# Data Extraction and Analysis for LC-MS Based Proteomics

Instructors
**Jake Jaffe[1], Deep Jaitly[2], and Matt Monroe[2]**

Co-Organizers
**Josh Adkins[2] and Gordon Anderson[2]**

[1] The Broad Institute, Cambridge, MA 02142
[2] Pacific Northwest National Laboratory, Richland, WA 99354

# Course Outline

- Introduction (Adkins)
  - Goals
  - Our Historical Perspective
  - Why Use an LC-MS Approach
  - Data and Tools Availability
- Part I: Overview of Label-Free Quantitative Proteomics (Jaffe)
- Part II: Feature discovery in LC-MS datasets (Monroe and Jaitly)
- Part III: PEPPeR, GenePattern and Real-world examples (Jaffe)
- Break
- AMT tag Pipeline Demo (general)
- Panel Discussion
  - Questions
  - Future Directions

# Course Goals

- Understand the reasons for developing and applying an LC-MS-based approach to proteomics

- Discuss considerations of experimental design for larger scale experiments

- Develop a sense of the source of information, its relative complexity and the algorithms required to make use of this approach

- See (and participate) in a demonstration of the critical tools applied to "real" data

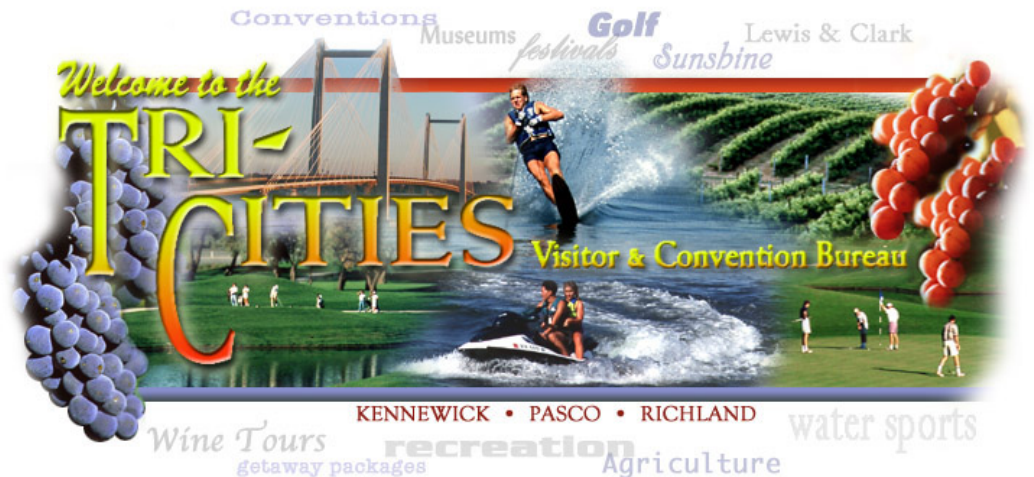- Learn where to get more information

# Pacific Northwest National Laboratory
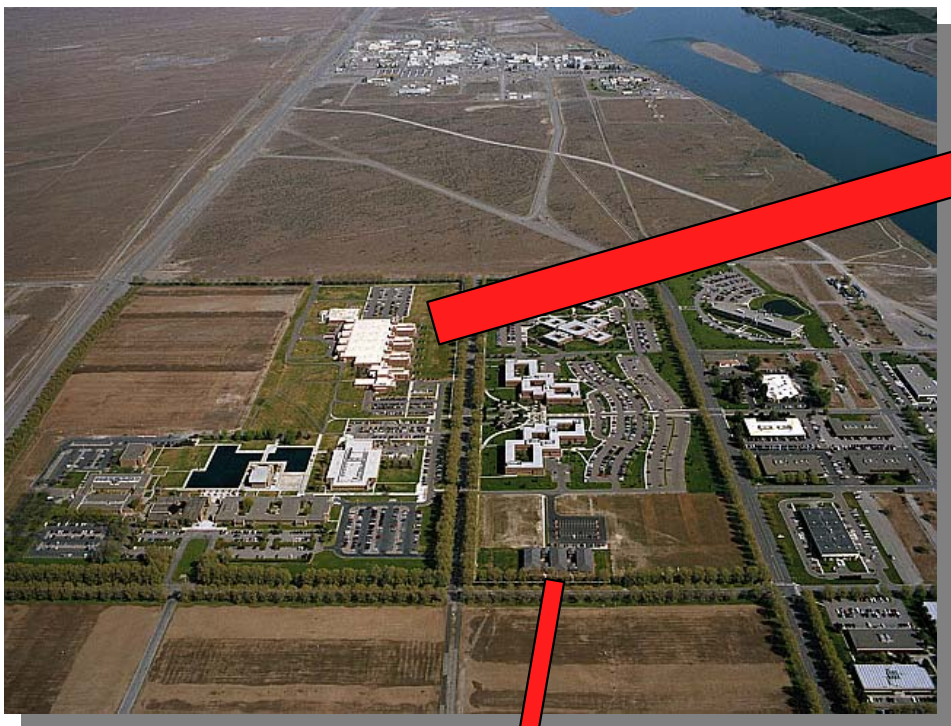


Environmental Molecular Sciences Laboratory

Washington Wine Country

# Pacific Northwest National Laboratory and EMSL

**PNNL** performs basic and applied research to deliver energy, environmental, and national security solutions for our nation.





**W.R. Wiley Environmental Molecular Sciences Laboratory**



**The Guest House at PNNL for EMSL Users**

## EMSL Mission

The W.R. Wiley Environmental Molecular Sciences Laboratory (EMSL), **a national scientific user facility** at Pacific Northwest National Laboratory, provides integrated experimental and computational resources for discovery and technological innovation in the environmental molecular sciences to support the needs of DOE and the nation.

To find out more and request access to the resource: **www.emsl.pnl.gov**

# BROAD INSTITUTE

*"Realizing the promise of the genome project for human health"*

A collaboration among MIT, Harvard, and affiliated teaching hospitals

## Programs

- Cellular Circuits
- Medical Genetics
- Chemical Biology
- Cancer Research

## Initiatives

- Metabolic Disease
- Infectious Disease
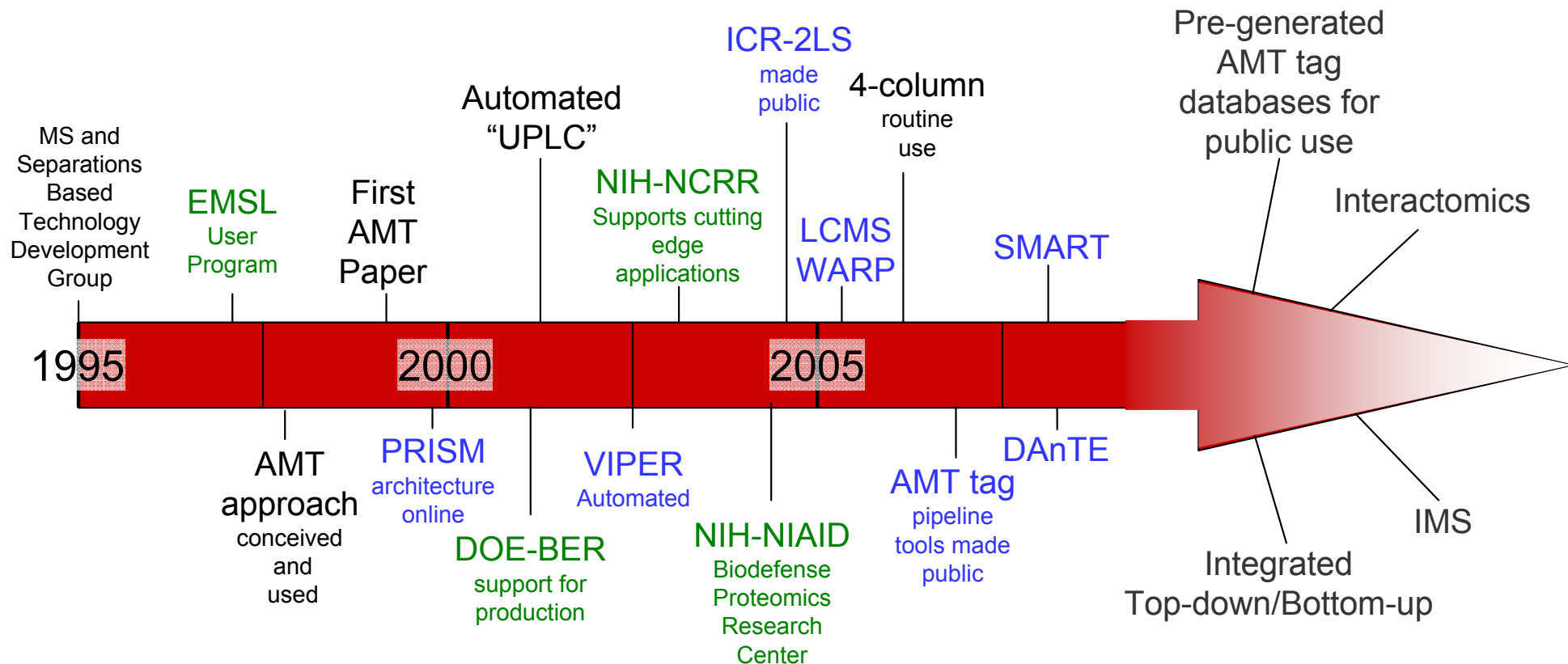- Psychiatric Disease
- Inflammatory Disease

## Platforms

- Sequencing
- Genotyping
- Chemical Synthesis and Screening
- Proteomics and Metabolite Profiling
- Image Analysis

- **Scientific mission:** Create **comprehensive, broadly available tools** for genomic medicine; pioneer **applications** toward disease understanding and treatment

- **Organizational mission:** Enable **collaborative projects** not readily done in individual labs; empower scientists through access to tools and approaches
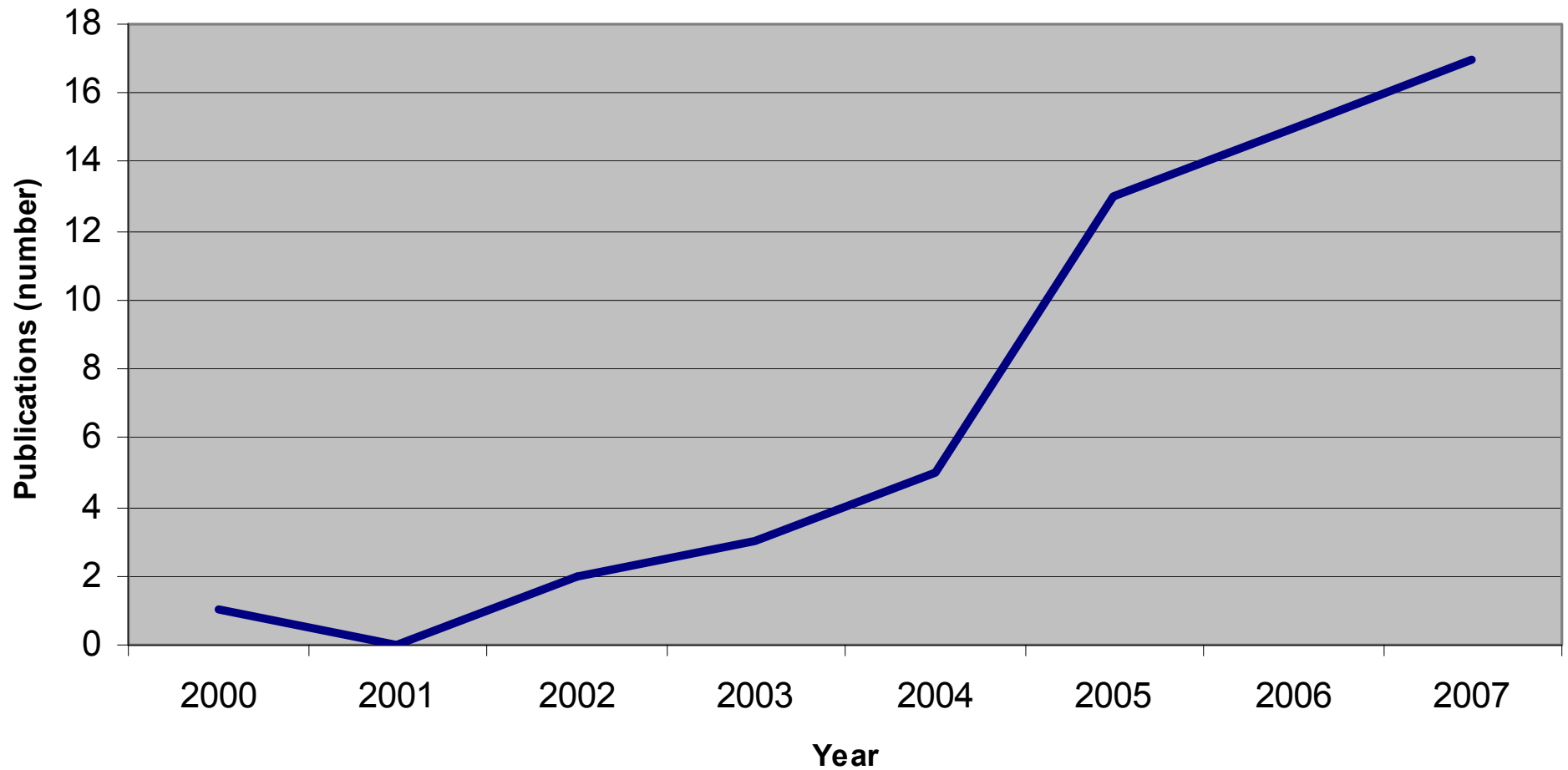
# History/Evolution of PNNL Proteomics



Key point: early access and experience with higher resolution LC and MS with ~1 ppm mass accuracy

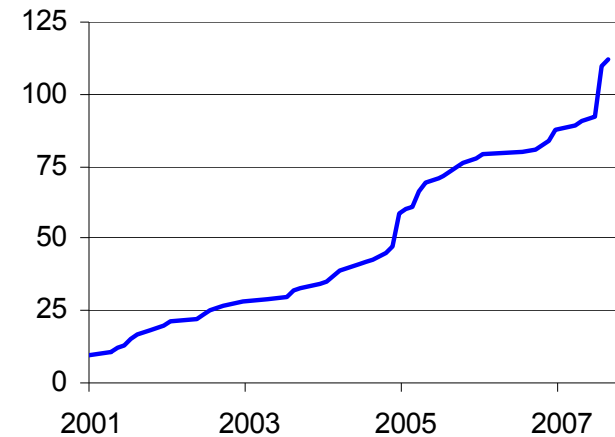# Peer-Reviewed Applications, Reviews, and Software Specific to the AMT tag Approach



Note: excludes non-AMT tag applications papers and excludes broader technology development papers
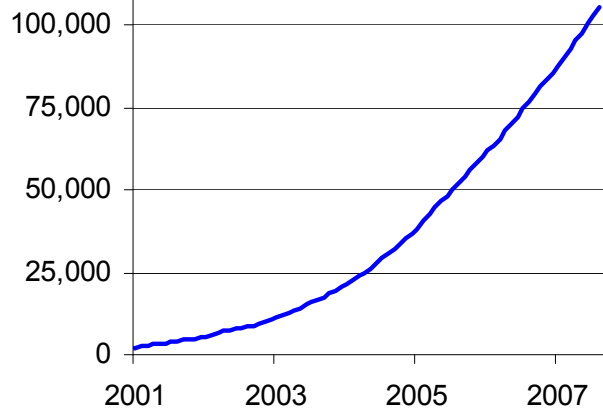
# PRISM Data Trends

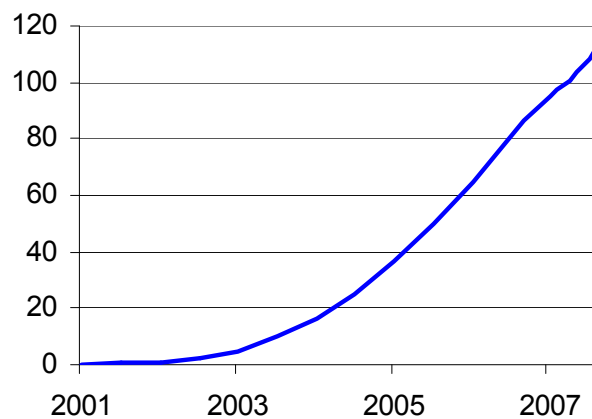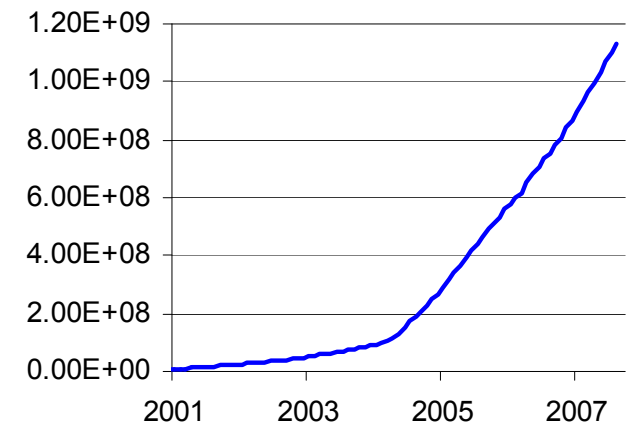| Organisms | 115 |
|---|---|
| Prepared Samples | >50,000 |
| LC-MS(/MS) Analyses | >105,000 |
| Automated Software Analyses | >277,000 |
| Data Files | 115 TB |
| Data in SQL Server databases | 1 TB |

**Organisms studied**

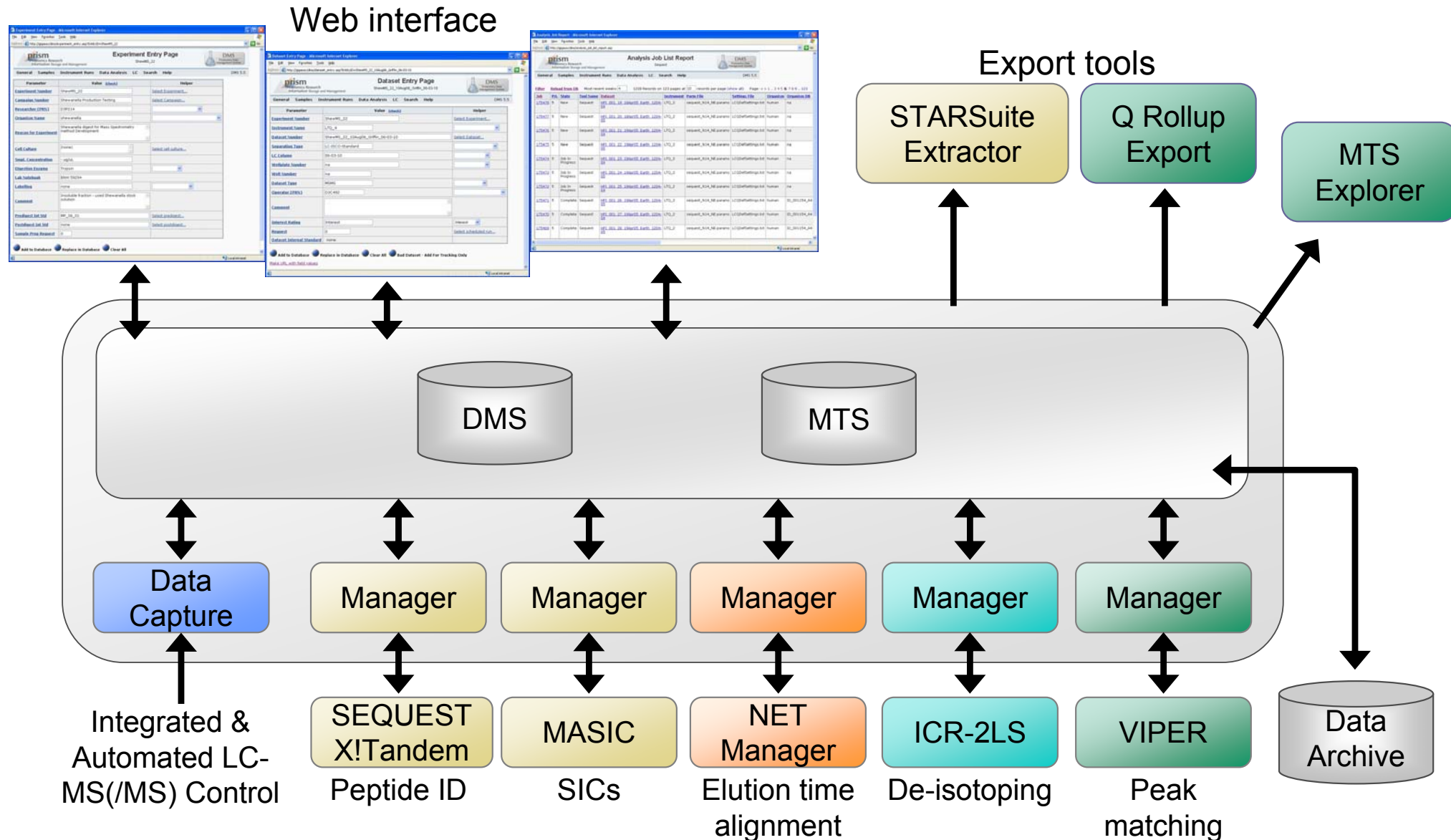**Datasets acquired (instrumental analyses)**

