	1.28	0.81	1.13	0.72	1.10	0.87
5	1.25	0.97	0.68	0.55	1.08	0.97	0.63	0.73	1.08	0.95
6	0.94	0.60	1.17	0.46	1.22	0.80	1.15	0.86	1.41	0.76
7	1.32	0.73	1.22	1.16	1.41	0.70	1.13	0.84	1.16	1.23
8	0.80	0.79	0.75	0.82	1.10	0.69	0.83	0.55	0.90	0.82
9	1.13	1.12	1.83	0.90	1.11	0.80	1.33	0.87	1.48	0.88
10	1.32	0.95	1.20	0.79	1.04	0.94	1.17	0.98	1.18	0.76
11	1.25	0.57	1.34	0.82	1.08	1.10	1.20	0.48	1.02	0.82
12	0.91	0.73	1.15	0.81	0.96	0.72	1.07	0.83	1.34	1.17
13	1.11	0.79	1.43	0.99	1.39	1.13	1.31	0.92	1.23	1.21
14	1.33	0.77	1.32	0.84	1.23	0.91	1.22	0.85	0.99	0.87
15	1.07	0.81	0.94	0.78	1.28	0.90	1.28	0.80	1.70	1.16
16	0.97	0.78	0.96	0.77	1.24	0.76	1.53	0.77	1.26	0.86
17	1.34	1.12	1.29	0.82	1.24	0.89	1.48	0.80	0.96	0.81
18	1.73	0.95	0.90	1.12	1.07	0.80	1.22	0.55	1.53	0.78
19	0.87	0.65	1.37	0.97	1.21	0.91	1.03	0.78	1.26	0.80
20	1.15	0.87	0.85	0.72	1.43	0.64	1.14	0.79	1.31	0.76
21	1.33	0.77	0.97	1.06	1.32	0.89	1.17	0.61	0.78	1.16
22	1.17	0.91	1.47	0.45	1.09	0.91	1.65	0.83	1.21	0.88

Phosphoproteome

Tumors tend to yield more phosphorylated material than NATs

	126C_	127N_	127C_	128N_	128C_	129N_	129C_	130N_	130C_	
plex	131	131	131	131	131	131	131	131	131	plex
1	1.30	0.98	0.90	0.92	0.99	0.97	1.07	0.90	1.21	1
2	1.50	1.10	1.17	0.67	0.89	0.80	1.08	1.14	0.99	2
3	1.28	1.16	1.17	1.02	1.08	0.90	1.10	1.03	1.12	3
4	1.26	1.01	1.25	1.11	1.00	0.91	1.15	1.01	1.00	4
5	1.33	1.26	0.83	0.91	0.93	0.81	1.12	0.97	1.12	5
6	1.20	1.05	1.09	1.10	0.89	1.03	1.21	1.00	1.34	6
7	1.27	1.11	1.15	0.93	1.19	0.81	1.12	0.90	0.89	7
8	1.24	1.31	1.06	0.96	1.05	0.88	1.01	1.00	1.08	8
9	1.17	0.59	1.11	0.92	0.90	0.80	0.95	0.56	1.08	9
10	1.22	1.07	1.14	0.90	0.91	0.88	1.01	0.87	1.12	10
11	1.21	1.16	1.22	0.83	1.13	0.70	1.16	0.71	0.92	11
12	1.01	0.88	1.04	0.68	1.16	0.99	1.03	0.64	0.90	12
13	1.09	0.98	0.82	0.87	0.90	0.96	0.98	0.85	0.94	13
14	1.16	1.16	0.98	0.91	1.12	0.84	1.05	0.99	1.10	14
15	1.35	1.04	1.13	1.01	1.15	0.90	1.11	1.01	1.21	15
16	1.41	1.03	1.13	0.95	1.18	0.93	1.07	0.95	1.01	16
17	1.11	1.11	1.20	1.15	1.05	1.08	1.08	0.92	1.10	17
18	1.16	1.08	1.06	1.07	1.07	0.87	1.00	0.66	1.06	18
19	1.29	1.14	1.20	1.01	0.97	0.64	1.11	0.93	1.31	19
20	1.30	1.25	1.24	1.13	1.19	0.94	0.99	0.89	1.07	20
21	1.47	1.14	1.18	0.96	1.20	0.87	1.15	1.03	1.22	21
22	1.37	0.99	1.29	0.96	1.08	0.57	1.43	0.77	1.46	22
23	1.22	0.84	1.09	0.79	0.93	0.80	1.24	0.99	1.18	23
24	1.00	1.02	1.06	0.88	1.02	0.60	1.11	0.86	0.82	24
25	0.96	0.86	1.09	0.85	1.39	0.89	1.14	0.90	0.92	25