**TB data stored in PRISM**

**Over 1 billion mass spectra acquired**

# Proteomics Informatics Architecture
## modular and loosely coupled for flexibility



Web interface

Export tools

STARSuite Extractor

Q Rollup Export

MTS Explorer

DMS

MTS

Data Capture

Manager

Manager

Manager

Manager

Manager

Integrated & Automated LC-MS(/MS) Control

SEQUEST X!Tandem
Peptide ID

MASIC
SICs

NET Manager
Elution time alignment

ICR-2LS
De-isotoping

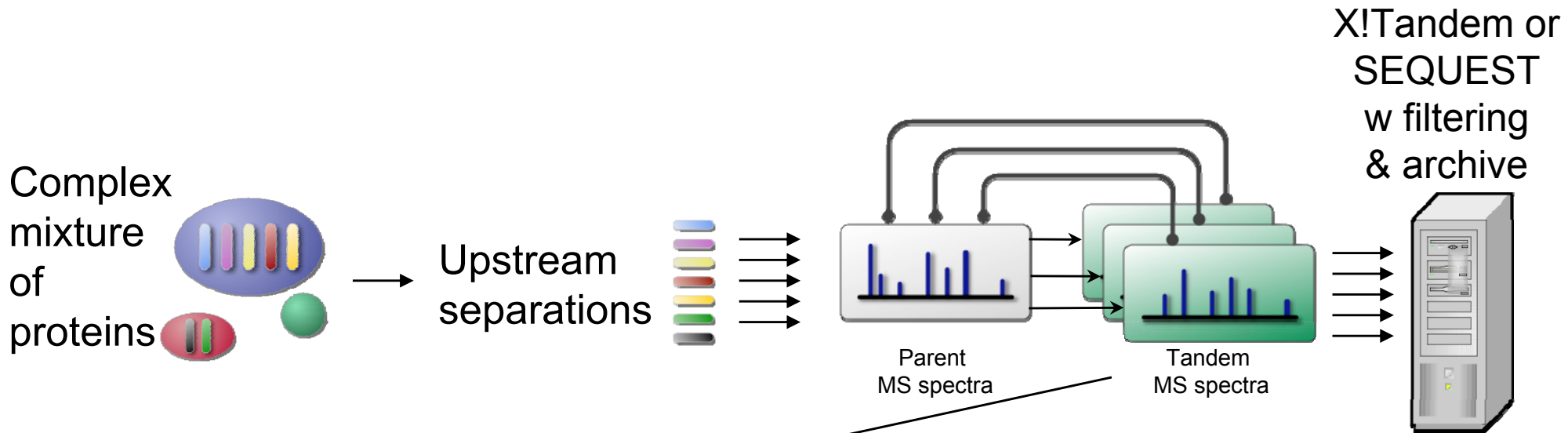VIPER
Peak matching

Data Archive

PRISM: G.R. Kiebel et. al. *Proteomics* **2006**, *6*, 1783-1790.
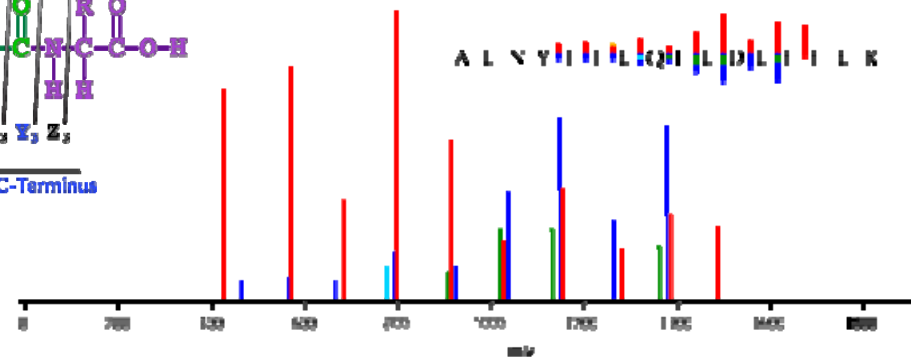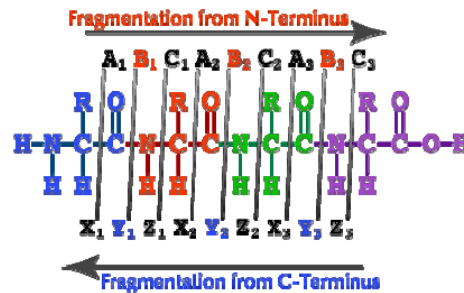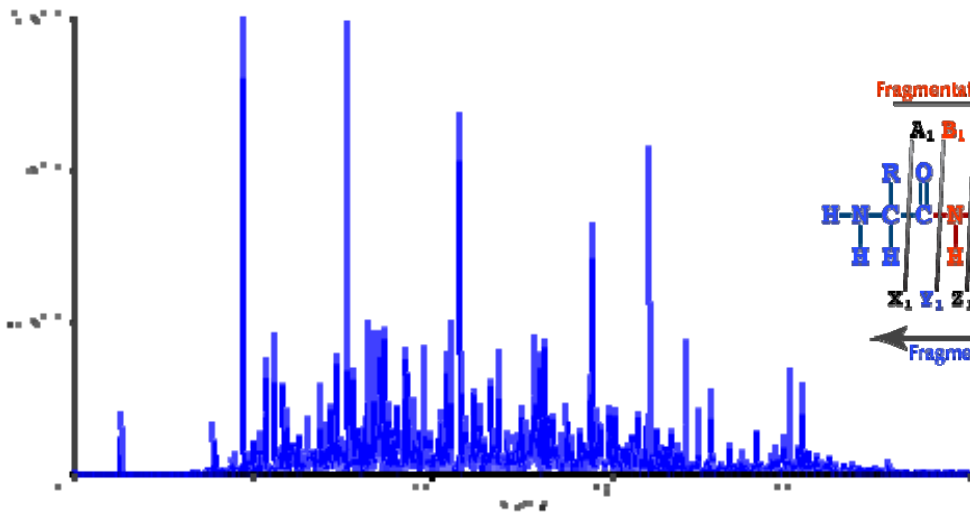
# Motivations for LC-MS Based Proteomics

- Throughput, sensitivity, and sampling efficiency
  - Compared to LC-MS/MS based approaches
- Shortcomings with chemical/labeling methods
  - Multiple species need to be sampled for each "peptide"
  - Potentially more sample preparation steps or increased cost
  - Multiple analyses still required for statistical assessment

- New challenges for experimental design
  - Blocking and randomization needs
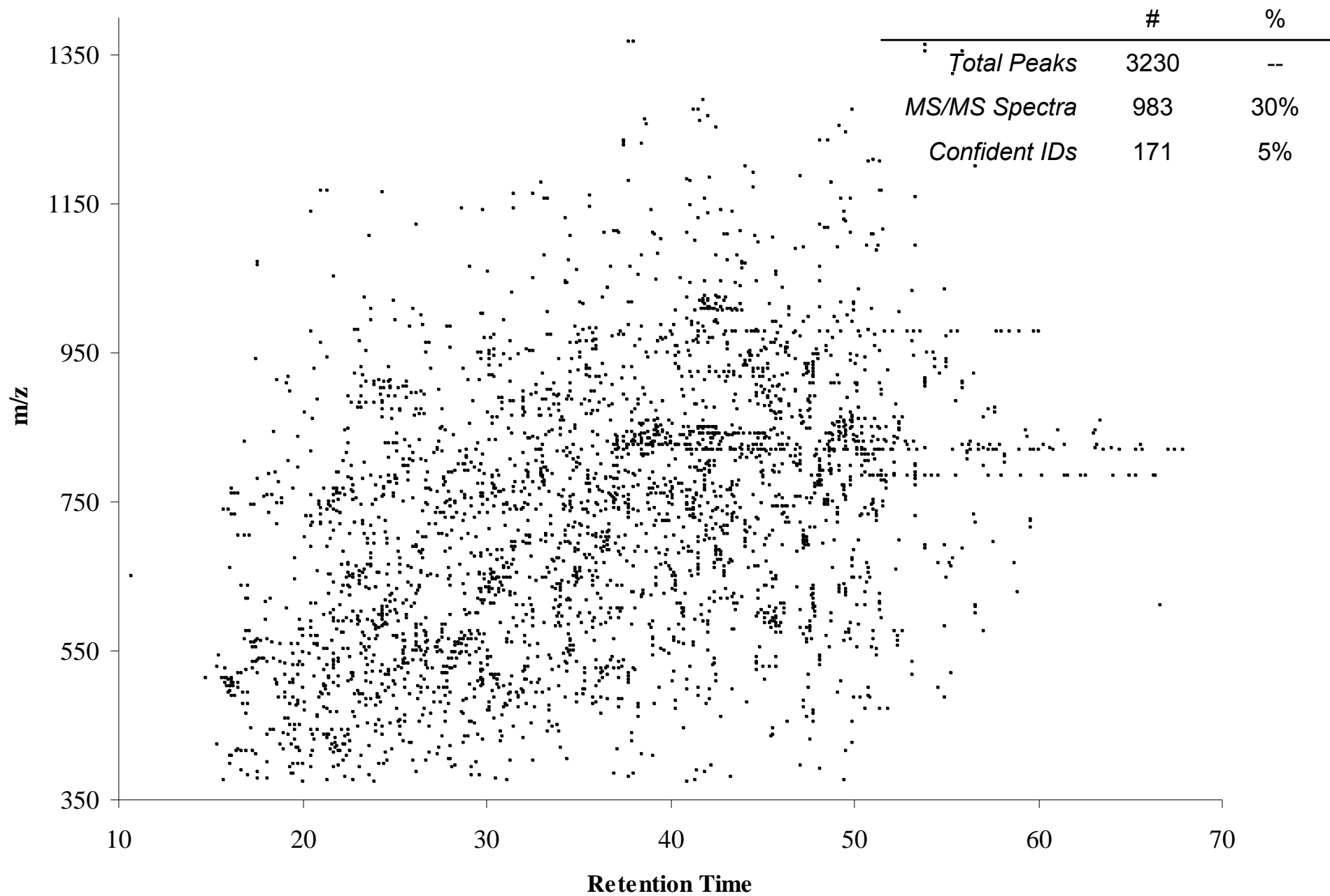
# Shotgun or MuDPIT Proteomics

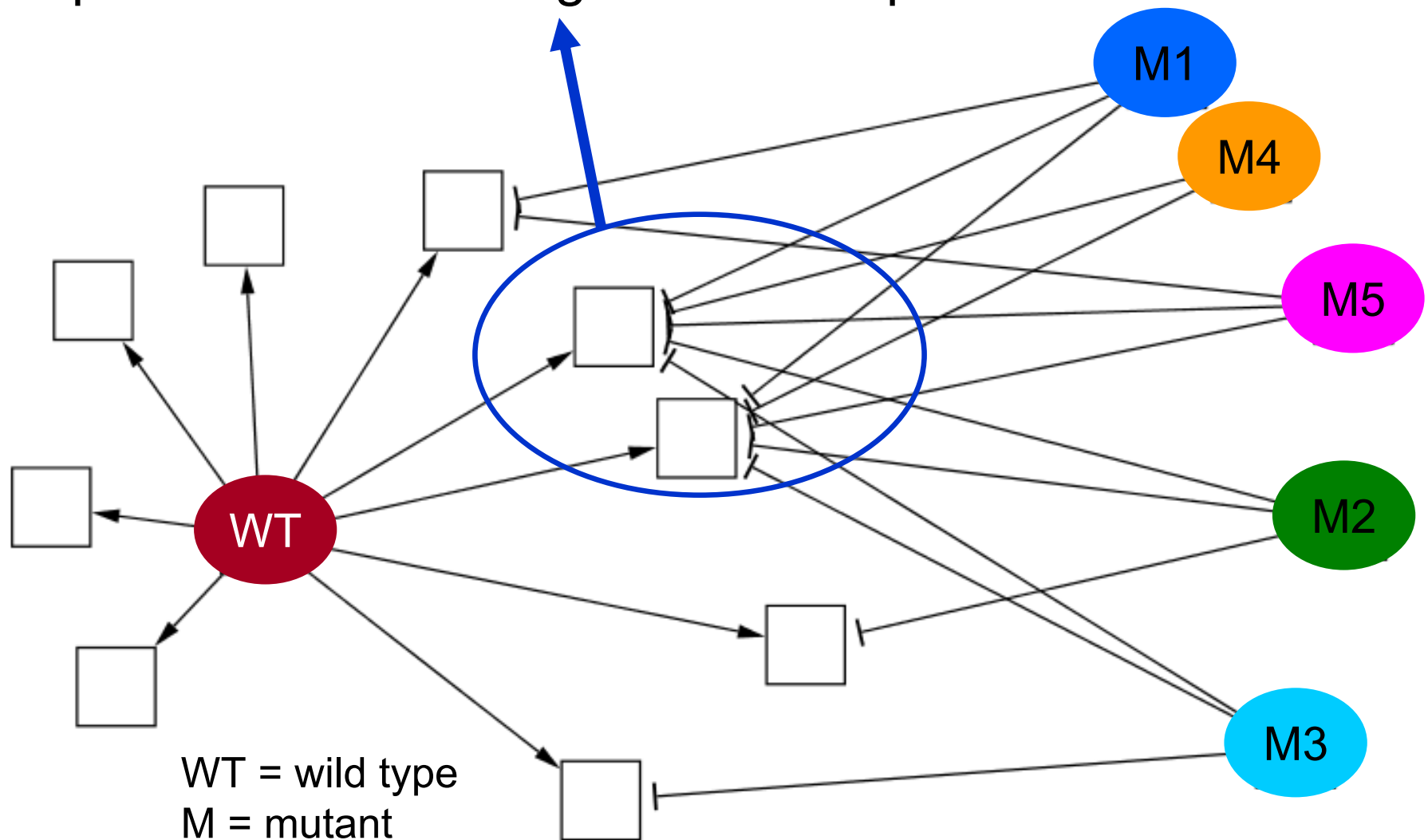# LCMS Information Gauntlet



| | # | % |
|---|---|---|
| Total Peaks | 3230 | -- |
| MS/MS Spectra | 983 | 30% |
| Confident IDs | 171 | 5% |

**m/z**

**Retention Time**

Courtesy Jake Jaffe
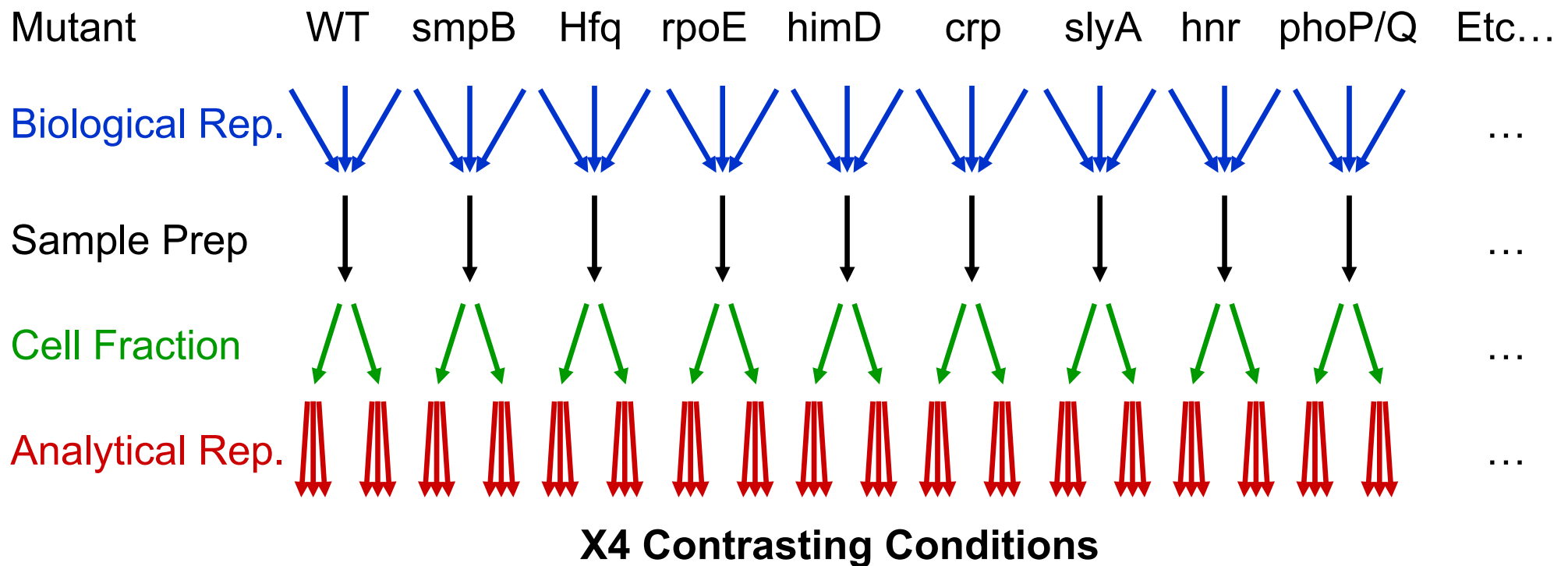
# An example need for increased throughput
# Analysis of Regulatory Mutants

Hypothesis: Knock-out regulatory proteins involved in pathogenesis and the commonly regulated proteins represent the best targets for therapeutics
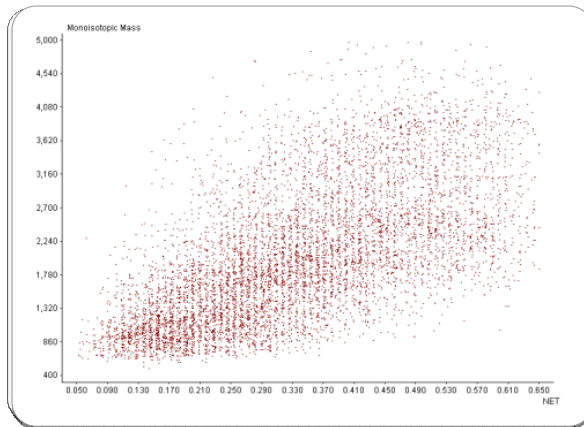


WT = wild type
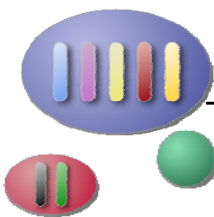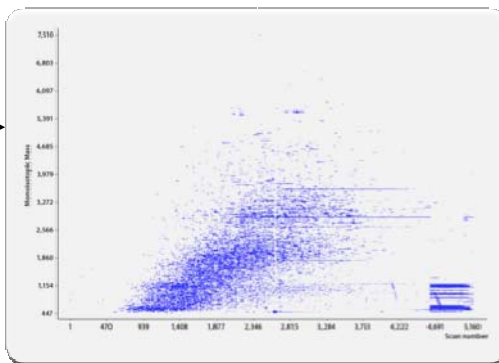M = mutant

# Accurate Mass and Time Tag Approach



SEQUEST and/or X!Tandem results
- Filtering
- Calculate exact mass
- Normalize observed rlution time

## High-throughput LC-FTICR-MS Analysis (AMT) tag



**Complex samples**

μLC- FTICR-MS

Peak-Matched Results

Compare abundances across samples

Example: V.A. Petyuk, et al. *Genome Research*. **2007**, *17* (3), 328-336.

# PEPPeR Pipeline

# New Concerns with Larger Comparisons

- Column effects (PNNL operates 4 column systems)
  - Elution time variability, potential for carryover, and stationary phase life span
- Electrospray emitters
  - Alignment, wear, clogging, etc.
- Mass Spectrometer
  - Calibration, detector response, tuning, etc.
- Samples
  - Oxidation, degradation, and other chemical modifications

# Accurate Mass and Time (AMT) Tag Data Processing Pipeline



J.S. Zimmer et. al. *Mass. Spectrom. Rev.* **2006**, *25* (3), 450-482.

# Recent Examples of Successful Applications using LC-MS Proteomics Approaches

- NIAID: *Salmonella* infecting host cells; small sample quantities → whole proteome coverage

    J.N. Adkins, et. al. *Mol. Cell. Proteomics*. **2006**, *5* (8), 1450-1461.

- Analysis of purified viral particles of Monkeypox and Vaccinia viruses

    N.P. Manes, et. al. *J. Proteome Res.* **2008**, *7* (3), 960-968.

- Analysis of "Voxels" from mouse brains to reveal protein abundance patterns in brain structures

    V.A. Petyuk, et al. *Genome Research*. **2007**, *17* (3), 328-336.

- Jake Jaffe will expand on a couple of examples such as primary tissue example; quantities too small for labeling

# Course Related Software & Data

**AMT tag Pipeline Software**



http://ncrr.pnl.gov

**PEPPeR, software within GenePattern**



http://www.broad.mit.edu/cancer/software/genepattern/

**PNNL's LCMS-based data repository**



http://omics.pnl.gov

Currently in open beta-testing
>1 Terabyte available
More coming soon!

***Salmonella typhimurium* data resource**



http://www.proteomicsresource.org

# Other Software Resources

- http://www.ms-utils.org/  (Magnus Palmblad)

- http://open-ms.sourceforge.net/index.php  (European consortium)

- http://tools.proteomecenter.org/SpecArray.php  (ISB)

- http://fiehnlab.ucdavis.edu/staff/kind/Metabolomics/Peak_Alignment/
  (Tobias Kind with Oliver Fiehn)

- http://www.proteomecommons.org/tools.jsp
  (Phil Andrews and Jayson Falkner)

# Example Data for the AMT tag Pipeline Demo

- *Salmonella typhimurium*, LC-MS/MS
  - Grown in LB (Luria-Bertani) up to log phase
  - Soluble portion of cell lysis
  - "Mini-AMT tag" database, composed of 25 SCX fractions analyzed by LC-MS/MS
  - Mass and time tag database composed from searches using X!Tandem (Log E_Value ≤ -2)
  - Linear alignment of datasets for AMT tag database
- LC-MS
  - Different sample, grown and prepared in the same conditions
- LC-FTICR-MS analysis (11T FTICR)
  - Non-linear alignment and peak matching to the database
- DAnTE data
  - Similar experiment with new growth condition

# Course Outline

- Introduction (Adkins)
- Part I: Overview of Label-Free Quantitative Proteomics (Jaffe)
  - When and why to use label free quantitative proteomics
  - Overview of the generic 'label free' pattern-based approach with guidelines
  - Discussion of alternate pipelines
- Part II: Feature discovery in LC-MS datasets (Monroe and Jaitly)
- Part III: PEPPeR, GenePattern and Real-world examples (Jaffe)
- Break
- AMT tag Pipeline Demo (general)
- Panel Discussion
  - Questions
  - Future Directions

# Part I: An Overview of Label-free quantitative proteomics
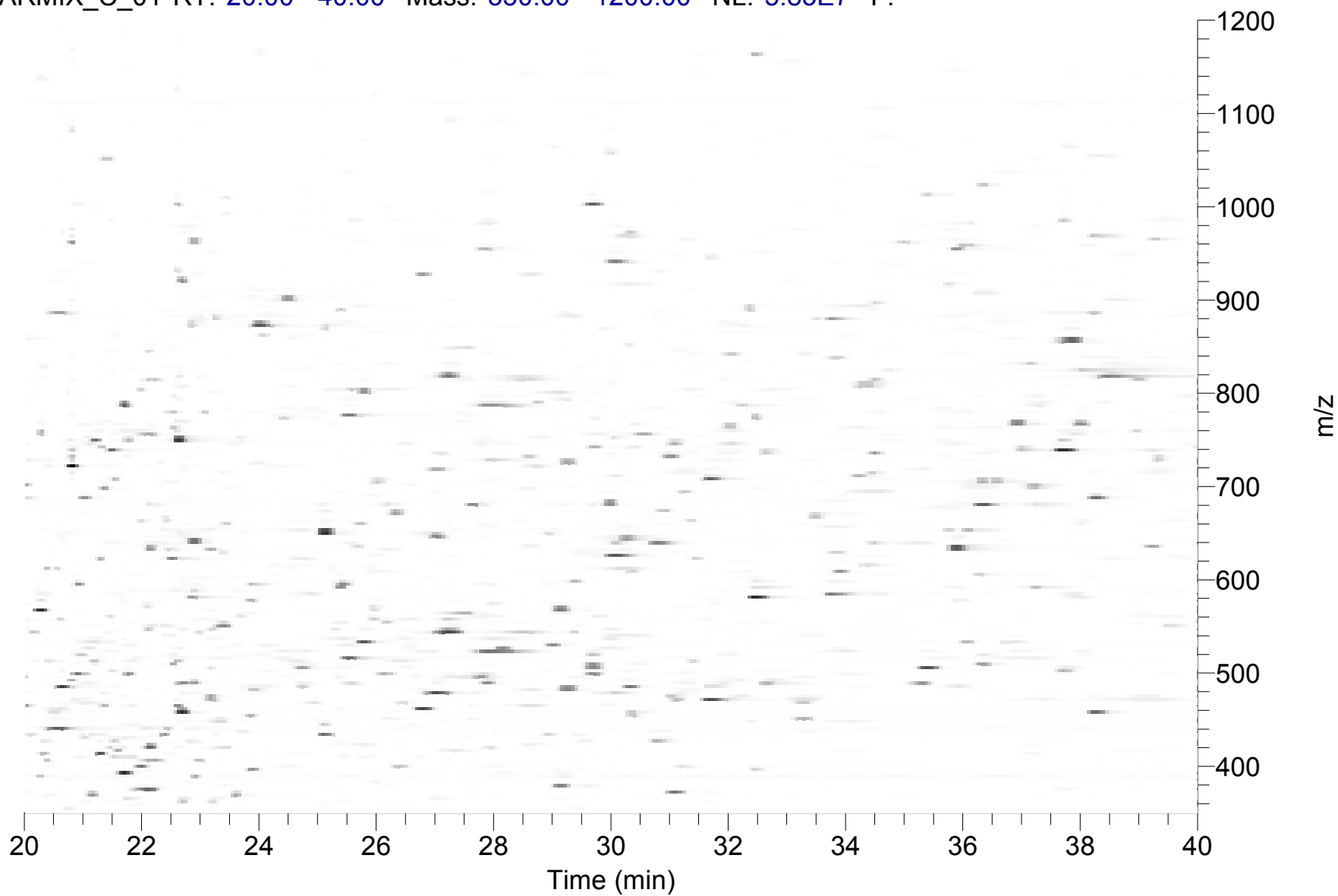
**Jacob D. Jaffe**

**The Broad Institute of Harvard and MIT**

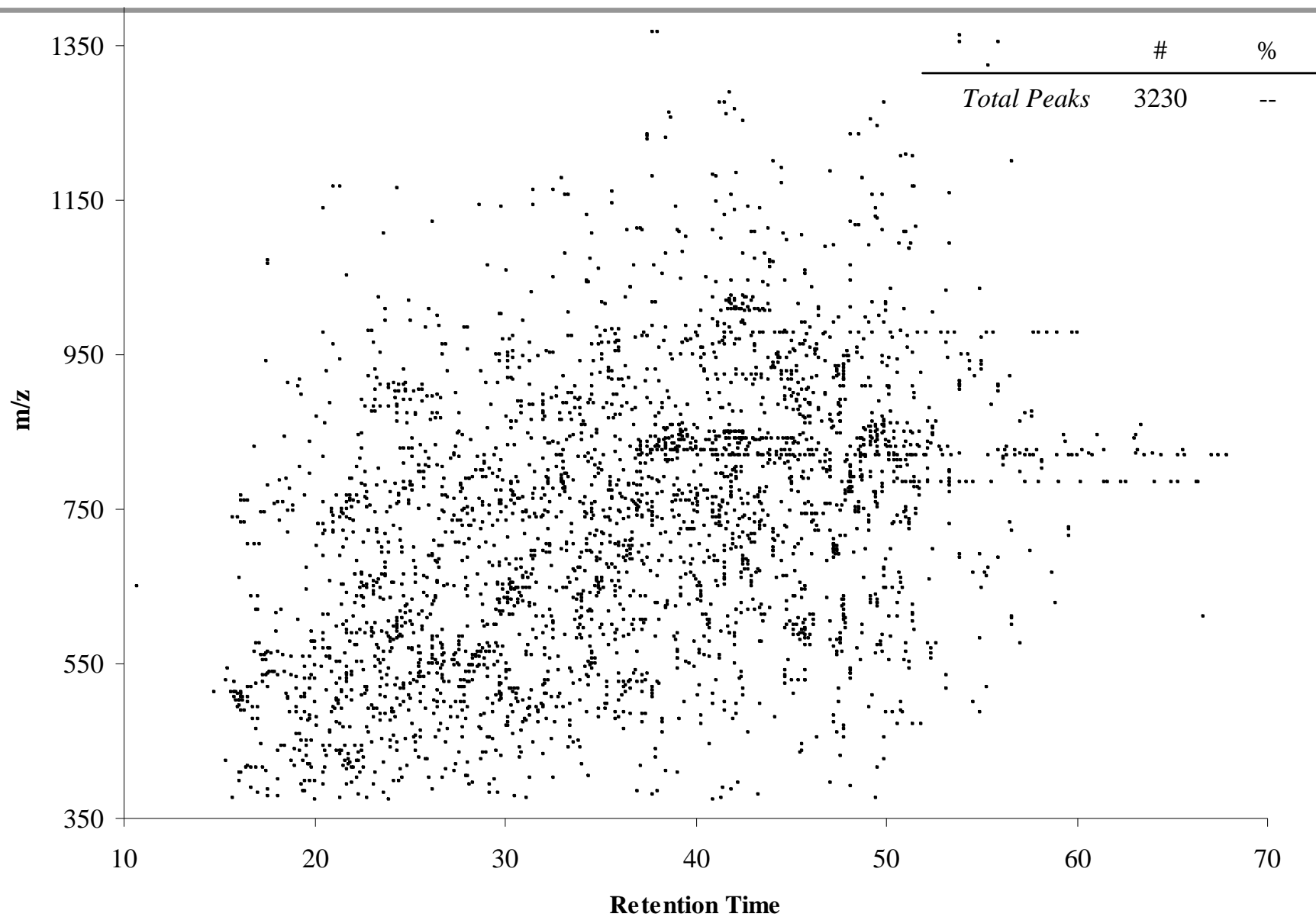**Proteomics Platform**

# Section Outline

- When and why to use label free quantitative proteomics

- Overview of the generic 'label free' pattern-based approach with guidelines

# A picture is worth 1000 parameters…

# LCMS Information Funnel – Total Peaks



| | # | % |
|---|---|---|
| *Total Peaks* | 3230 | -- |

# LCMS Information Funnel – MS/MS Sampling



| .          | #    | %   |
|------------|------|-----|
| Total Peaks | 3230 | --  |
| MS/MS Done | 983  | 30% |

# LCMS Information Funnel – MS/MS Identified



| | # | % |
|---|---|---|
| *Total Peaks* | 3230 | -- |
| *MS/MS Done* | 983 | 30% |
| *Confident IDs* | 171 | 5% |

# Definition of **<u>Label-free Quantitative Proteomics</u>**

- Use of raw mass spectral signal intensity (peaks) as a surrogate for the abundance of a peptide and/or protein

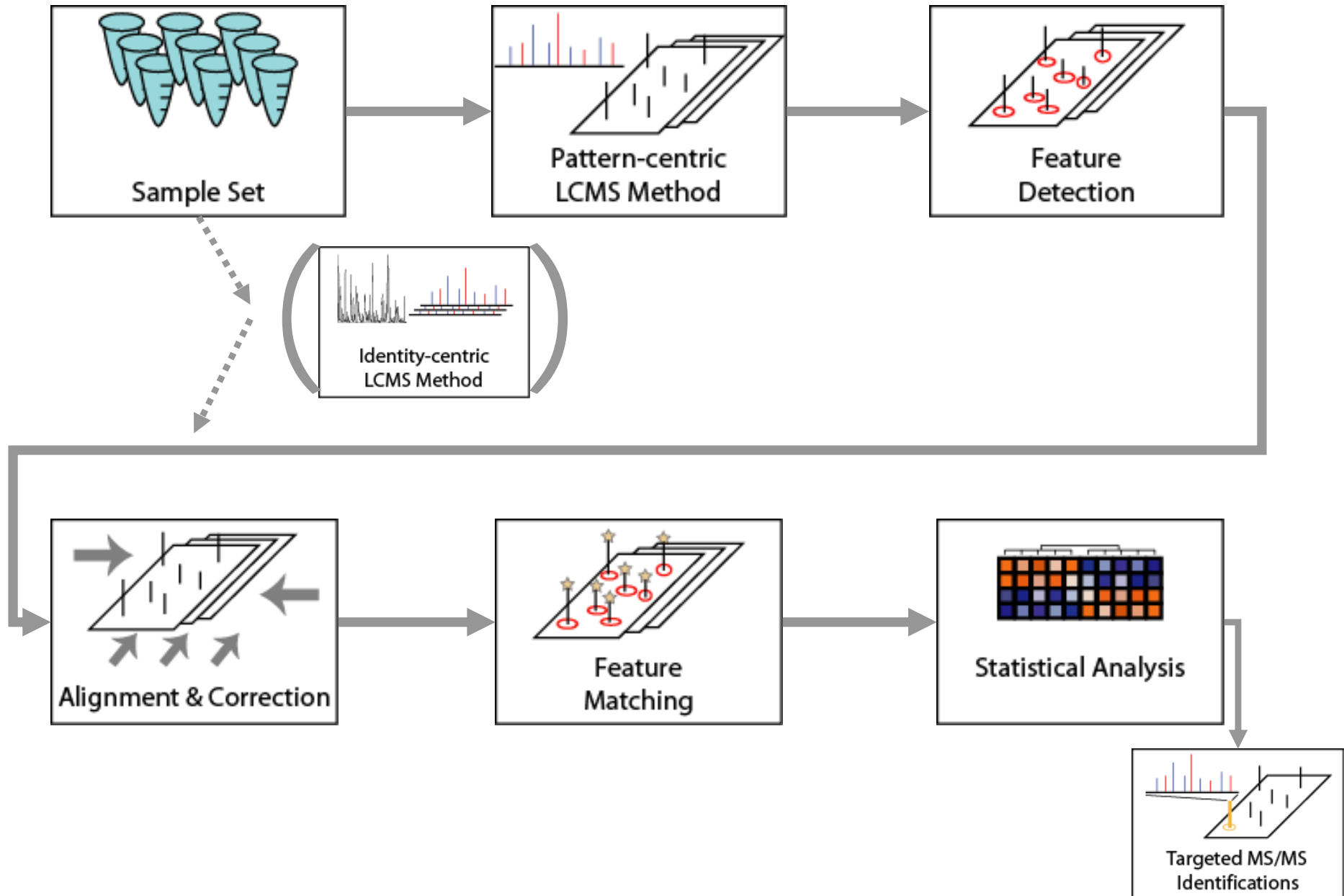- Signal intensity from the same analyte is compared across multiple experimental conditions as the basis for quantitation

- When coupled to LC, peaks have dimensions in retention time and well as *m/z* and intensity

- Careful experimental processing and computational methods are required to extract quantitative information in label-free proteomics

# Motivations for Label-free Quantitative Proteomics

- **Microarray envy**
  - Well-defined experiment, well-defined tools

- **Differential detection and _quantification_ of proteins**
  - Biomarker discovery and pattern recognition
  - Biological insight into the real actors in the cell: proteins
  - Time course analysis

- **$MS^2$ independent but friendly**
  - SILAC and iTRAQ (labeling methods) require $MS^2$ ID for entry
  - Comprehensive!  Quantify all the spots!  Even the faint ones!

- **Minimal sample workup**
  - Primary tissue OK
  - No artifacts from labeling efficiency

# The Generic 'Label-free' Workflow



Sample Set

Pattern-centric LCMS Method

Feature Detection

Identity-centric LCMS Method

Alignment & Correction

Feature Matching

Statistical Analysis

Targeted MS/MS Identifications

# Best Practices: Getting Started

- Team approach
  - LCMS expert experimentalist
  - Computer scientist/programmer
  - Statistician

- Planning
  - Statistical power of study (consult statistician)
  - Identification of reliable sample sources
  - Instrument / Computational / Storage infrastructure

- Execution
  - Patience
  - Consistency

# Best Practices: Samples

- BEST POSSIBLE **SAMPLES** AND **CONTROLS**
  - Relevant to disease or study target
  - Proximal to the source of differential markers
  - Consistent in composition
  - Controls appropriately matched (same subject if possible)
  - Enriched in likely differential markers
  - GARBAGE IN, GARBAGE OUT

- Sample processing pipeline TESTED and CONSISTENT
  - Abundant protein depletion (serum proteomics)?
  - Fractionation required?
  - Measure yields – are they consistent?
  - CLEAN!!!

- Collect more than you need – outlier removal!

# Best Practices: Data Acquisition

- **Resolution! Resolution! Resolution!**
  - FTICR or Orbitrap recommended > 60,000 resolution
  - More 'channels'

- **Accuracy! Accuracy! Accuracy!**
  - Calibrate mass often
  - Downstream recognition of "same" feature easier
  - Statistical confidence

- **Consistency**
  - LCMS methods and instrumentation
  - LC column and length

- **Common Sense**
  - $MS^1$ sampling rate -> chromatographic resolution
  - Tolerances and dynamic exclusion for $MS^2$ sampling
  - Carry over testing and sample randomization
  - SAVE THE SAMPLES!!!!!

# Best Practices: Feature Picking

- Understand the method
  - No method is demonstrably 'best'
  - Consult with expert help
  - All methods have parameters and tolerances that have to be tailored to your operating characteristics
  - There is no magic 'black box'

- Patience
  - You will spend a long time collecting data; expect to spend at least as much time extracting and analyzing data
  - Budget time and resources to explore parameters on a subset of your data before doing feature picking *en masse*

# Best Practices: Experiment Alignment

- Consistency in experimental execution
  - Makes life easier, less computational correction

- Pay attention to output of aligners
  - Methods may have metrics of alignment quality
  - Large corrections may signal outlier experiments
    - Consider discarding

- Intensity normalization
  - Total ion current (TIC)?
  - TIC of all features?
  - Subset of 'housekeeping' features?
  - Medians, means, etc?

# Best Practices: Feature Assignment and Matching

- **Assignment: annotation of an LCMS feature with a peptide identity (sequence)**
  - Derived from external or embedded MS$^2$ data that has been searched against a database (i.e. Sequest, Mascot, etc)
  - AMT-based assignment (importance of mass accuracy)
  - Look for statistics!

- **Matching: recognition that a feature is the same across multiple experiments irrespective of an identity assignment**
  - Assignments can help
  - Accuracy and alignment are paramount
  - Take care with user-adjustable tolerances
  - Look for statistics!

# Best Practices: Statistical Analysis

- Intensity normalization of features must be done prior to statistical analysis
  - Also address handling of missing values

- Understand what you are doing or seek assistance
  - Know your $p$-values from your $q$-values (and FDRs)

- Have a well-formulated statistical question
  - Most statistical tests are measured vs. the 'null' hypothesis
  - Decide in advance what levels of false discovery are acceptable
  - Significance level $\neq$ priority for follow-up

- There are many tools available
  - Some are more proteomics-amenable
    - Handling of intensity normalization
    - Handling of proteins as combinations of peptides

# Best Practices: Following-up

- Targeted reinterrogation of samples for identification of 'unidentified' features

- Literature mining
  - Possible connections to your biological questions
  - Helps with prioritization

- Targeted assessment of interesting features in alternative matrices
  - I.e., discovered in tissue, but is it present in blood?
  - Methods other than mass spec, too!