	126C	127N	127C	128N	128C	129N	129C	130N	130C	
plex										plex
1	1.38	0.68	0.69	0.77	1.03	0.93	1.23	0.79	1.42	1
2	1.78	0.72	1.17	0.48	0.75	0.42	1.31	0.95	0.77	2
3	1.28	0.86	1.32	0.64	1.08	0.51	1.50	0.97	1.33	3
4	1.02	0.62	1.63	0.87	1.16	0.95	1.36	0.89	0.78	4
5	1.18	0.85	0.81	0.80	0.71	0.48	1.01	0.84	1.33	5
6	1.05	0.63	1.41	0.82	0.87	0.73	1.26	0.80	1.55	6
7	1.63	0.85	0.88	0.58	1.24	0.49	1.28	0.65	0.98	7
8	1.43	0.92	1.12	0.86	1.21	0.72	1.16	0.63	1.36	8
9	1.34	0.37	1.39	0.72	1.22	0.67	1.07	0.25	1.29	9
10	1.03	0.76	1.40	0.61	1.28	0.66	0.92	0.39	1.70	10
11	1.19	1.01	1.30	0.59	1.15	0.53	1.66	0.35	0.73	11
12	1.21	0.61	1.38	0.35	1.46	0.90	0.83	0.27	0.91	12
13	1.12	0.36	0.58	0.44	1.13	0.90	1.45	0.65	1.18	13
14	1.24	0.98	1.15	0.53	1.47	0.63	1.27	0.75	1.32	14
15	1.85	0.87	1.48	0.86	1.56	0.89	1.11	0.93	1.41	15
16	1.91	1.08	1.74	0.59	1.75	0.69	1.39	0.67	1.25	16
17	1.13	0.97	1.52	0.95	1.21	1.02	1.38	0.87	1.65	17
18	1.14	0.88	1.48	0.93	1.14	0.77	1.04	0.35	0.89	18
19	1.33	1.06	1.29	0.89	1.04	0.34	1.03	0.83	0.88	19
20	1.34	1.24	1.76	1.14	1.66	0.90	1.33	0.80	1.30	20
21	1.68	1.03	1.23	0.62	1.20	0.69	1.08	1.03	1.37	21
22	1.54	0.80	1.64	0.98	0.91	0.28	1.65	0.69	1.20	22
23	1.59	0.60	1.46	0.57	1.15	0.71	1.77	0.80	1.52	23
24	0.95	0.87	1.03	0.75	1.32	0.39	1.01	0.74	0.70	24
25	1.15	0.46	1.18	0.82	1.25	0.82	1.48	0.74	0.74	25

T N T N T N T N T T/N

T N T N T N T N T/N

## SM quality metric – TMT channel balance

PSM level reporter ion ratio  
(summed intensity of all PSMs  
Each reporter ion/common ref )

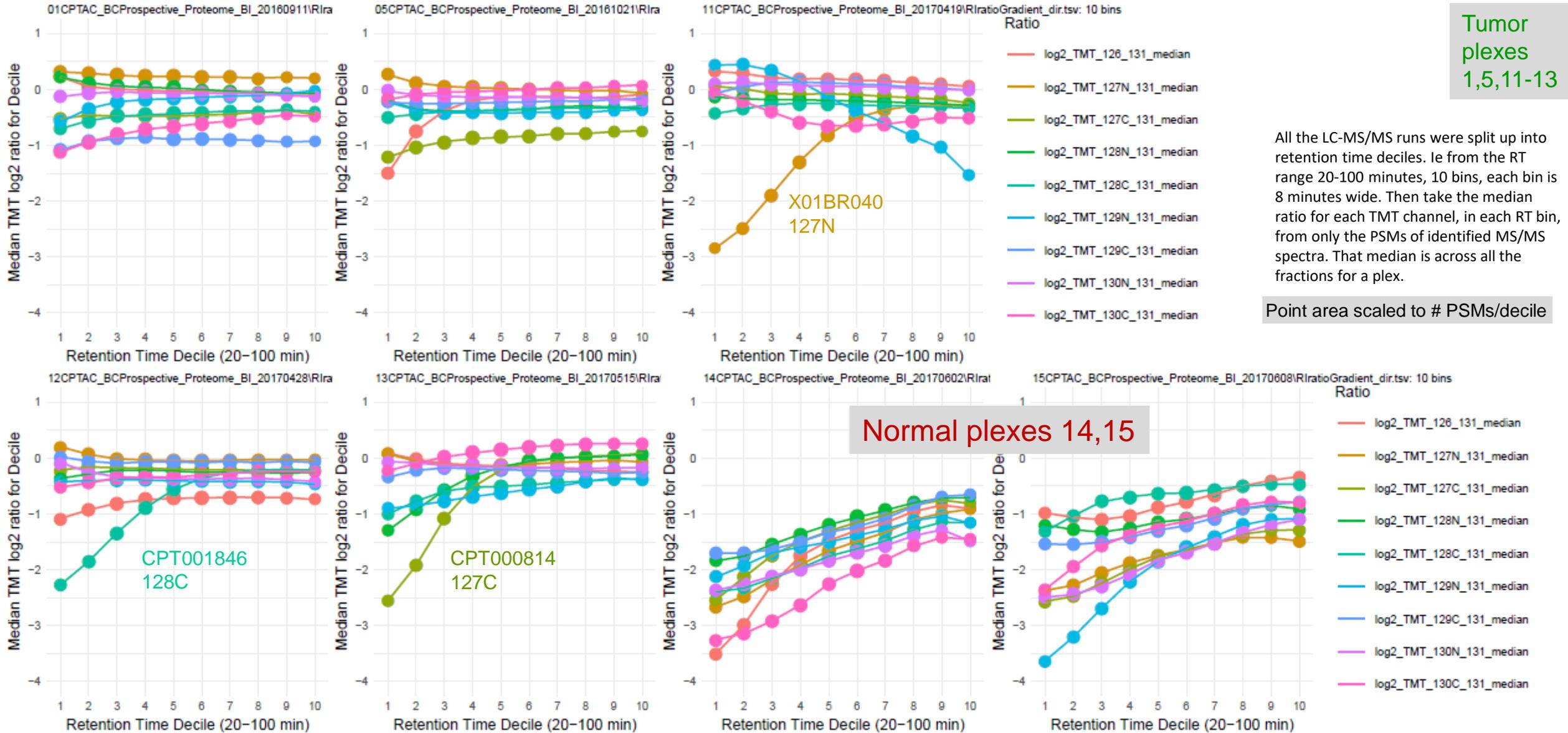
SM quality metric table values, coloring done in Excel



# SM Quality Metrics Plots – TMT Ratio vs Retention Time Decile

Proteome

Tumor plexes 1,5,11-13



All the LC-MS/MS runs were split up into retention time deciles. I.e from the RT range 20-100 minutes, 10 bins, each bin is 8 minutes wide. Then take the median ratio for each TMT channel, in each RT bin, from only the PSMs of identified MS/MS spectra. That median is across all the fractions for a plex.