# Reference Chart of Label-free Platforms

| | PNNL Pipeline | PEPPeR | msInspect | SuperHirn | CRAWDAD |
|---|---|---|---|---|---|
| Lab | PNNL | Broad Institute | FHCRC | IMSB (Swiss) | Univ. Wash. |
| Feature Picker | Decon2LS/Viper | Mapquant (or any other) | msInspect | SuperHirn | CRAWDAD |
| Method | Spectrum de-isotoping then clustering | Image Analysis then de-isotoping | Wavelet decomposition then de-isotoping | Spectrum de-isotoping then merging | m/z channel binning |
| RT Alignment | Normalization, then linear or LCMSWARP | Relative, then linear, or LOESS (exp) | Iterative non-linear transformation | LOESS modeling | Dynamic time warping |
| *m/z* recalibration | Yes (dynamic) | Yes (quadratic) | No | No | No |
| Assignment of IDs to features | AMT database, normalized elution times | AMT database, relative elution order (Landmarks) | AMT database through user interaction | Yes, but not well documented at present | Yes, for differences only if they exist |
| Statistical Evaluation of assignment | Mass shift decoy and/or Bayesian Statistics | Bayesian Statistics | No | No | No |
| Unidentified Feature Recognition | Stored in database for later analysis | Data-dependent tolerance-based clustering | User specified tolerance-based clustering | Tolerance-based merging, heuristics | Difference mapping only |
| Runs on | Windows with GUI | Web-based (Linux or Windows install bases) | Java with GUI | Linux | Linux/Windows |

# Part II: LC-MS Feature Discovery

- Introduction (Adkins)
- Part I: Overview of Label-Free Quantitative Proteomics (Jaffe)
- Part II: Feature discovery in LC-MS datasets (Monroe and Jaitly)
  - Structure of LC-MS Data
  - Feature discovery in individual spectra (deisotoping)
  - Feature definition over elution time
  - Identifying LC-MS Features using an AMT tag DB
  - Extending the AMT tag approach for feature based analyses
  - Estimating confidence of identified LC-MS features
  - Downstream quantitative analysis with DAnTE
- Part III: PEPPeR, GenePattern and Real-world examples (Jaffe)
- Break
- AMT tag Pipeline Demo (general)
- Panel Discussion
  - Questions
  - Future Directions

# Part II: Feature discovery in LC-MS datasets

**Navdeep Jaitly and Matthew E. Monroe**

**Pacific Northwest National Laboratory**

# Structure of LC-MS Data

- Mass spectra capture the changing composition of peptides eluting from a chromatographic column

  - Complex peptide mixture on a column is separated by liquid chromatography over a period of time

  - Changing composition of the mobile phase causes different peptides to elute at different times

  - The components eluting from a column are sampled continuously by sequential mass spectra



QC Standards
(12 protein digest)

# Structure of LC-MS Data

- Each compound is observed as an <u>isotopic pattern</u> in a mass spectrum

  - The pattern is dependent on the compound's chemical composition, charge, and resolution of instrument

## Theoretical Profile



Peptide: *VKHPSEIVNVGDEINVK*

Parent Protein: *gi|16759851 30S ribosomal protein S1*

Charge: *2+*
m/z: *939.0203*
Monoisotopic Mass: *1876.0054 Da*

# Structure of LC-MS Data

- A mass spectrum of a complex mixture contains overlaid distributions of several different compounds



scan 1844

# Structure of LC-MS Data

- With LC as the first dimension, each compound is observed over multiple spectra, showing a three-dimensional pattern of m/z, elution time and abundance

Salmonella typhimurium dataset



Peptide: *VKHPSEIVNVGDEINVK*

Parent Protein: *gi|16759851 30S ribosomal protein S1*

Charge: 2+
m/z: *939.0203*
Monoisotopic Mass: *1876.0054 Da*

Elution range: *Scans 1539 - 1593*

# Feature Discovery in LC-MS data

- Goal: Infer *(mass, elution time, intensity)* of compounds that are present in data obtained from an LC-MS dataset
  - Compounds are termed LC-MS features since they are inferred from a three dimensional pattern, yet identity is unknown

2D view of an LC-MS analysis of Salmonella typhimurium

# Feature Discovery in LC-MS data

- Sequential process of finding features in each mass spectrum is followed by grouping of features over multiple spectra together

2D views of an LC-MS dataset in different stages of processing



raw data → deisotoping → Collapsed monoisotopic features in all spectra → Elution profile discovery → LC-MS features

# Feature discovery in individual spectra

- Deisotoping
  - Process of converting a mass spectrum (*m/z, intensity*) into a list of species (*mass, abundance, charge*)

Deisotoping a mass spectrum of 4 overlapping species



| charge | Monoisotopic MW | abundance |
|---|---|---|
| 2 | 1546.856603 | 533467 |
| 2 | 1547.705048 | 194607 |
| 2 | 1547.887682 | 671947 |
| 2 | 1548.799612 | 426939 |

# Deisotoping an Isotopic Distribution

- Decon2LS deisotoping algorithm compares theoretical isotopic patterns with observed patterns



Observed spectrum

Theoretical spectrum

Fitness value

Charge detection algorithm[2]

charge = 2

avg. mass = 1876.02

Averagine[3]

estimated empirical formula:

$C_{83} H_{124} N_{23} O_{25} S_1$

Mercury[4]

1. Horn, D.M., Zubarev, R.A., McLafferty, F.W. Automated Reduction and Interpretation of High Resolution Electrospray Mass Spectra of Large Molecules. *J. Am. Soc. Mass Spectrom.* **2000**, *11*, 320-332.

2. Senko, M. W.; Beu, S. C.; McLafferty, F. W. Automated assignment of charge states from resolved isotopic peaks for multiplycharged ions. *J. Am. Soc. Mass Spectrom.* **1995,** *6*, 52–56.

3. Senko, M. W.; Beu, S. C.; McLafferty, F. W. Determination of monoisotopic masses and ion populations for large biomolecules from resolved isotopic distributions. *J. Am. Soc. Mass Spectrom.* **1995,** *6*, 229–233.

4. Rockwood, A. L.; Van Orden, S. L.; Smith, R. D. Rapid Calculation of Isotope Distributions. *Anal. Chem.* **1995,** *67,* 2699–2704.

# Deisotoping an Isotopic Distribution

- Patterson (Autocorrelation) algorithm to detect charge of a peak in a complex spectrum

- Mercury algorithm used to guess an average empirical formula for a given mass

  - Averagine empirical formula of $C_{4.9384}$ $H_{7.7583}$ $N_{1.3577}$ $O_{1.4773}$ $S_{0.0417}$ $\rightarrow$ $C_{83}$ $H_{124}$ $N_{23}$ $O_{25}$ S for 1876.02 Da

- Fitness (fit) functions to quantitate quality of match between theoretical and observed profiles

- For additional details, see the slides presented at 2007 US HUPO, available at http://ncrr.pnl.gov/training/workshops/

# $^{16}O/^{18}O$ Mixtures

- Overlapping isotope patterns are separated by 4 Da
  - Creates challenges for deisotoping, particularly for charge states of 3+ or higher

# Isotopic Composition

- Deviation from natural abundances

  - In $^{13}$C, $^{15}$N depleted media, isotopic composition of atoms is different from those found in nature

  - E.g., sulfur isotopes predominate the distribution at right

  - Constrast with an isotopic distribution of a peptide with similar mass and charge (16+), but a natural atomic distribution (below)

# Isotopic Composition

- Decon2LS supports changing the isotope composition

# Part II: LC-MS Feature Discovery

- Introduction (Adkins)
- Part I: Overview of Label-Free Quantitative Proteomics (Jaffe)
- Part II: Feature discovery in LC-MS datasets (Monroe and Jaitly)
    - ✓ Structure of LC-MS Data
    - ✓ Feature discovery in individual spectra (deisotoping)
    - Feature definition over elution time
    - Identifying LC-MS Features using an AMT tag DB
    - Extending the AMT tag approach for feature based analyses
    - Estimating confidence of identified LC-MS features
    - Downstream quantitative analysis with DAnTE
- Part III: PEPPeR, GenePattern and Real-world examples (Jaffe)
- Break
- AMT tag Pipeline Demo (general)
- Panel Discussion

# Feature definition over elution time

- Deisotoping collapses original data into data lists

| scan num | charge | abundance | mz | fit | average mw | monoiso mw | most abu. mw | fwhm | signal noise |
|---|---|---|---|---|---|---|---|---|---|
| 1500 | 1 | 2772933 | 759.0649 | 0.0716 | 758.5222 | 758.0576 | 758.0576 | 0.0106 | 718.83 |
| 1500 | 1 | 2614913 | 1103.033 | 0.1111 | 1102.698 | 1102.026 | 1102.026 | 0.0222 | 74.04 |
| 1500 | 1 | 2422829 | 864.4919 | 0.0156 | 864.0073 | 863.4846 | 863.4846 | 0.0137 | 74.75 |
| 1500 | 2 | 2297822 | 563.3253 | 0.012 | 1125.322 | 1124.636 | 1124.636 | 0.006 | 77.94 |
| 1500 | 1 | 1213607 | 943.9815 | 0.1025 | 943.5518 | 942.9742 | 942.9742 | 0.0165 | 120.36 |
| 1500 | 3 | 988761 | 675.0246 | 0.02 | 2023.375 | 2022.052 | 2023.0549 | 0.0086 | 79.22 |
| 1500 | 2 | 734070 | 688.392 | 0.0384 | 1375.694 | 1374.77 | 1374.7695 | 0.009 | 92.09 |
| 1500 | 2 | 663954 | 642.3243 | 0.0253 | 1283.417 | 1282.634 | 1282.6341 | 0.0076 | 109.01 |
| 1500 | 1 | 661477 | 730.1117 | 0.024 | 729.5461 | 729.1045 | 729.1045 | 0.0096 | 39.06 |
| 1500 | 2 | 630657 | 689.3645 | 0.0446 | 1377.64 | 1376.715 | 1376.7145 | 0.0088 | 57.52 |
| 1500 | 2 | 569896 | 591.8343 | 0.0198 | 1182.379 | 1181.654 | 1181.6541 | 0.0065 | 111.2 |
| 1500 | 2 | 503993 | 757.8854 | 0.0706 | 1513.762 | 1512.753 | 1512.7533 | 0.0105 | 80.4 |
| 1500 | 2 | 451007 | 936.9389 | 0.0296 | 1873.091 | 1871.863 | 1872.8662 | 0.0156 | 46.74 |

- Goal: Given series of deisotoped mass spectra, group related data across elution time
  - Look for repeated monoisotopic mass values in sequential spectra, allowing for missing data
  - Can also look for expected chromatographic peak shape

# Feature definition over elution time

- Can visualize deisotoped data in two-dimensions



Monoisotopic Mass    1 ■ 2 ■ 3 ■ 4 ■ 5 ■ Other ■

- Plotting monoisotopic mass
- Color is based on charge of the original data point seen
- Monoisotopic Mass = (m/z x charge) - 1.00728 x charge

Mass

Time

Scan number

# Feature definition over elution time

- ● Zoom-in view of species
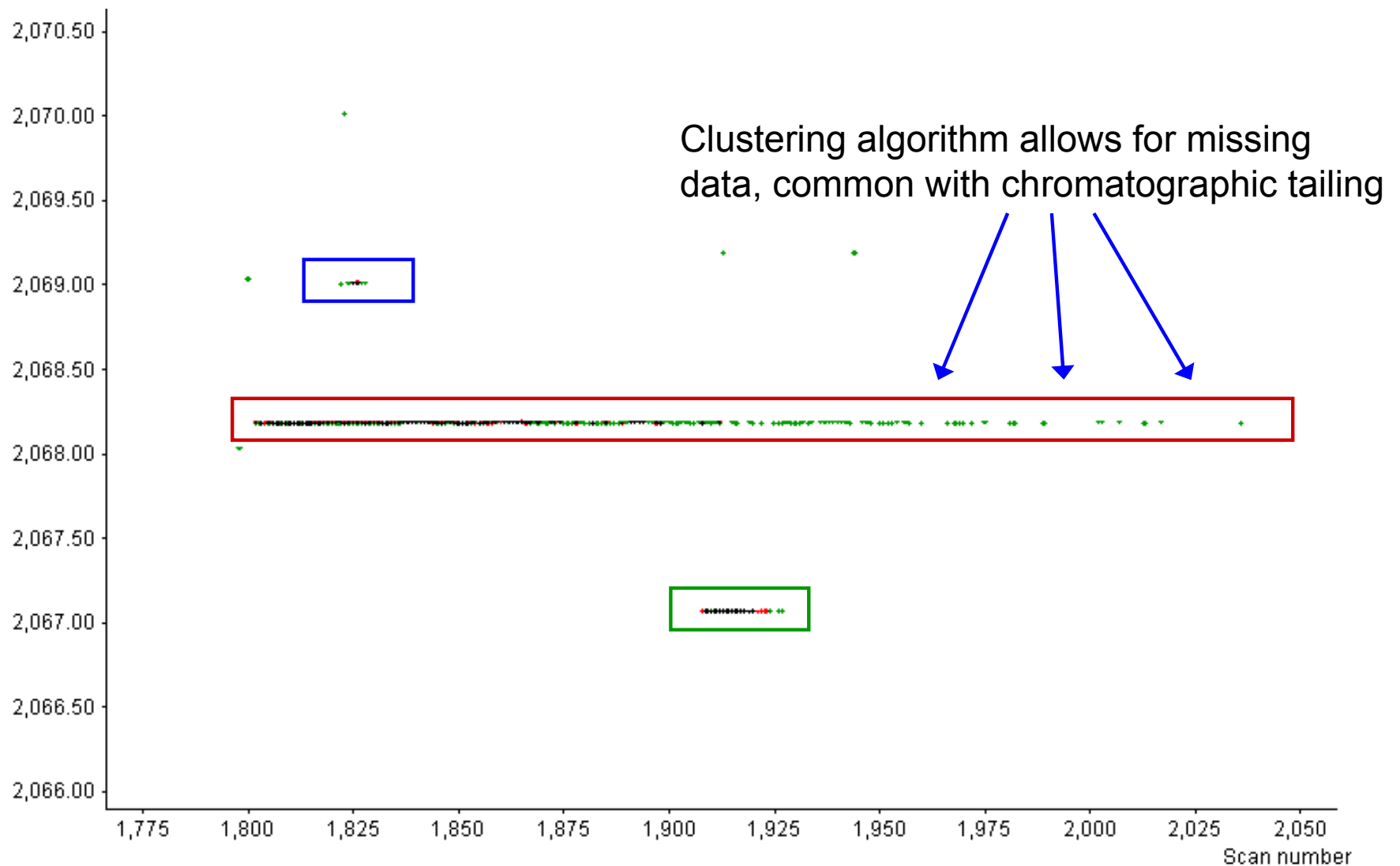  - ● Same species in multiple spectra need to be grouped together



Related peaks found using a weighted Euclidean distance; considers:
- ● Mass
- ● Abundance
- ● Elution time
- ● Isotopic Fit

# Feature definition over elution time

- **Feature detail**
  - Median Mass: 1904.9399 Da  (more tolerant to outliers than average)
  - Elution Time: Scan 1757 (0.363 NET)
  - Abundance: $1.7 \times 10^7$ counts (area under 2+ SIC)
    - See both 2+ and 3+ data
    - Stats typically come from the most abundant charge state

# Feature definition over elution time

- **Second example**
  - LC-MS feature eluting over 7.5 minutes



Clustering algorithm allows for missing data, common with chromatographic tailing

Scan number

# Feature definition over elution time

- Second example, feature detail
  - Median Mass: 2068.1781 Da
  - Elution Time: Scan 1809 (0.380 NET)
  - Abundance: $8.7 \times 10^7$ counts (area under 3+ SIC)
    - This example has primarily 3+ data; previous had even mix of 2+ and 3+ data

# Feature definition over elution time

- Example: *S. typhimurium* dataset on 11T FTICR

  - 100 minute LC-MS analysis (3360 mass spectra)
  - 67 cm, 150 μm I.D. column with 5 μm $C_{18}$ particles
  - 78,641 deisotoped peaks
  - Group into 5910 LC-MS Features

# Isotopic Pairs Processing

- Paired features typically have identical sequences, with and without an isotopic label
  - e.g. $^{16}O/^{18}O$ pairs have 4 Da spacing due to two $^{18}O$ atoms

# Isotopic Pairs Processing

- Paired feature example: $^{16}O/^{18}O$ data
  - Compute AR using ratio of areas, or
  - Compute AR scan-by-scan, then average AR values (members must co-elute)



Monoisotopic Mass

4.0085 Da

AR = 1.78 (Light$_{Area}$÷Heavy$_{area}$); or
AR = 1.34 ± 0.2 (scan-by-scan)



Monoisotopic Mass

4.0085 Da

AR = 0.13 (Light$_{Area}$÷Heavy$_{area}$); or
AR = 0.12 ± 0.02 (scan-by-scan)

# Feature definition over elution time

- Numerous options in VIPER for clustering data to form LC-MS features and for finding paired features

# Part II: LC-MS Feature Discovery

- Introduction (Adkins)
- Part I: Overview of Label-Free Quantitative Proteomics (Jaffe)
- Part II: Feature discovery in LC-MS datasets (Monroe and Jaitly)
  - ✅ Structure of LC-MS Data
  - ✅ Feature discovery in individual spectra (deisotoping)
  - ✅ Feature definition over elution time
  - Identifying LC-MS Features using an AMT tag DB
  - Extending the AMT tag approach for feature based analyses
  - Estimating confidence of identified LC-MS features
  - Downstream quantitative analysis with DAnTE
- Part III: PEPPeR, GenePattern and Real-world examples (Jaffe)
- Break
- AMT tag Pipeline Demo (general)
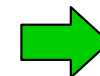- Panel Discussion
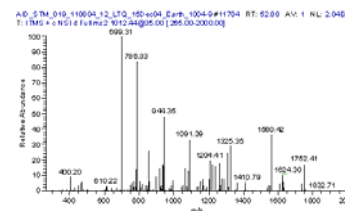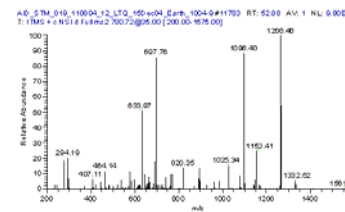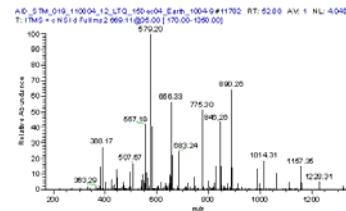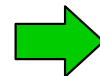
# Assembling an AMT tag DB

- Accurate Mass and Time (AMT) tag
  - Unique peptide sequence whose monoisotopic mass and normalized elution time are accurately known
  - AMT tags also track any modified residues in peptide
- AMT tag DB
  - Collection of AMT tags
- AMT tag approach articles
  - R.D. Smith et. al., *Proteomics* **2002**, *2*, 513-523.
  - J.S. Zimmer, M.E. Monroe et. al., *Mass Spec. Reviews* **2006**, *25*, 450-482.
  - L. Shi, J.N. Adkins, et. al., *J. of Biological Chem.* **2006**, *281*, 29131-29140.

# Assembling an AMT tag DB

- **What can we use an AMT tag DB for?**
  - Query LC-MS/MS data to answer questions
    - How many distinct peptides were observed passing filter criteria?
    - Which peptides were observed most often by LC-MS/MS?
    - How many proteins had 2 or more partially or fully tryptic peptides?
  - Correlate LC-MS features to the AMT tags
    - Analyze multiple, related samples by LC-MS using a high mass accuracy mass spectrometer
      - e.g. Time course study, 5 data points with 3 points per sample
    - Characterize the LC-MS features
      - Deisotope to obtain monoisotopic mass and charge
      - Cluster in time dimension to obtain abundance information
    - Match to AMT tags to identify peptides
      - Align in mass and time dimensions
      - Match mass and time of LC-MS features to mass and time of AMT tags

# Assembling an AMT tag DB

- **Characterizing AMT tags**
  - Analyze samples by LC-MS/MS
    - 10 minute to 180 minute LC separations
    - Obtain 1000's of MS/MS fragmentation spectra for each sample
  - Analyze spectra using SEQUEST, X!Tandem, etc.
    - SEQUEST: http://www.thermo.com/bioworks/
    - X!Tandem: http://www.thegpm.org/TANDEM/
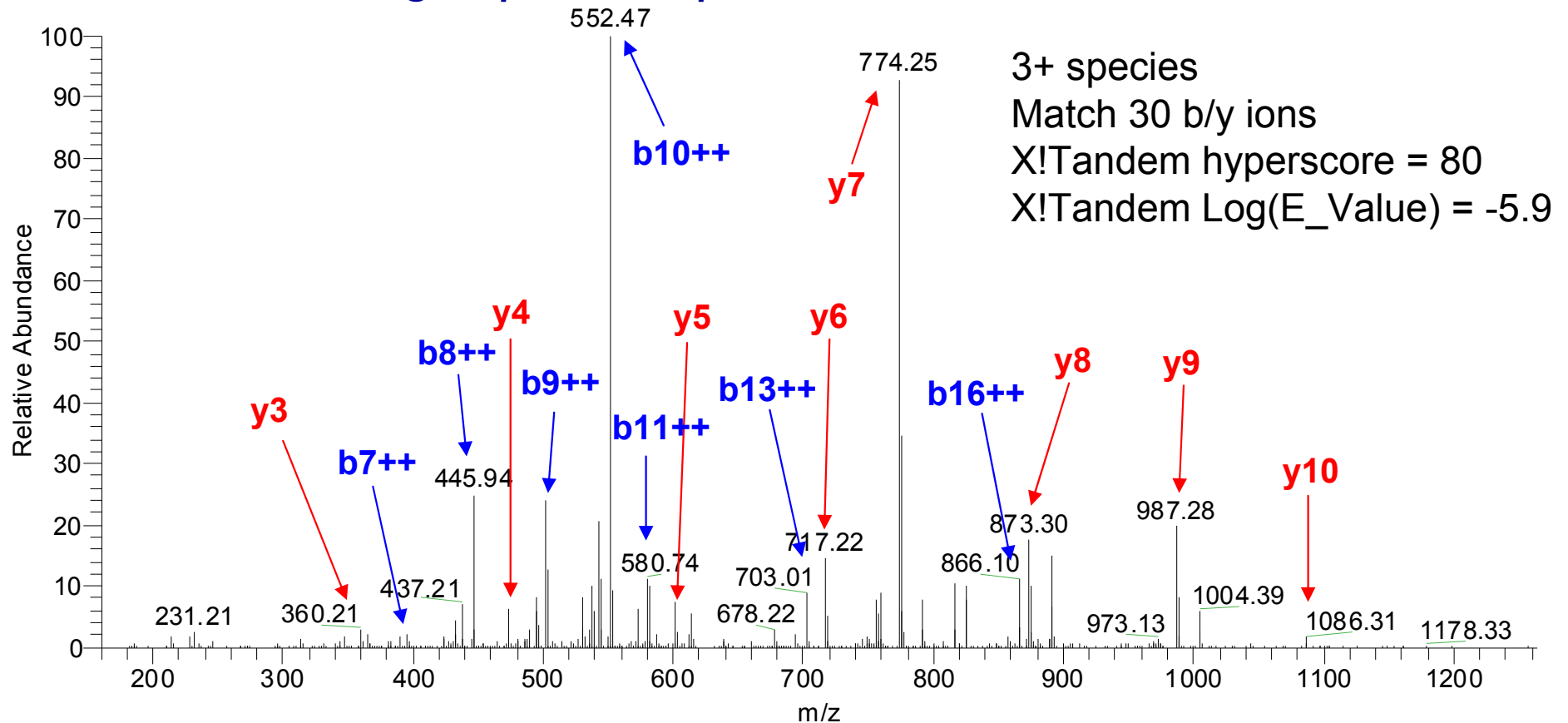  - Collate results