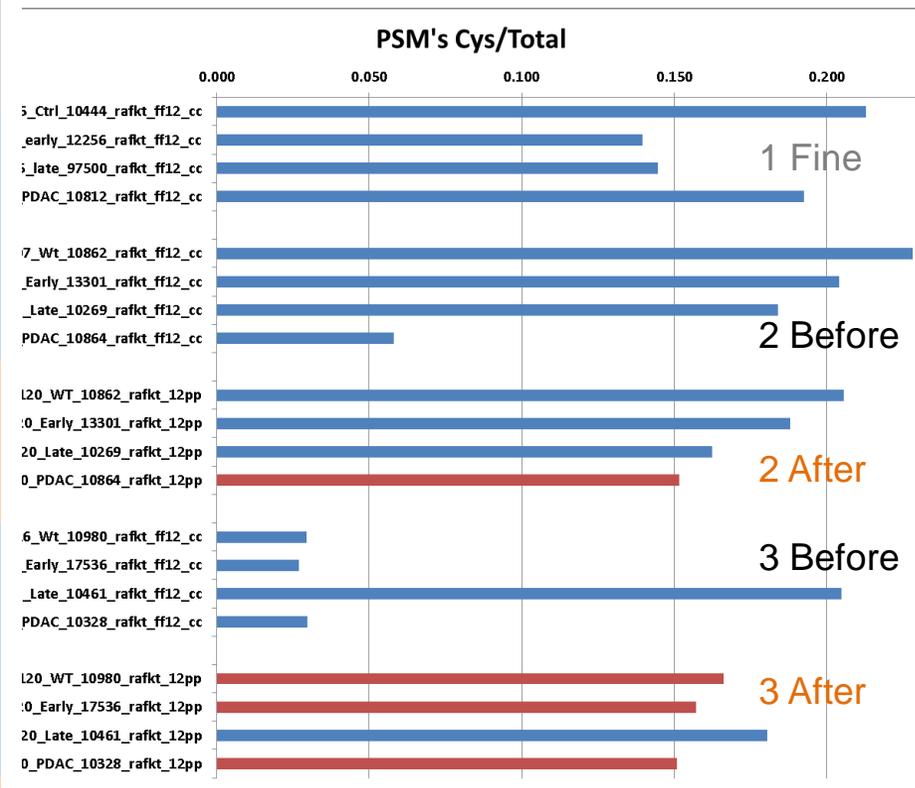
Point area scaled to # PSMs/decile



# Reduction/Alkylation Problem

PSMs					
Containing PSMs no Var Mods (#)	Containing PSMs Containing m (#)	PSMs Containing c (#)	PSMs Containing p (#)	PSMs Containing C (#)	PSM's Cys/Total
5101	285	133	918	1527	0.213
3857	219	61	1051	778	0.140
4147	225	77	1107	866	0.145
5852	309	129	1178	1580	0.193
5204	522	223	1247	1855	0.228
4785	470	185	1072	1475	0.204
4161	287	157	1212	1170	0.184
6395	424	26	1371	533	0.058
7617	1199	303	1752	2431	0.206
7408	1078	259	1469	2084	0.188
6685	764	212	1524	1613	0.162
14035	946	113	1770	2770	0.152
1184	200	2	1144	77	0.029
1928	258	3	1004	91	0.027
4483	580	117	1381	1495	0.205
3684	447	5	1323	177	0.030
4017	434	22	1571	1061	0.166
5540	578	28	1434	1273	0.157
12714	1058	234	2129	3127	0.180
11774	936	70	1934	2376	0.151

Reduction/alkylation repeated due to low #'s of PSM's Containing Cys



SM quality metric table values, coloring & plotting done in Excel

# Future Directions

- Process Report
  - ✓ Export GCT, vertical reporter-sample template
  - Export of reporter ion intensity, de-multiplex precursor intensity.
- Quality Metrics
  - **More graphic display of metrics**
  - Report FDR at VM-site level
  - Integrate Retention Time prediction
- Autovalidation
  - Integrate Percolator/ML for more sensitive use of SM scores
- ✓ Personalized sequence databases – routine use for SAAV's and spliceforms
  - ✓ Generalize features that are hard-coded for CPTAC
  - ✓ Automate conversion of .fa/.maf for neoantigen candidates
  - Support SAAV's that are common in the population
    - minor allele frequency > 0.5%
    - ICGC-TCGA pancancer shared mutations
- 2-10x faster searches for long peptides with many variable mods.
  - faster fragment ion pre-matching to more rapidly exclude poor matching sequences.
- Fully integrate Sherenga de novo MS/MS interpretation
  - improve accuracy of HLA peptide antigen identification
  - optimize FDR in large-scale projects.
- Spectrum Matcher - propagate identifications of related spectra
- Protein Quantitation
  - dynamic allocation of shared peptides, beyond subgroup-specific, subgroup-top
  - automate identification of isoform regulation (exon splicing, pre-propeptide processing)
  - store protein grouping for faster report re-generation
- Elastic processing capacity
  - Try Google Cloud-based VMs, evaluate access to and I/O time for extraction from Google Cloud Storage bucket.
    - Create SM server image for virtual-distribution
  - Dockerize key modules of SM to enable Terra integration and more facile processing of CPTAC data from multiple PCC's.

# Load Instrument Gauges

## Spectrum Mill - Quality Metrics Dashboard

[Spectrum Mill](#) [Quality Metrics](#) [Help - slides covering overall design](#)

The following instruments are fully populated with metrics from start of Q4-2020 to present: Franklin, Galileo, Hubble2, Kati, Lacks, Beaker, Daly, Vera, Vera\_10ng

Choose an option: [Harvest Instrument QM Logs](#) **[Load Instrument Gauges](#)** [Graph Metrics](#)

### Instrument Jurkat QC Status

Gauges show **latest total peptide count** with respect to historical value ranges of **2nd quartile** and **3rd quartile**.



[Bullet chart above produced by Plotly JavaScript library](#)

[Consider switching the underlying code to D3 instead of Plotly](#)

# Graph Metrics

## Spectrum Mill - Quality Metrics Dashboard

[Spectrum Mill](#) [Quality Metrics](#) [Help - slides covering overall design](#)

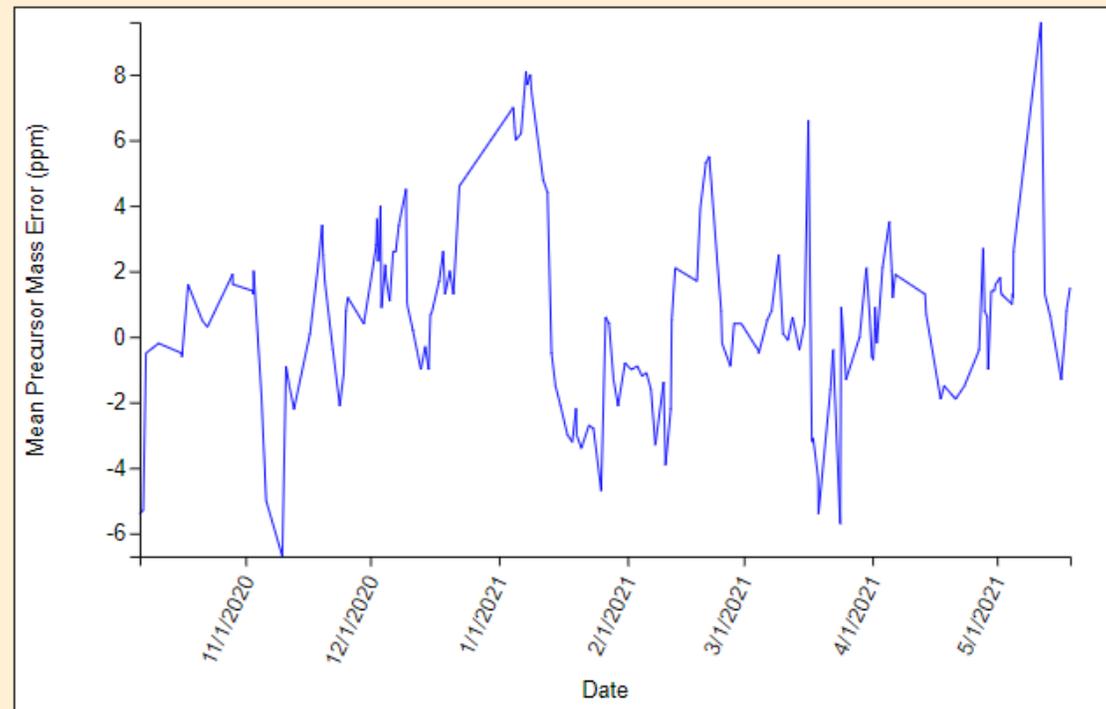
The following instruments are fully populated with metrics from start of Q4-2020 to present: Franklin, Galileo, Hubble2, Kati, Lacks, Beaker, Daly

Choose an option: [Harvest Instrument QM Logs](#) [Load Instrument Gauges](#) **[Graph Metrics](#)**

### Line Charts

Click and drag to zoom. Double click to reset zoom.

Instrument:  Quality Metric:



[Line chart produced by D3 JavaScript library using Support Vector Graphics \(SVG\)](#)

# Report to Plots - prototype

## Spectrum Mill - Report To Plots - HLA\PSM\_z1-3\_len8-11

Spectrum Mill | Extractor | MS/MS Search | Autovalidation | Quality Metrics & FDR | Protein/Peptide Summary | Workflows | Tool Belt | Protigy | Help

### Report to Plots

**Make Plots** | Save As... | Load...

### Data Directories

Select ...  Karl/PatientsIP\_snvIndel\_nuORF\_v3/Mel02\_13240-002\_20190419\_10IPs

### Peptide Filtering

Min Precursor Charge: 1 ▾ Max Precursor Charge: 3 ▾  
Min Sequence Length: 8 ▾ Max Sequence Length: 11 ▾