List of peptide and protein matches

# Assembling an AMT tag DB

- **AMT tag example**
  - R.VKHPSEIVNVGDEINVK.V
  - Observed in scan 11195 of dataset #19 in an SCX fractionation series

# Assembling an AMT tag DB

- **AMT tag example**
  - R.VKHPSEIVNVGDEINVK.V
  - Observed in scan 11195 of dataset #19 in an SCX fractionation series

| # | Immon. | b | b++ | Seq. | y | y++ | # |
|---|--------|------|------|------|---------|--------|----|
| 1 | 72.08 | | | V | 1877.01 | 939.01 | 17 |
| 2 | 101.11 | 228.17 | | K | 1777.94 | 889.48 | 16 |
| 3 | 110.07 | 365.23 | | H | 1649.85 | 825.43 | 15 |
| 4 | 70.07 | 462.28 | | P | 1512.79 | 756.90 | 14 |
| 5 | 60.04 | 549.31 | | S | 1415.74 | 708.37 | 13 |
| 6 | 102.06 | 678.36 | 339.68 | E | 1328.71 | 664.86 | 12 |
| 7 | 86.10 | 791.44 | 396.22 | I | 1199.66 | 600.33 | 11 |
| 8 | 72.08 | 890.51 | 445.76 | V | 1086.58 | 543.79 | 10 |
| 9 | 87.06 | 1004.55 | 502.78 | N | 987.51 | 494.26 | 9 |
| 10 | 72.08 | 1103.62 | 552.31 | V | 873.47 | 437.24 | 8 |
| 11 | 30.03 | 1160.64 | 580.83 | G | 774.40 | 387.70 | 7 |
| 12 | 88.04 | 1275.67 | 638.34 | D | 717.38 | 359.19 | 6 |
| 13 | 102.06 | 1404.71 | 702.86 | E | 602.35 | 301.68 | 5 |
| 14 | 86.10 | 1517.80 | 759.40 | I | 473.31 | | 4 |
| 15 | 87.06 | 1631.84 | 816.42 | N | 360.22 | | 3 |
| 16 | 72.08 | 1730.91 | 865.96 | V | 246.18 | | 2 |
| 17 | 101.11 | | | K | 147.11 | | 1 |

3+ species
Match 30 b/y ions
X!Tandem hyperscore = 80
X!Tandem Log(E_Value) = -5.9

# Assembling an AMT tag DB

- Align related datasets using elution times of observed peptides
  - One option: utilize NET prediction algorithm to create theoretical dataset to align against
    - NET prediction uses position and ordering of amino acid residues to predict normalized elution time

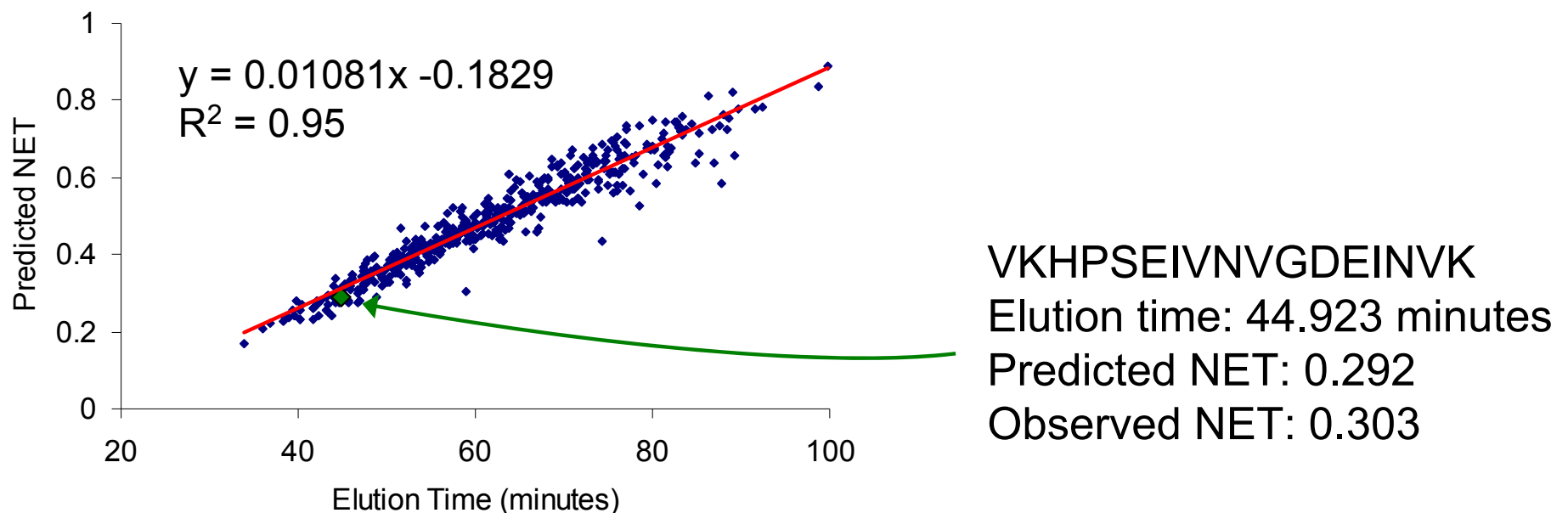| Peptide | X!Tandem Log (E_Value) | Elution Time | Predicted NET |
|---|---|---|---|
| R.AARPAKYSYVDENGETK.T | -6.1 | 33.958 | 0.167 |
| R.LVHGEEGLVAAKR.I | -8.8 | 36.915 | 0.224 |
| R.GIIKVGEEVEIVGIK.E | -8.2 | 53.003 | 0.415 |
| K.RFNDDGPILFIHTGGAPALFAYHPHV.- | -7.3 | 62.583 | 0.519 |
| K.KTGVLAQVQEALKGLDVR.E | -11.6 | 62.803 | 0.438 |
| R.KVAAQIPNGSTLFIDIGTTPEAVAHALLGHSNLR.I | -8.9 | 73.961 | 0.589 |
| R.TFAISPGHMNQLRAESIPEAVIAGASALVLTSYLVR.C | -6.5 | 88.043 | 0.764 |

K. Petritis, L.J. Kangas, P.L. Ferguson, et al., *Analytical Chemistry* **2003**, *75*, 1039-1048.
K. Petritis, L.J. Kangas, B. Yan, et al., *Analytical Chemistry* **2006**, *78*, 5026-5039.

# Assembling an AMT tag DB

- Align related datasets using elution times of observed peptides
  - One option: utilize NET prediction algorithm to create theoretical dataset to align against
    - NET prediction uses position and ordering of amino acid residues to predict normalized elution time
  - Alignment yields NET values based on observed elution times
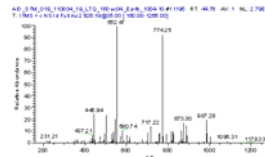    - Observed NET = Slope×(Observed Elution Time) + Intercept

Example: 506 unique peptides used for alignment; Log(E_Value) ≤ -6

$y = 0.01081x - 0.1829$
$R^2 = 0.95$

Predicted NET (y-axis, 0 to 1)
Elution Time (minutes) (x-axis, 20 to 100)

VKHPSEIVNVGDEINVK
Elution time: 44.923 minutes
Predicted NET: 0.292
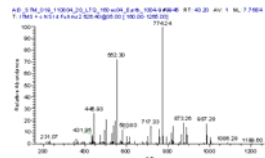Observed NET: 0.303

# Assembling an AMT tag DB

- **AMT tag example**
  - R.VKHPSEIVNVGDEINVK.V
  - Observed in 7 (of 25) LC-MS/MS datasets in the SCX fractionation series
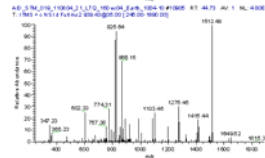
Analysis 1, scan 11195 ⟶ 3+, hyperscore 80, Obs. NET 0.303
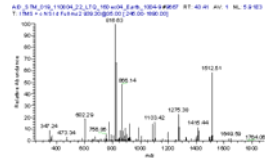
Analysis 2, scan 9945 ⟶ 3+, hyperscore 69, Obs. NET 0.298

Analysis 3, scan 10905 ⟶ 2+, hyperscore 74, Obs. NET 0.301

Analysis 4, scan 9667 ⟶ 2+, hyperscore 77, Obs. NET 0.302

⋮                            ⋮

Compute monoisotopic mass: 1876.0053 Da
Average Normalized Elution Time: 0.3021 (StDev 0.0021)

# Assembling an AMT tag DB

- **Mass and Time Tag Database**
  - Repository for AMT tags
    - Mass, elution time, modified residues, and supporting information for each AMT tag
  - Allows samples of unknown composition to be matched quickly and efficiently, without needing to perform tandem MS
  - Assembled by analyzing a control set of samples, cataloging each peptide identification until subsequent analyses no longer provide new identifications

| MT Tag ID | Peptide | LC-MS/MS Obs. Count | Calculated Monoisotopic Mass | Average Observed NET | Observed NET StDev |
|---|---|---|---|---|---|
| 1662039 | MTGRELKPHDR | 1 | 1338.6826 | 0.143 | 0.000 |
| 17683899 | SSALNTLTNQK | 3 | 1175.6146 | 0.235 | 0.005 |
| 36609588 | HRDLLGATNP…TLR | 5 | 1960.0602 | 0.379 | 0.002 |
| 36715875 | WVKVDGWDN…FER | 11 | 2590.2815 | 0.459 | 0.011 |
| 36843675 | MYGHLKGEVA…QER | 8 | 2533.2304 | 0.557 | 0.005 |

# Assembling an AMT tag DB

- **Mini AMT tag DB**
  - Database constructed from a relatively small number of datasets
  - e.g. 25 SCX fractionation samples from S. typhimurium, each analyzed by LC-MS/MS and then by X!Tandem
  - Protein database: S_typhimurium_LT2_2004-09-19
    - 4550 proteins and 1.4 million residues

>STM1834 putative YebN family transport protein (yebN) {Salmonella typhimurium LT2}

MFAGGSDVFNGYPGQDVVMHFTATVLLAFGMSMDAFAASIGKGATLHKPKFSEALRTGLI

FGAVETLTPLIGWGLGILASKFVLEWNHWIAFVLLIFLGGRMIIEGIRGGSDEDETPLRR

HSFWLLVTTAIATSLDAMAVGVGLAFLQVNIIATALAIGCATLIMSTLGMMIGRFIGPML

GKRAEILGGVVLIGIGVQILWTHFHG

>STM1835 23S rRNA m1G745 methyltransferase (rrmA) {Salmonella typhimurium LT2}

MSFTCPLCHQPLTQINNSVICPQRHQFDVAKEGYINLLPVQHKRSRDPGDSAEMMQARRA

FLDAGHYQPLRDAVINLLRERLDQSATAILDIGCGEGYYTHAFAEALPGVTTFGLDVAKT

AIKAAAKRYSQVKFCVASSHRLPFADASMDAVIRIYAPCKAQELARVVKPGGWVVTATPG

PHHLMELKGLIYDEVRLHAPYTEQLDGFTLQQSTRLAYHMQLTAEAAVALLQMTPFAWRA

RPDVWEQLAASAGLSCQTDFNLHLWQRNR

# Assembling an AMT tag DB

- **Database Relationships**
  - Minimum information required:
    - Single table with Mass and NET
  - Expanded schema:

| T_Mass_Tags | |
|---|---|
| PK | Mass_Tag_ID |
| | Peptide<br>Monoisotopic_Mass<br>NET |

| T_Mass_Tags | |
|---|---|
| PK | Mass_Tag_ID |
| | Peptide<br>Monoisotopic_Mass |

| T_Mass_Tags_to_Protein_Map | |
|---|---|
| PK,FK1<br>PK,FK2 | Mass_Tag_ID<br>Ref_ID |
| | |

| T_Proteins | |
|---|---|
| PK | Ref_ID |
| | Reference<br>Description |

| T_Mass_Tags_NET | |
|---|---|
| PK,FK1 | Mass_Tag_ID |
| | Avg_GANET<br>Cnt_GANET<br>StD_GANET |

PK := Primary Key
FK := Foreign Key

# Assembling an AMT tag DB

- ● Microsoft Access DB Relationships
  - ● Full schema to track individual peptide observations

**T_Analysis_Description**

| PK | Job |
|----|-----|
| | Dataset |
| | Dataset_ID |
| | Dataset_Created_DMS |
| | Dataset_Acq_Time_Start |
| | Dataset_Acq_Time_End |
| | Dataset_Scan_Count |
| | Experiment |
| | Campaign |
| | Organism |
| | Instrument_Class |
| | Instrument |
| | Analysis_Tool |
| | Parameter_File_Name |
| | Settings_File_Name |
| | Organism_DB_Name |
| | Protein_Collection_List |
| | Protein_Options_List |
| | Completed |
| | ResultType |
| | Separation_Sys_Type |
| | ScanTime_NET_Slope |
| | ScanTime_NET_Intercept |
| | ScanTime_NET_RSquared |
| | ScanTime_NET_Fit |

**T_Peptides**

| PK | Peptide_ID |
|----|-----------|
| FK1 | Analysis_ID |
| | Scan_Number |
| | Number_Of_Scans |
| | Charge_State |
| | MH |
| | Multiple_Proteins |
| | Peptide |
| FK2 | Mass_Tag_ID |
| | GANET_Obs |
| | Scan_Time_Peak_Apex |
| | Peak_Area |
| | Peak_SN_Ratio |

**T_Mass_Tags**

| PK | Mass_Tag_ID |
|----|-------------|
| | Peptide |
| | Monoisotopic_Mass |
| | Multiple_Proteins |
| | Created |
| | Last_Affected |
| | Number_Of_Peptides |
| | Peptide_Obs_Count_Passing_Filter |
| | High_Normalized_Score |
| | High_Peptide_Prophet_Probability |
| | Mod_Count |
| | Mod_Description |
| | PMT_Quality_Score |

**T_Mass_Tags_to_Protein_Map**

| PK,FK1 | Mass_Tag_ID |
|--------|-------------|
| PK,FK2 | Ref_ID |
| | Mass_Tag_Name |
| | Cleavage_State |
| | Fragment_Number |
| | Fragment_Span |
| | Residue_Start |
| | Residue_End |
| | Repeat_Count |
| | Terminus_State |
| | Missed_Cleavage_Count |

**T_Score_Sequest**

| PK,FK1 | Peptide_ID |
|--------|-----------|
| | XCorr |
| | DelCn |
| | Sp |
| | DelM |

**T_Score_XTandem**

| PK,FK1 | Peptide_ID |
|--------|-----------|
| | Hyperscore |
| | Log_EValue |
| | DeltaCn2 |
| | Y_Score |
| | Y_Ions |
| | B_Score |
| | B_Ions |
| | DelM |
| | Intensity |
| | Normalized_Score |

**T_Mass_Tags_NET**

| PK,FK1 | Mass_Tag_ID |
|--------|-------------|
| | Min_GANET |
| | Max_GANET |
| | Avg_GANET |
| | Cnt_GANET |
| | StD_GANET |
| | StdError_GANET |
| | PNET |

**T_Proteins**

| PK | Ref_ID |
|----|--------|
| | Reference |
| | Description |
| | Protein_Sequence |
| | Protein_Residue_Count |
| | Monoisotopic_Mass |
| | Protein_Collection_ID |
| | Last_Affected |

**V_Filter_Set_Overview_Ex**

| | |
|--|--|
| | Filter_Type |
| | Filter_Set_ID |
| | Extra_Info |
| | Filter_Set_Name |
| | Filter_Set_Description |

**T_Score_Discriminant**

| PK,FK1 | Peptide_ID |
|--------|-----------|
| | Peptide_Prophet_FScore |
| | Peptide_Prophet_Probability |

# Assembling an AMT tag DB

- Example data

## T_Mass_Tags

| Mass_Tag_ID | Peptide | Monoisotopic_Mass |
|---|---|---|
| 24847 | VKHPSEIVNVGDEINVK | 1876.00533 |

## T_Mass_Tags_NET

| Mass_Tag_ID | Avg_GANET | Cnt_GANET | StD_GANET |
|---|---|---|---|
| 24847 | 0.3021 | 7 | 2.11E-03 |

## T_Peptides

| Peptide_ID | Peptide | Mass Tag ID | Job | Scan Number | Charge State |
|---|---|---|---|---|---|
| 53428 | R.VKHPSEIVNVGDEINVK.V | 24847 | 206386 | 11195 | 3 |
| 57461 | R.VKHPSEIVNVGDEINVK.V | 24847 | 206387 | 9945 | 3 |
| 61511 | R.VKHPSEIVNVGDEINVK.V | 24847 | 206388 | 10905 | 2 |
| 65386 | R.VKHPSEIVNVGDEINVK.V | 24847 | 206389 | 9667 | 2 |
| 69081 | R.VKHPSEIVNVGDEINVK.V | 24847 | 206390 | 9118 | 2 |
| 72556 | R.VKHPSEIVNVGDEINVK.V | 24847 | 206391 | 9159 | 2 |
| 76263 | R.VKHPSEIVNVGDEINVK.V | 24847 | 206392 | 9421 | 2 |

## T_Score_XTandem

| Peptide_ID | Hyperscore | Log(E_Value) |
|---|---|---|
| 53428 | 80.2 | -5.89 |
| 57461 | 69.2 | -4.92 |
| 61511 | 74 | -12.85 |
| 65386 | 77.2 | -12.80 |
| 69081 | 69 | -12.82 |
| 72556 | 78 | -13.77 |
| 76263 | 60.3 | -11.27 |

# Assembling an AMT tag DB

- Processing steps

| Flowchart | Description |
|---|---|
| **Thermo-Finnigan LTQ .Raw files** | Convert to .Dta files or single _Dta.txt file using DeconMSn.exe. DeconMSn is similar to Thermo's Extract_MSn but has better support for data from LTQ-Orbitrap or LTQ-FT instruments. |
| ↓ | |
| **MS/MS spectra files** | Process _Dta.txt file with X!Tandem or .Dta files with SEQUEST.  Use the Peptide File Extractor to convert SEQUEST .Out files to Synopsis (_Syn.txt) files. |
| ↓ | |
| **Peptide ID Results** | Convert X!Tandem .XML output files or SEQUEST _Syn.txt file to tab-delimited files using the Peptide Hit Results Processor (PHRP) application. |
| ↓ | |
| **Tab delimited text files** | |
| ↓ | Align datasets using the MTDB Creator application |
| **Summarized result files** | |
| ↓ | Load into database using MTDB Creator |
| **Microsoft Access DB** | |

# DeconMSn

- Determines the monoisotopic mass and charge state of each parent ion chosen for fragmentation on a hybrid LC-MS/MS instrument using Decon2LS algorithms

- Replacement for the Extract_MSn.exe tool provided with SEQUEST and Bioworks

# Assembling an AMT tag DB

- Peptide Hit Results Processor (PHRP) relationships

| Results_Info | |
|---|---|
| PK | **Result_ID** |
| | |
| FK1 | **Unique_Seq_ID** |
| | Group_ID |
| | Scan |
| | Charge |
| | Peptide_MH |
| | Peptide_Hyperscore |
| | Peptide_Expectation_Value_Log(e) |
| | Multiple_Protein_Count |
| | Peptide_Sequence |
| | DeltaCn2 |
| | y_score |
| | y_ions |
| | b_score |
| | b_ions |
| | Delta_Mass |
| | Peptide_Intensity_Log(I) |

| Result_To_Seq_Map | |
|---|---|
| PK,FK1 | **Unique_Seq_ID** |
| PK,FK2 | **Result_ID** |
| | |

| Seq_Info | |
|---|---|
| PK | **Unique_Seq_ID** |
| | |
| | Mod_Count |
| | Mod_Description |
| | Monoisotopic_Mass |

| Mod_Details | |
|---|---|
| PK,FK1 | **Unique_Seq_ID** |
| | |
| | Mass_Correction_Tag |
| | Position |

| Seq_to_Protein_Map | |
|---|---|
| PK,FK1 | **Unique_Seq_ID** |
| PK | **Protein_Name** |
| | |
| | Cleavage_State |
| | Terminus_State |
| | Protein_Expectation_Value_Log(e) |
| | Protein_Intensity_Log(I) |

# MTDB Creator

- ## MTDB Creator application

  - Allows external researchers to align multiple LC-MS/MS analyses, run PeptideProphet (for SEQUEST data) and create a standalone AMT tag database

# Assembling an AMT tag DB

- Database histograms – filtered on Log(E_Value) ≤ -2



Peptide Mass Histogram



NET Histogram



X!Tandem Hyperscore Histogram

# AMT Tag DB Growth Trend

- **Trend for Mini AMT tag DB**
  - 25 SCX fractionation datasets of a single growth condition



Filtered on Log(E_Value) ≤ -2

*(Y-axis: Peptide Count, 0 to 20000; X-axis: Dataset Count, 0 to 25)*

- **Trend for Mature AMT tag DB**
  - 521 different samples from ~25 different conditions
  - Slope of curve decreases as more datasets are added and as fewer new peptides are seen



Filtered on Peptide Prophet Probability ≥ 0.99

*(Y-axis: Peptide Count, 0 to 60000; X-axis: Dataset Count, 0 to 600)*

# Identifying LC-MS Features

- VIPER software
  - Visualize and find features in LC-MS data
  - Match features to peptides (AMT tags)
  - Graphical User Interface and automated analysis mode

# Identifying LC-MS Features

- **Peak Matching Steps**
  - ✔ Load LC-MS peak lists from Decon2LS
  - ✔ Filter data
  - ✔ Feature definition over elution time
  - Select AMT tags to match against
  - Optionally, find paired features (e.g. $^{16}O/^{18}O$ pairs)
  - Align LC-MS features to AMT tags using LCMSWarp
  - Broad AMT tag DB search
  - Search tolerance refinement
  - Final AMT tag DB search
  - Report results

VIPER

File   Steps   Edit   Info   View   Tools   Special   Window   Help

1a. Load peak list file ▶
1b. Filter
2. Find LC-MS Features (UMCs) ▶
3. Select MT Tags (Connect to DB)
4. Find Pairs ▶
5. Align LC-MS Features to MT Tags ▶
6. Database Search ▶
7. Mass Calibration and Tolerance Refinement
8. Database Search using Pairs ▶
9. Save QC Plots...

View Analysis History Log

# Identifying LC-MS Features

- AMT Tag database selection



Connect to mass tag system (MTS) if inside PNNL or use standalone Microsoft Access DB

# Alignment using LCMSWarp

- Align scan number (i.e. elution time) of features to NETs of peptides in given AMT tag database
  - Match mass and NET of AMT tags to mass and scan number of MS features
  - Use LCMSWarp algorithm to find optimal alignment to give the most matches

Calculated
monoisotopic mass

AMTs

Average observed NET

Deisotoped
monoisotopic mass

LC-MS Features

Observed scan number

# Alignment using LCMSWarp

- LCMSWarp computes a similarity score from conserved local mass and retention time patterns



Best score = 0.00681
Scan = 1113
Shift = 113

N. Jaitly, M.E. Monroe et. al.,
*Analytical Chemistry* **2006**, *78*,
7397-7409.

Alignment Score

Scan number

# Alignment using LCMSWarp

- Similarity scores between LC-MS features and AMT tags are used to generate a score graph of similarity

- Best alignment is found using a dynamic programming algorithm that determines the transformation function with maximum likelihood



*S. typhimurium* on 11T

AMT tag NET

MS Scan Number

Heatmap of similarity score between LC-MS features and AMT tags (z-score representation)

Alignment Function

N. Jaitly, M.E. Monroe et. al., *Analytical Chemistry* **2006**, 78, 7397-7409.

# Alignment using LCMSWarp

- Transformation function is used to convert from scan number to NET
  - Features centered at same scan number get the same obs. NET value
  - When matching LC-MS features to AMTs, we will search +/- a NET tolerance, which effectively allows for LC-MS features to shift around a little in elution time

| LC-MS Feature Scan | Matching AMT tag NET | LC-MS Feature NET |
|---|---|---|
| 1011 | 0.1519 | 0.1569 |
| 1019 | 0.1626 | 0.1589 |
| 1019 | 0.1507 | 0.1589 |
| 1021 | 0.1653 | 0.1594 |
| 1027 | 0.1509 | 0.1609 |
| 1037 | 0.1519 | 0.1633 |
| 1042 | 0.183 | 0.1645 |
| 1055 | 0.1652 | 0.1677 |
| 1056 | 0.1862 | 0.1679 |
| 1056 | 0.1697 | 0.1679 |
| 1056 | 0.1682 | 0.1679 |

# Alignment using LCMSWarp

- **NET Residual Plots**
  - Difference between NET of LC-MS feature and NET of matching AMT tag
    - Indicates quality of alignment between features and AMT tags
  - This data shows nearly linear alignment between features and AMTs, but the algorithm can easily account for non-linear trends



*S. typhimurium* on 11T

AMT tag NET

MS Scan Number

NET Residuals if a linear mapping is used



NET Residuals after LCMSWarp

# Alignment using LCMSWarp

- **Non-linear alignment example #1**
  - Identical LC separation system, but having column flow irregularities

AMT tag NET

*S. typhimurium* on 9T

MS Scan Number

NET Residuals if a linear mapping is used

NET Residuals after LCMSWarp

# Alignment using LCMSWarp

- Non-linear alignment example #2
  - AMT Tag DB from $C_{18}$ LC-MS/MS analyses using ISCO-based LC (**exponential dilution gradient**)
  - LC-MS analysis used $C_{18}$ LC-MS via Agilent **linear gradient** pump

*S. oneidensis* on LTQ-Orbitrap

NET Residuals if a linear mapping is used

NET Residuals after LCMSWarp

# Alignment using LCMSWarp

- Non-linear alignment example #3
  - AMT Tag DB from $C_{18}$ LC-MS/MS analyses using ISCO-based LC
  - LC-MS analysis used $C_{18}$ LC-MS via Agilent linear gradient pump

NET Residuals if a linear mapping is used



QC Standards (12 protein digest) on LTQ-Orbitrap



NET Residuals after LCMSWarp

# Alignment using LCMSWarp

- **LCMSWarp Features**
  - Fast and robust
    - Previous method used least-squares regression, iterating through a large range of guesses (slow and often gave poor alignment)
  - Requires that a reasonable number of LC-MS features match the AMT Tag DB



*S. typhimurium* on 11T
match against 18,617 *S. typhimurium* PMTs

*S. typhimurium* on 11T
match against 65,193 *S. oneidensis* PMTs

# Alignment using LCMSWarp

- In addition to aligning data in time, we can also recalibrate the masses of the LC-MS features
  - Possible because mass and time values are available for both LC-MS features and AMT tags
- Two options for mass re-calibration
  - Bulk linear correction
  - Piece-wise correction via LCMSWarp
- Visualize mass differences using mass error histogram or mass residual plot

# Mass Error Histogram

- **List of binned mass error values**
  - Difference between feature's mass and matching AMT tag's mass
  - Bin values to generate a histogram
  - Typically observe background false positive level

| LC-MS Feature Mass (Da) | AMT Tag Mass (Da) | Delta Mass (Da) | Mass Error (ppm) |
|---|---|---|---|
| 1570.9005 | 1570.883 | 0.01745 | 11.1 |
| 1571.74325 | 1571.726 | 0.01770 | 11.3 |
| 1571.8498 | 1571.831 | 0.01912 | 12.2 |
| 1571.9107 | 1571.892 | 0.01848 | 11.8 |
| 1573.8381 | 1573.832 | 0.00569 | 3.6 |

Match Tolerances

Mass: ±25 ppm
NET: ±0.05 NET

Count (LC-MS Features)

Likely true positive identifications

Likely false positive identifications

Mass Error (ppm)

# Mass Calibration

- Option 1: Bulk linear correction
  - Use location of peak in mass error histogram to adjust masses of all features
  - Shift by ppm mass; absolute shift amount increases as monoisotopic mass increases

Peak Center of mass: 11.6 ppm
Peak Width: 2 ppm at 60% of max
Peak Height: 404 counts/bin
Noise level: 19 counts/bin

Count (LC-MS Features)

400

300

200

100

11.6 ppm

-10    0    10    20

Mass Error (ppm)

Shift all masses -11.6 ppm:

$$\Delta_{mass} = -11.6\text{ppm} \times \frac{mass_{old}}{1 \times 10^6 \text{ ppm/Da}}$$

For 1+ feature at 1570.9005 Da, $\Delta_{mass}$ = -0.0182 Da

For 3+ feature at 2919.4658 Da, $\Delta_{mass}$ = -0.0339 Da

# Mass Calibration

- Option 2: Piece-wise correction via LCMSWarp
  - Use smoothing splines to determine a smooth calibration curve which is a function of scan number



Mass Residual

Mass Error (ppm) vs. Scan Number

MS Scan Number

Mass Error (ppm) vs. Scan Number after correction

*S. typhimurium* on 11T

MS Scan Number

# Mass Calibration

- Option 2: Piece-wise correction via LCMSWarp
  - Use a smoothing spline calibration which is a function of m/z
  - LCMSWarp utilizes a hybrid correction based on both mass error vs. time and mass error vs. m/z



Mass Residual

Mass Error (ppm) vs. m/z

Mass Error (ppm) vs. m/z after correction

*S. typhimurium* on 11T

m/z

m/z

# Mass Calibration

- Comparison of the three methods
  - Mass error histogram gets taller, narrower, and more symmetric
    - Linear → Mass error vs. m/z → Mass error vs. time → Hybrid
  - Not all datasets show the same trends, but Hybrid mass recalibration is generally superior



*S. typhimurium* on 11T

*S. oneidensis* on LTQ-FT

# Identifying LC-MS Features

- Match Features to LC-MS/MS IDs

- *S. typhimurium* DB, from 25 LC-MS/MS analyses

  - 18,617 AMT tags, all fully or partially tryptic

  - Look for AMT tags within a broad mass range,
    e.g., ±25 ppm and ±0.05 NET of each feature

*S. typhimurium* AMT Tag Database



18,617 AMT tags

Average observed NET

*S. typhimurium* on 11T FTICR



5,934 features
4,678 features have match,
matching 6,242 AMT tags

Observed NET

# Search tolerance refinement

- Can use mass error and NET error histograms to determine optimal search tolerances



Examine distribution of errors to determine optimal tolerance using expectation maximization algorithm

±1.76 ppm

# Identifying LC-MS Features

- **Repeat search with final search tolerances**
  - 5,934 features
  - 3,866 features with matches
  - 3,958 out of 18,617 AMT tags matched using ±1.76 ppm

# Identifying LC-MS Features

- Caveat: given feature can match more than one AMT tag
  - Need measure of ambiguity

Match Tolerances

Mass: ±4 ppm
NET: ±0.02 NET

| AMT Tag ID | Peptide | Mass (Da) | NET |
|---|---|---|---|
| 35896216 | T.RALMQLDEALRPSLR.S | 1767.9777 | 0.373 |
| 105490 | K.DLETIVGLQTDAPLKR.A | 1767.9730 | 0.380 |
| 36259992 | R.SIGIAPDVLICRGDRAI.P | 1767.9664 | 0.392 |

1767.9727 Da
NET: 0.383

Monoisotopic Mass

Δ mass = 2.8 ppm
Δ NET = -0.010

Δ mass = 0.17 ppm
Δ NET = -0.003

1.6 ppm

Δ mass = -3.5 ppm
Δ NET = 0.009

NET

# Identifying LC-MS Features

$$d_{ij}^2 = \frac{(m_i - \mu_{mj})^2}{\sigma_{mj}^2} + \frac{(t_i - \mu_{tj})^2}{\sigma_{tj}^2}$$

$$\sigma_{mj} = 4 \text{ ppm}, \ \sigma_{tj} = 0.025$$

$$p_{ij} = \frac{(\sigma_{mj}\sigma_{tj})^{-1} \exp(-d_{ij}^2/2)}{\left( \sum_{k=1}^{N} (\sigma_{mk}\sigma_{tk})^{-1} \exp(-d_{ik}^2/2) \right)}$$

<u>Match Tolerances</u>

Mass: ±4 ppm
NET: ±0.02 NET

| AMT Tag ID | Mass (Da) | NET | $d_{ij}^2$ | Numerator | $p_{ij}$ |
|---|---|---|---|---|---|
| 35896216 | 1767.9777 | 0.373 | 3.012 | 6273.3 | 0.16 |
| 105490 | 1767.9730 | 0.380 | 0.090 | 27042.5 | 0.70 |
| 36259992 | 1767.9664 | 0.392 | 3.267 | 5521.4 | 0.14 |
| | | | Sum: | 38837.2 | |

K.K. Anderson, M.E. Monroe, and D.S. Daly. *Proteome Science* **2006**, *4*, 1.

# Identifying LC-MS Features

- VIPER reports a score that measures the uniqueness of each match

| AMT Tag ID | Peptide | Mass (Da) | NET | SLiC Score | Average XCorr | Avg Disc Score |
|---|---|---|---|---|---|---|
| 35896216 | T.RALMQLDEALRPSLR.S | 1767.9777 | 0.373 | 0.16 | 3.13 | 0.61 |
| 105490 | K.DLETIVGLQTDAPLKR.A | 1767.9730 | 0.380 | 0.70 | 3.68 | 0.97 |
| 36259992 | R.SIGIAPDVLICRGDRAI.P | 1767.9664 | 0.392 | 0.14 | 2.15 | 0.06 |



K.K. Anderson, M.E. Monroe, and D.S. Daly. *Proteome Science* **2006**, *4*, 1.

# Search tolerance refinement

- Effect of search tolerances on Mass Error histogram
  - If mass error plot not centered at 0, then narrow mass windows exclude valid data
  - Decreasing mass and/or NET tolerance reduces background false positive level



Mass error histograms with linear mass correction

Mass error histograms with LCMSWarp mass correction

# Automated Peak Matching

- Automated processing using VIPER
  - Processing steps and parameters defined in .Ini file
    - Separate .Ini file for $^{14}N/^{15}N$ pairs and $^{16}O/^{18}O$ pairs

# Peak Matching Results

- Browsable result folders for visual QC of each dataset
  - *S. typhimurium* on 11T FTICR



Data Searched



Data With Matches

- **2D Plot Metrics**
  - Reasonable number of matches
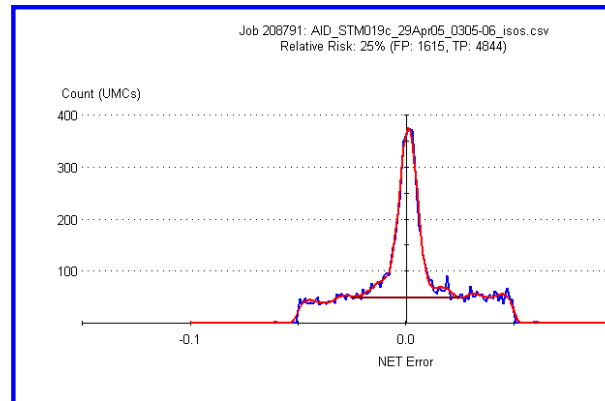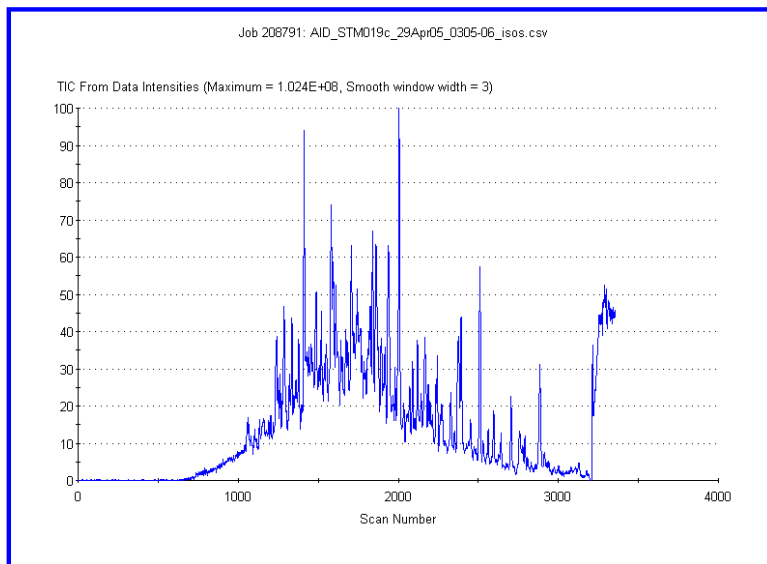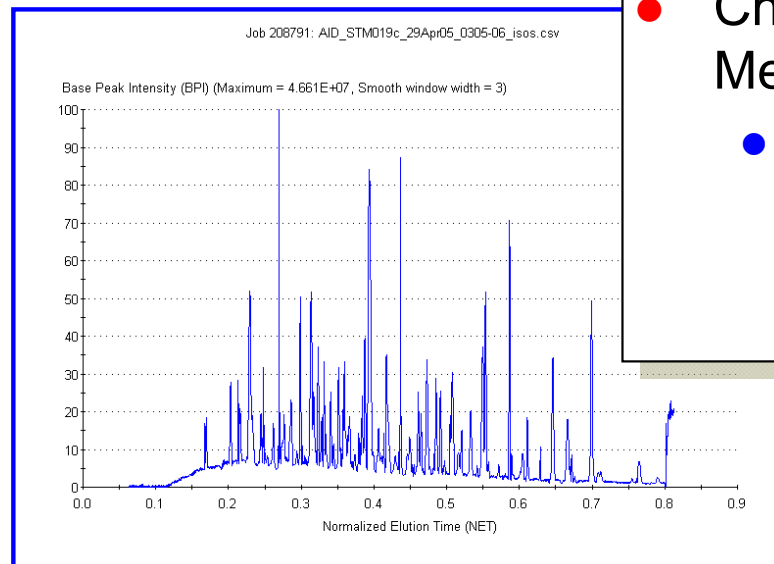  - NET range ≈ 0 to 1



Mass Errors Before Refinement



Mass Errors After Refinement

- **Mass Error Histogram Metrics**
  - Well defined, symmetric mass error peak centered at 0 ppm

# Peak Matching Results

- Browsable result folders for visual QC of each dataset
  - *S. typhimurium* on 11T FTICR



NET Errors Before Refinement



NET Errors After Refinement

- NET Error Histogram Metrics
  - Well defined, symmetric NET error peak centered at 0
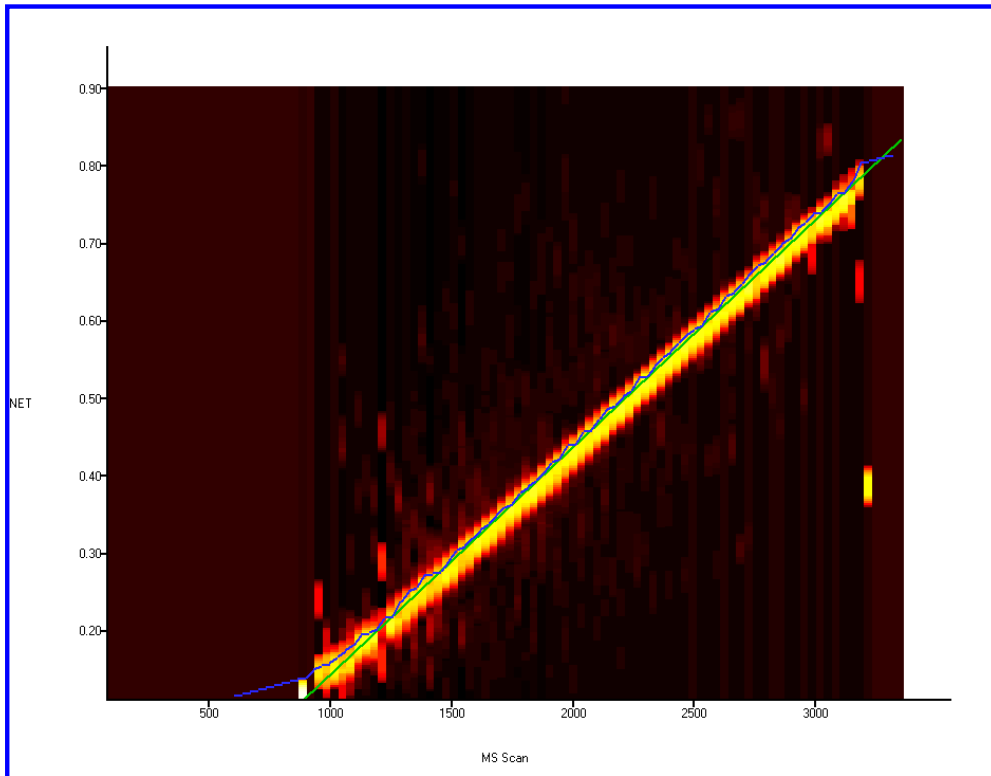