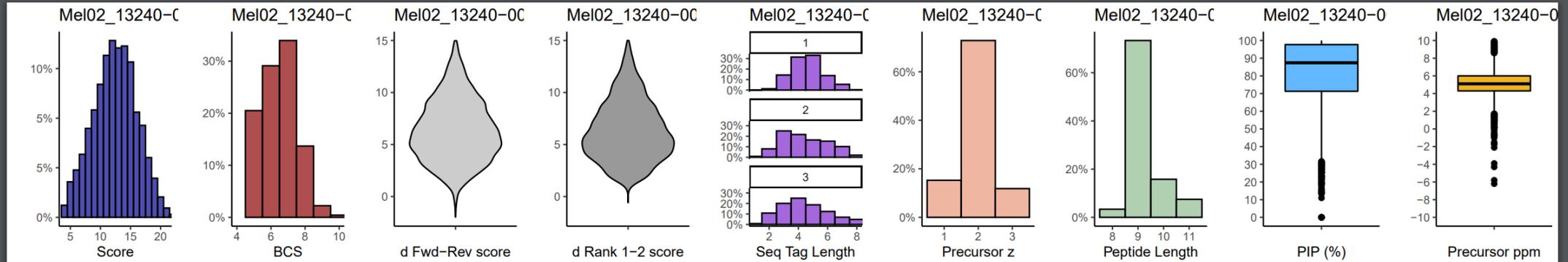
### Plot Parameters

SM report type: PSM ▾  
 Plot Score Metrics

### Results

#### R Graphics Output

1 / 1



# Similar Depth of Proteins & Phosphosites per