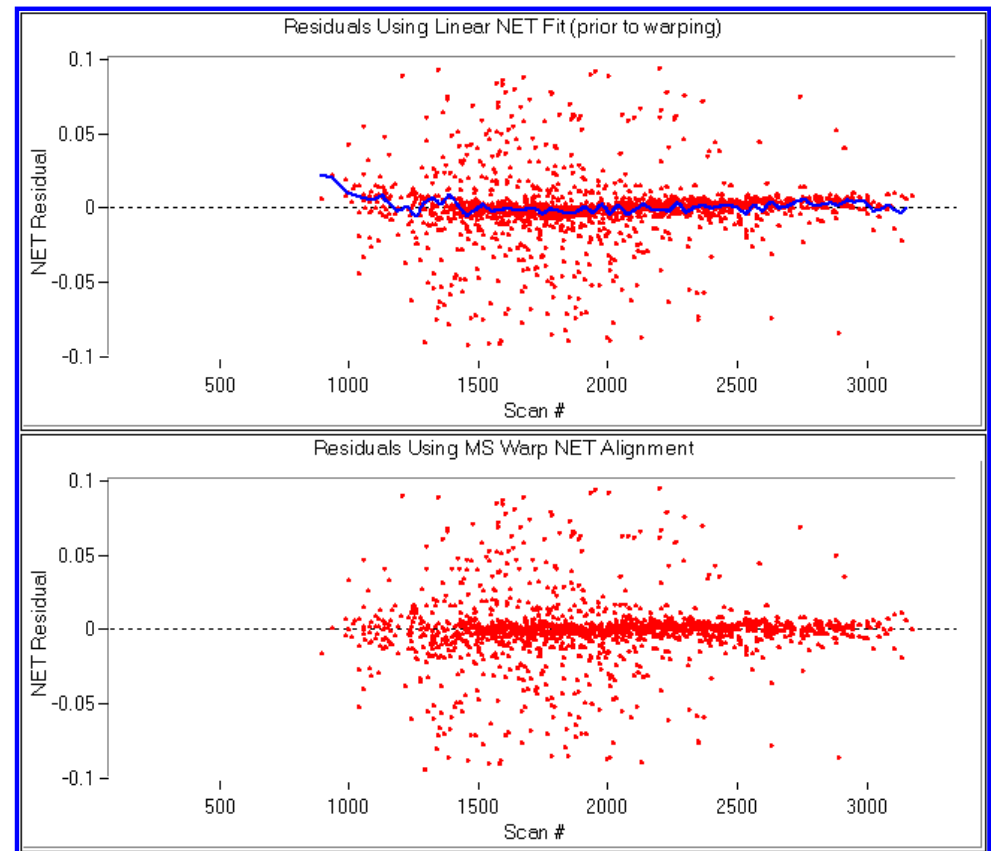Total Ion Chromatogram (TIC)



Base Peak Intensity (BPI) Chromatogram

- Chromatogram Metrics
  - Narrow peaks evenly distributed throughout separation window

# Peak Matching Results

- Browsable result folders for visual QC of each dataset
  - *S. typhimurium* on 11T FTICR



- NET Alignment Surface Metrics
  - Should show a smooth, bright yellow, diagonal line

- NET Alignment Residual Metrics
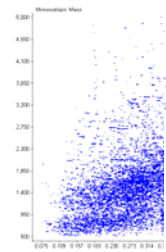  - Data after recalibration should be narrowly distributed around zero

# Part II: LC-MS Feature Discovery

- Introduction (Adkins)
- Part I: Overview of Label-Free Quantitative Proteomics (Jaffe)
- Part II: Feature discovery in LC-MS datasets (Monroe and Jaitly)
  - ✓ Structure of LC-MS Data
  - ✓ Feature discovery in individual spectra (deisotoping)
  - ✓ Feature definition over elution time
  - ✓ Identifying LC-MS Features using an AMT tag DB
  - Extending the AMT tag approach for feature based analyses
  - Estimating confidence of identified LC-MS features
  - Downstream quantitative analysis with DAnTE
- Part III: PEPPeR, GenePattern and Real-world examples (Jaffe)
- Break
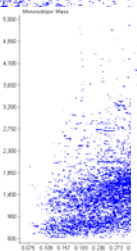- AMT tag Pipeline Demo (general)
- Panel Discussion

# Current AMT Tag Pipeline

- Individual LC-MS datasets are aligned to an AMT tag database independently

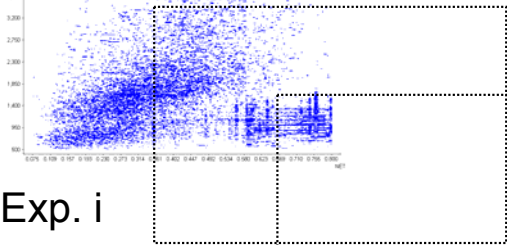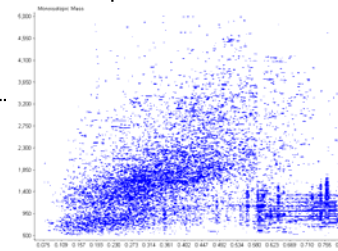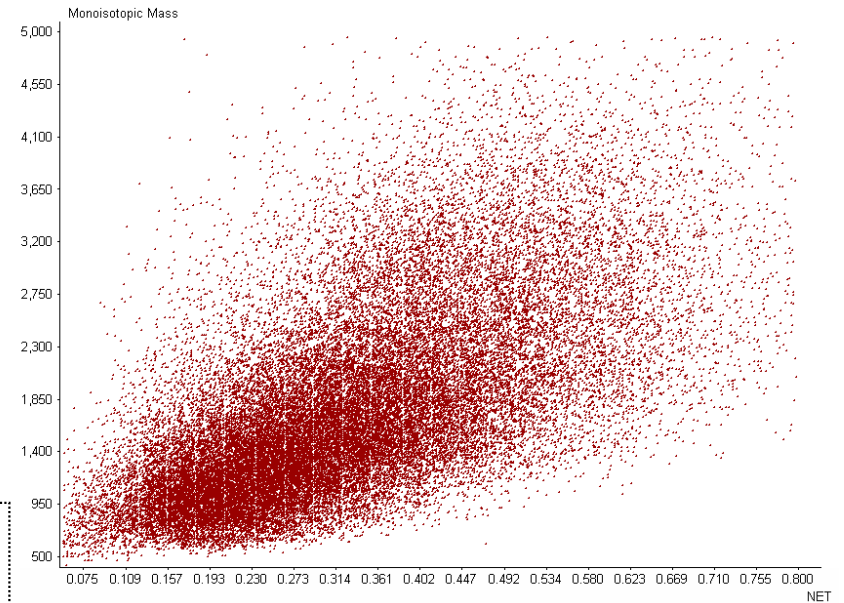- Results are combined together after independent processing

AMT tags from LC-MS/MS



LC-MS
Exp. 1

Exp. 2

Exp. i

Exp. 1600

# Current AMT Tag Pipeline
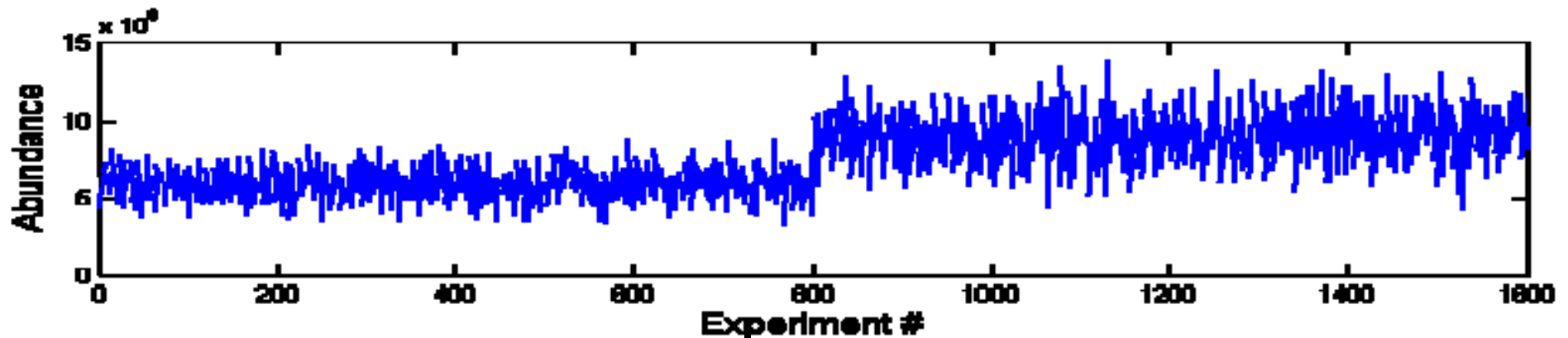
- For each peptide identified by peak matching, find the abundance of that peptide in all the peak matchings to create a profile

### LC-MS

| Experiment # | Scan # | Mass | Abundance |
|---|---|---|---|
| 1 | 2027 | 1063.56 | 3320000 |
| 2 | 2300 | 1063.56 | 3524300 |
| 3 | - | - | - |
| 1600 | 2400 | 1063.56 | 481000 |

### LC-MS/MS

| Peptide | NET | Mass | ORFName |
|---|---|---|---|
| TPHPALTEAK | 0.18 | 1063.57 | P006|BGAL_ECOLI |
| TPHPALTEAK | 0.18 | 1063.57 | P006|BGAL_ECOLI |
| - | - | - | - |
| TPHPALTEAK | 0.18 | 1063.57 | P006|BGAL_ECOLI |

Collate Abundances

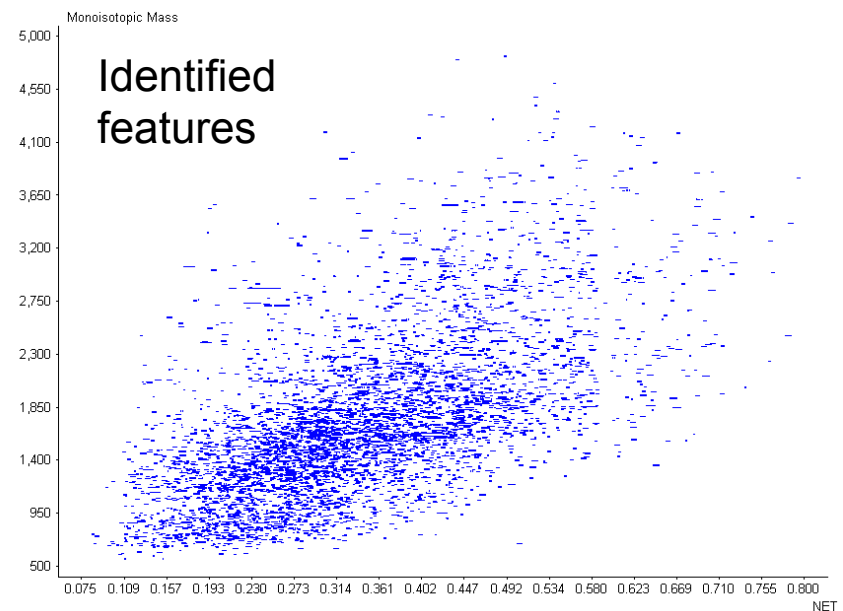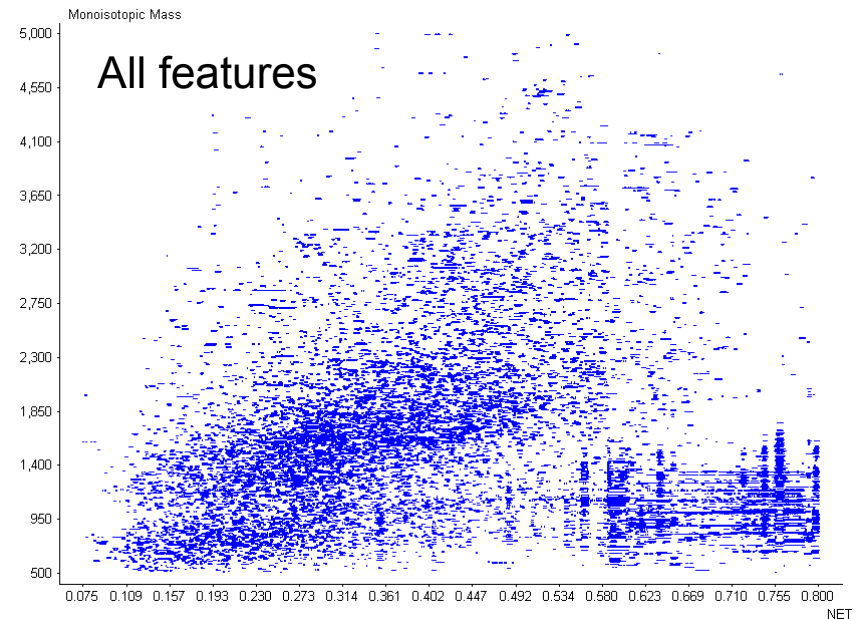| Peptide | NET | Mass | ORFName | Exp 1 | Exp 2 | Exp 3 | Exp i | Exp 1600 |
|---|---|---|---|---|---|---|---|---|
| TPHPALTEAK | 0.18 | 1063.57 | P006|BGAL_ECOLI | 3320000 | 3524300 | - | 381000 | 381000 |

# Current AMT Tag Pipeline

● LC-MS features without matches may represent useful information, but are effectively ignored
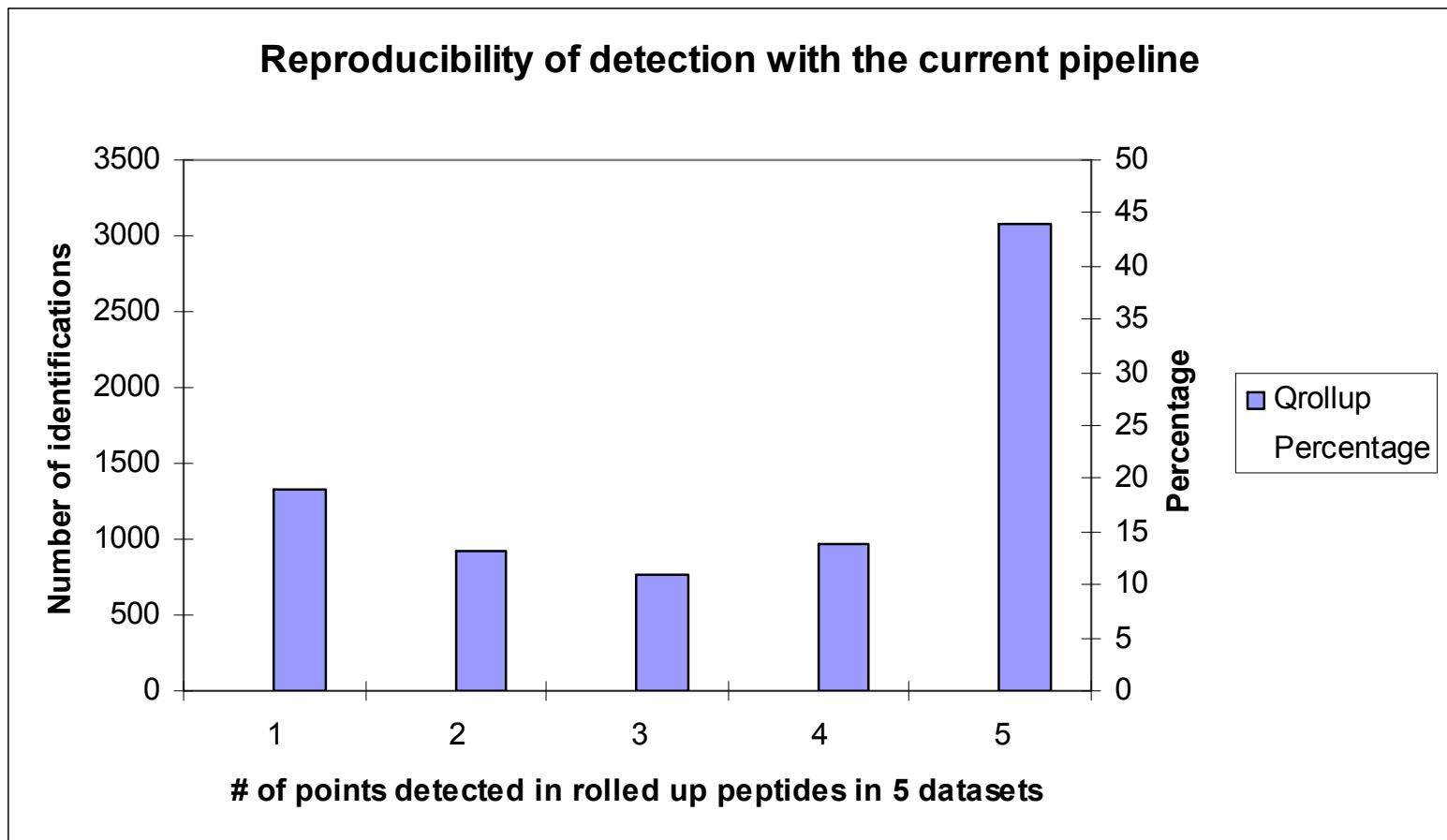
AMT tags from LC-MS/MS



LC-FTICR-MS
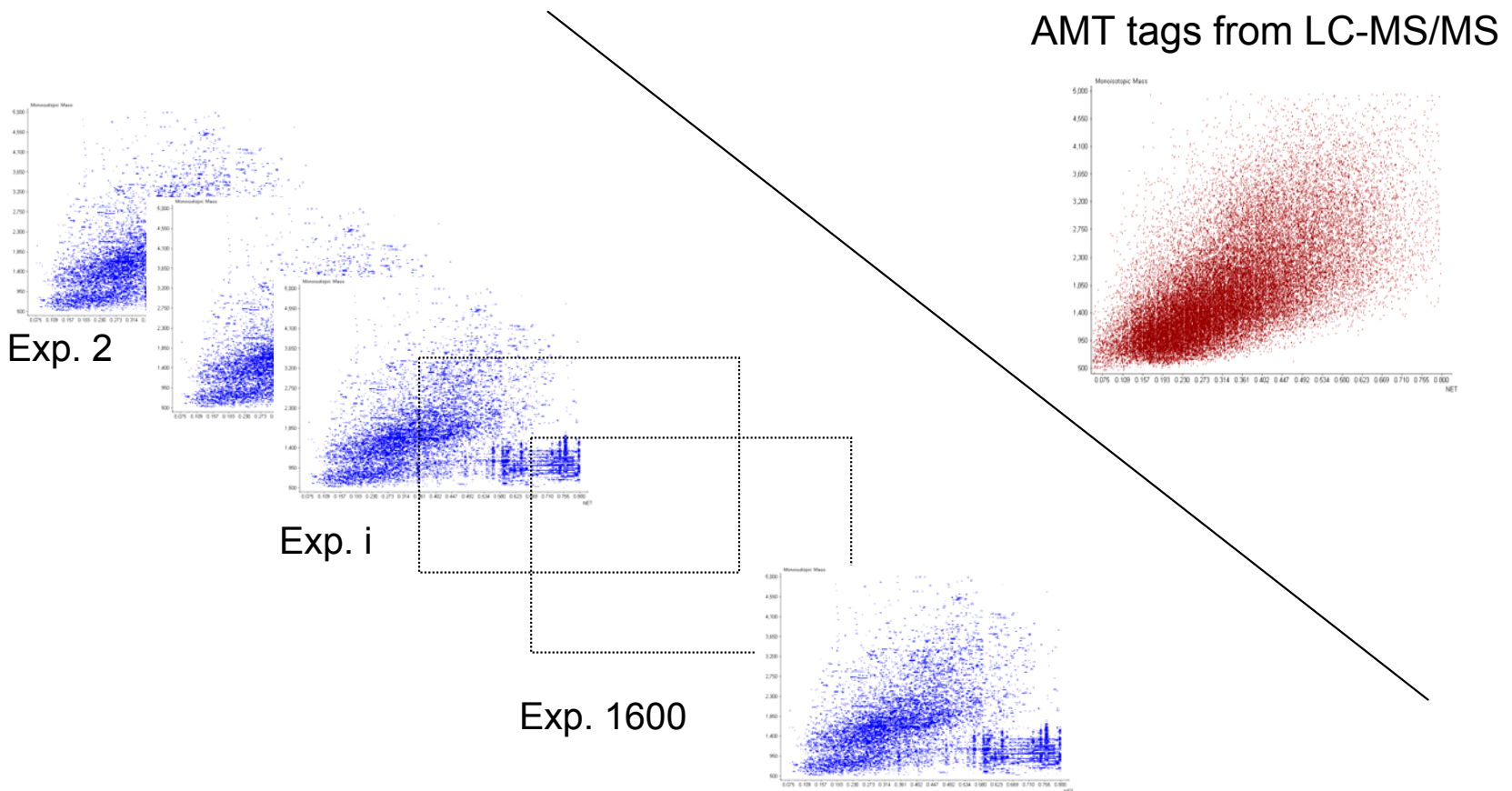
All features



Identified features

# Other issues

- Independent processing of each dataset results in more missing data, because of the lack of statistics
- Lower abundance features suffer more, but are not the only casualties