TMT plex & channel

## LSCC more uniform loading per TMT channel

LSCC

LUAD

Proteins  
Median # / plex

11,612

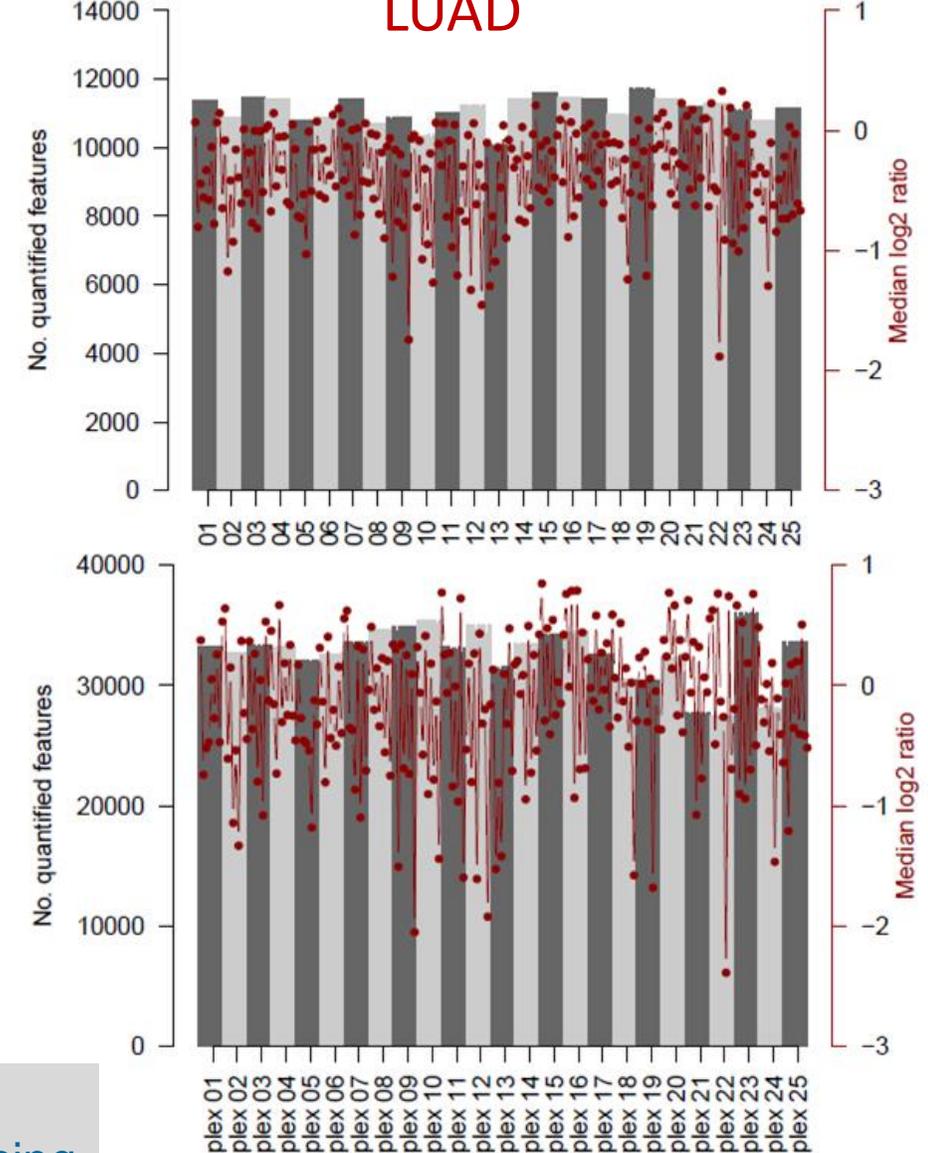
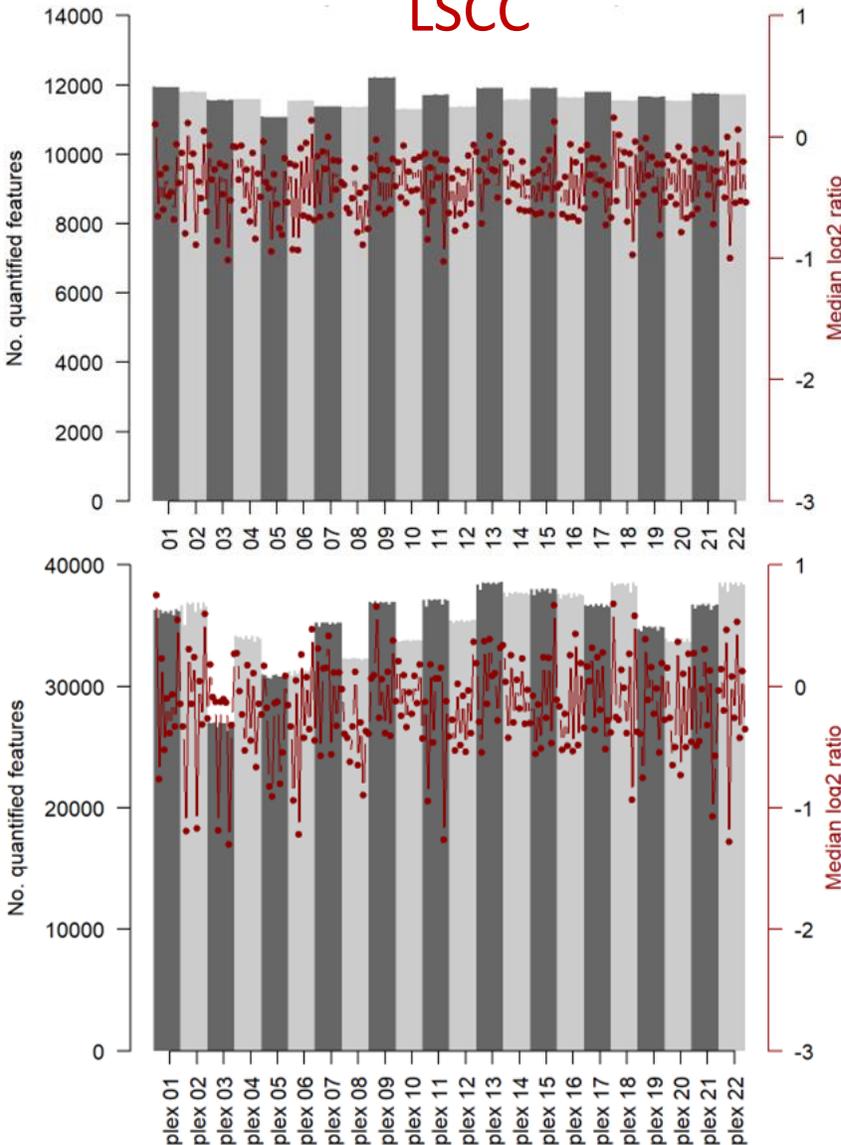
11,180

Phospho  
Sites  
Median # / plex

36,282

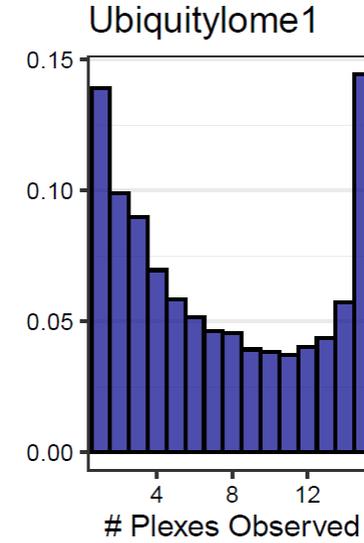
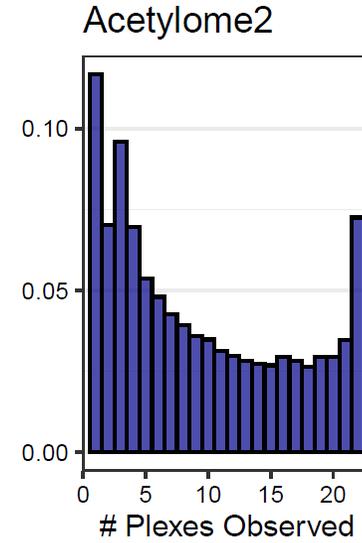
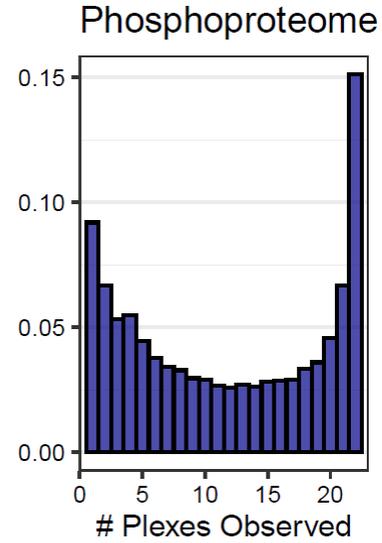
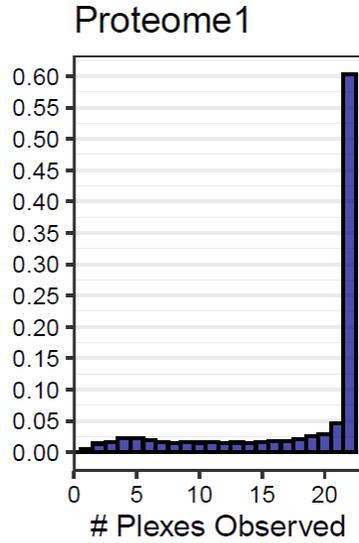
33,245

Generalize  
hardcoded CPTAC directory parsing



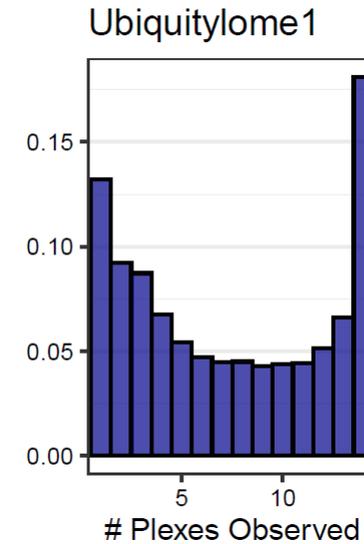
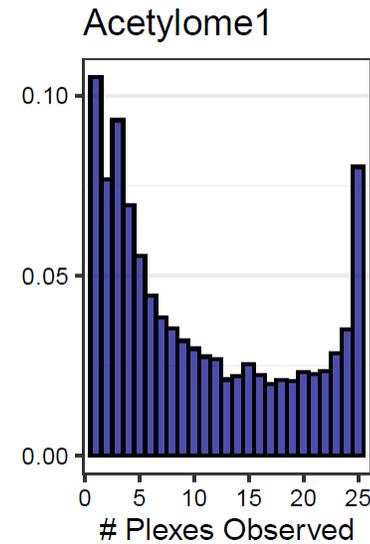
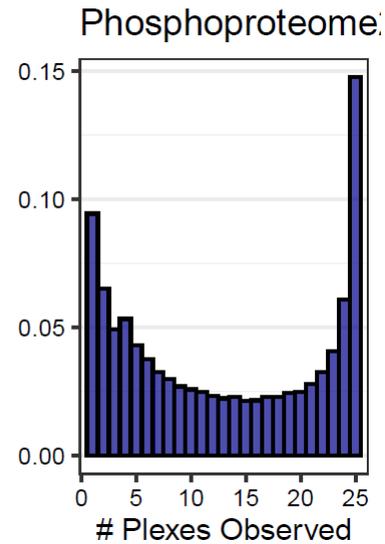
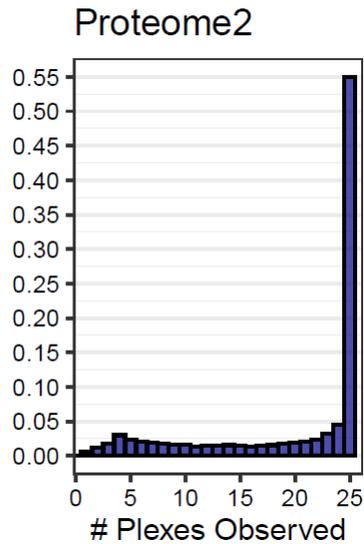
# Missing Value Histograms - prototype

LSCC



LSCC

LUAD



LUAD



# Acknowledgments



Vlado Dancik  
Broad Institute



Joe Roark  
Agilent Technologies

- DR Mani
- Karsten Krug
- Jake Jaffe
- Namrata Udeshi