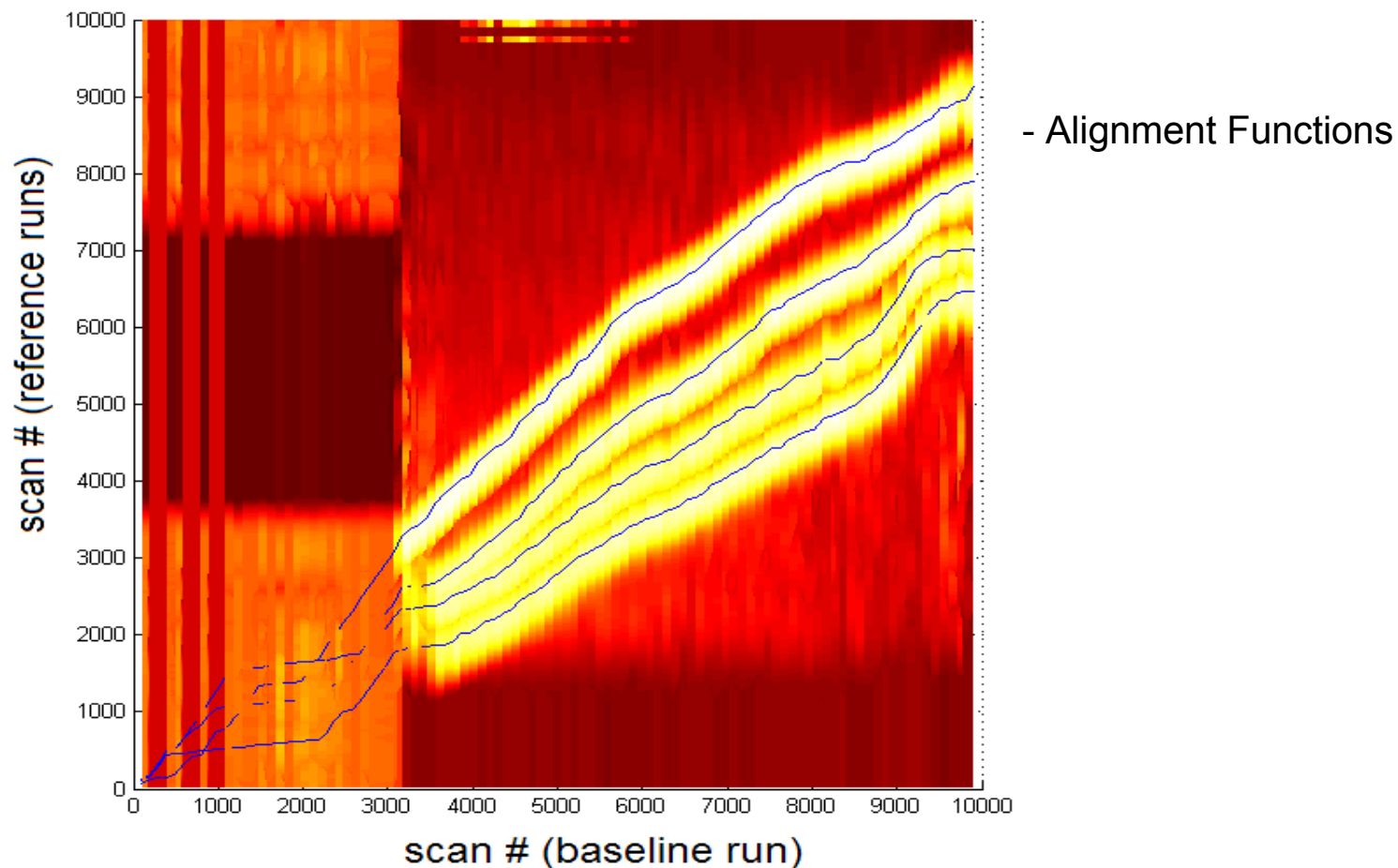**Reproducibility of detection with the current pipeline**

# Extended AMT Tag method

- Find common features based on mass and time patterns in all datasets first (with or without the AMT tag database)
- Align resulting groups of features to the AMT tag database using statistics from a larger number of features
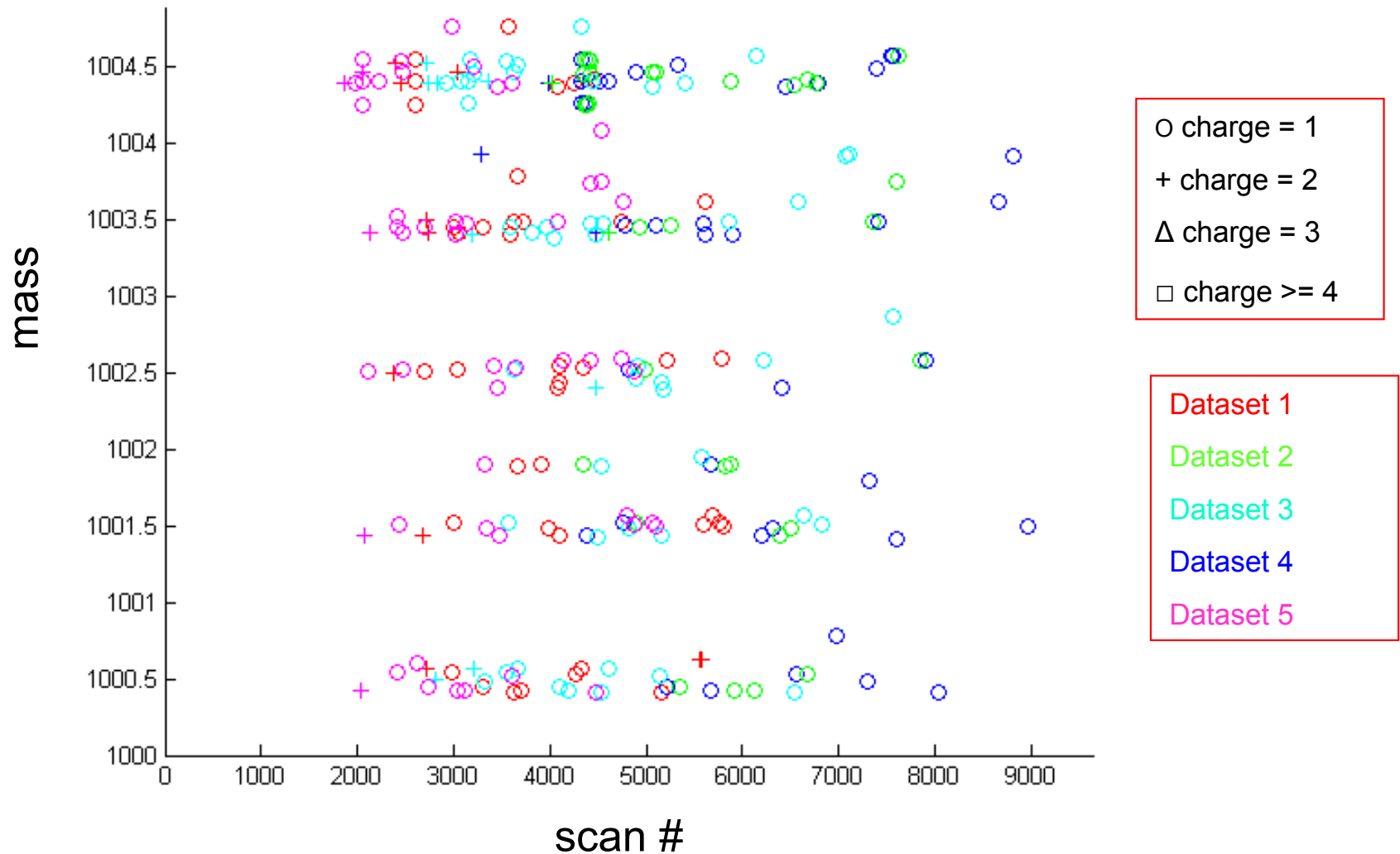


AMT tags from LC-MS/MS

LC-MS Exp. 1

Exp. 2

Exp. i

Exp. 1600

# Align all datasets to common baseline



- Alignment Functions

Score plots for alignment of 4 datasets against arbitrary baseline run

N. Jaitly, M.E. Monroe et. al., *Analytical Chemistry* **2006**, *78*, 7397-7409.
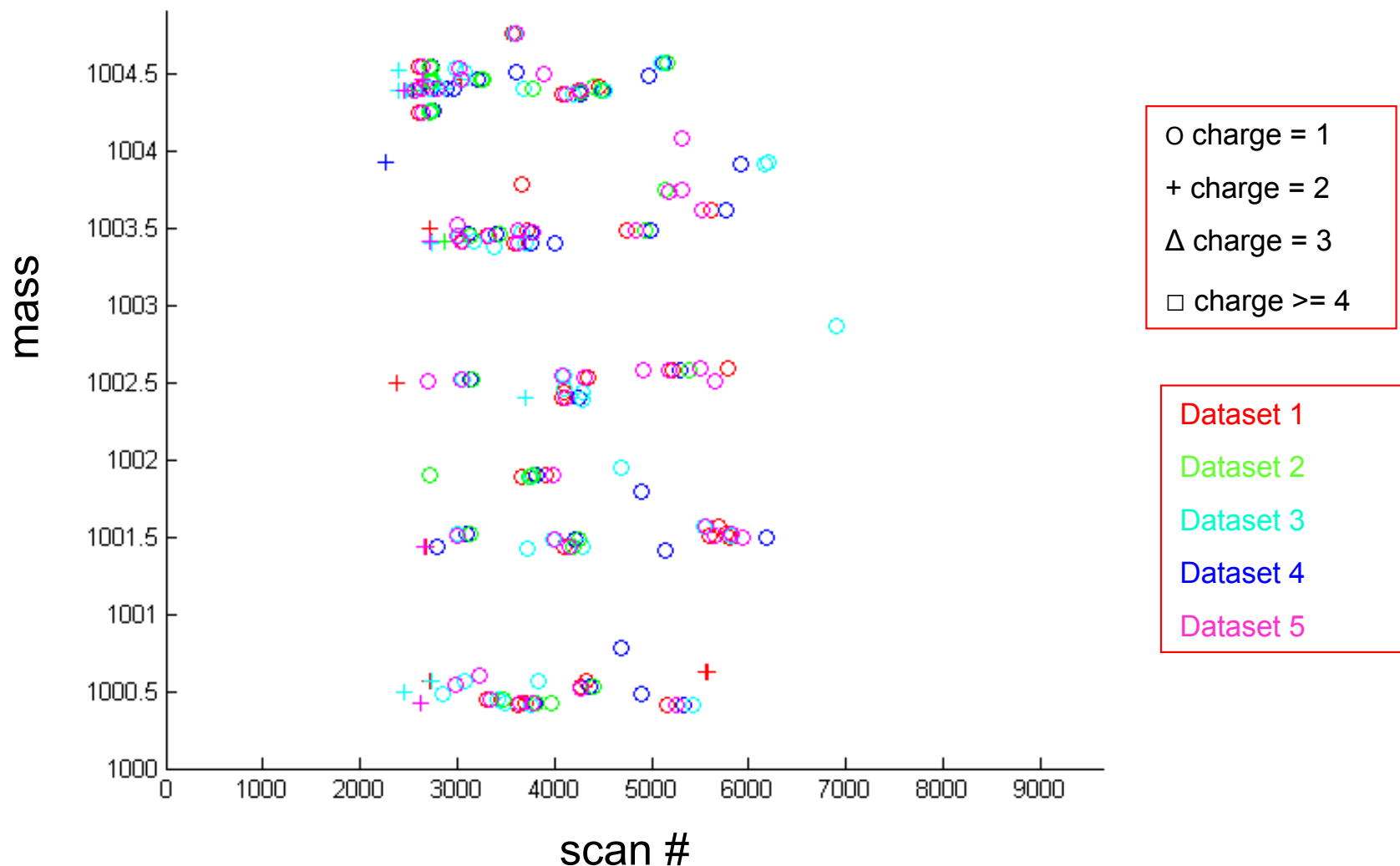
# Alignment of Multiple LC-MS Datasets

- Obvious need for alignment before finding common features
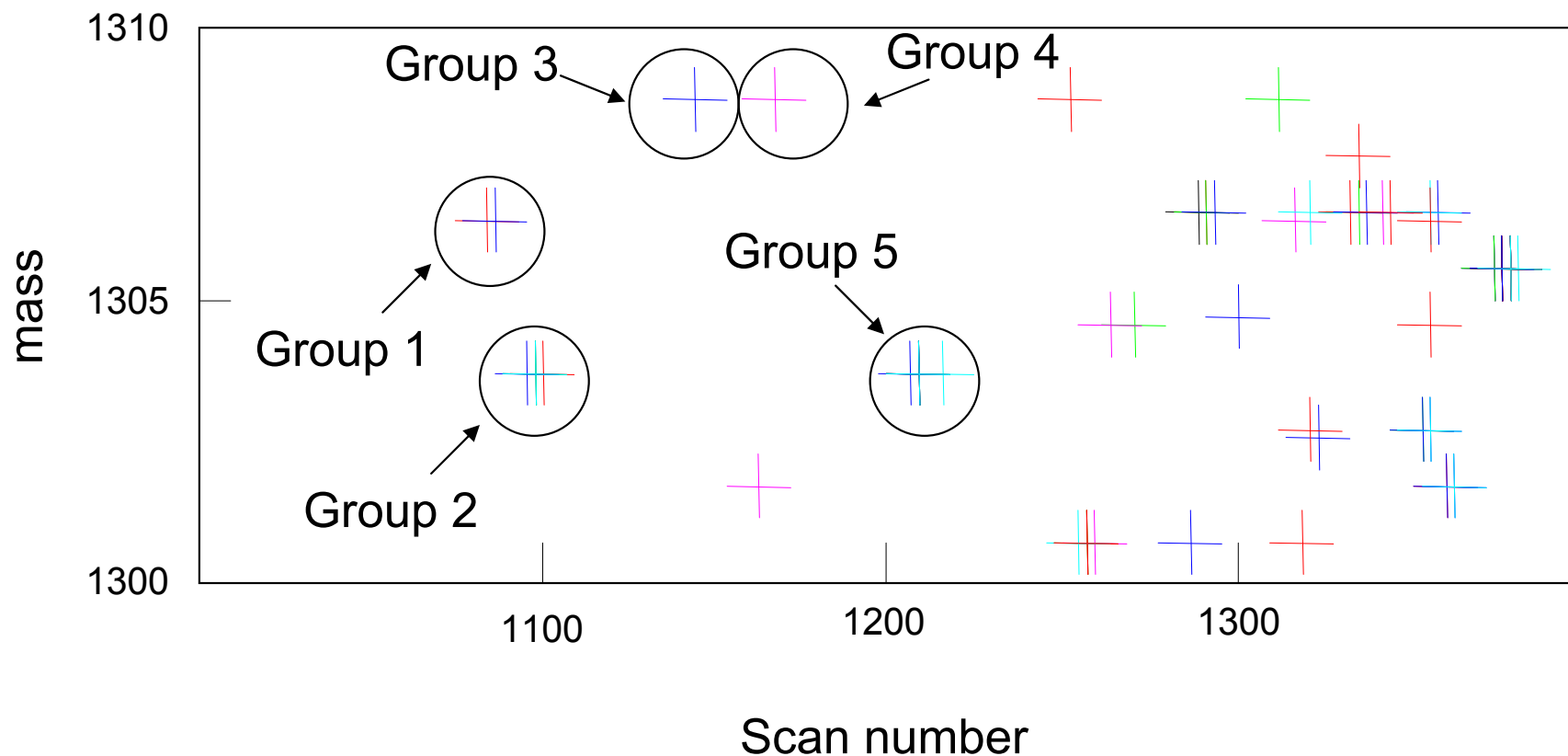  - Mass section of 5 LC-MS datasets before LC alignment

# Alignment of Multiple LC-MS Datasets

- Obvious need for alignment before finding common features
  - Mass section of 5 LC-MS datasets after LC alignment



○ charge = 1
+ charge = 2
Δ charge = 3
□ charge >= 4

Dataset 1
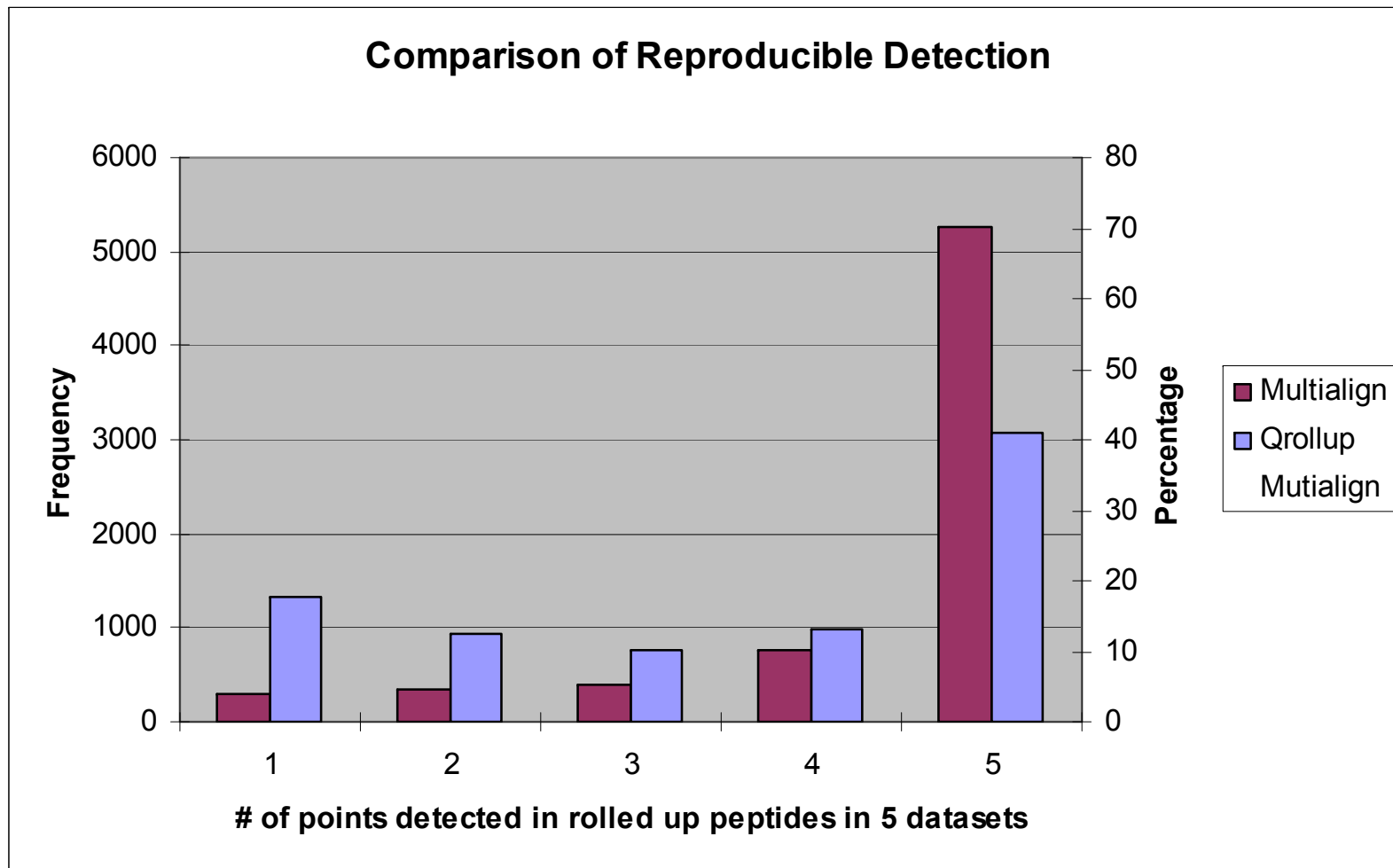Dataset 2
Dataset 3
Dataset 4
Dataset 5

# Clustering Features

- Create abundance profiles by finding similar features (using mass and retention time) across all LC-MS datasets, rather than analyzing each dataset separately and then collating results
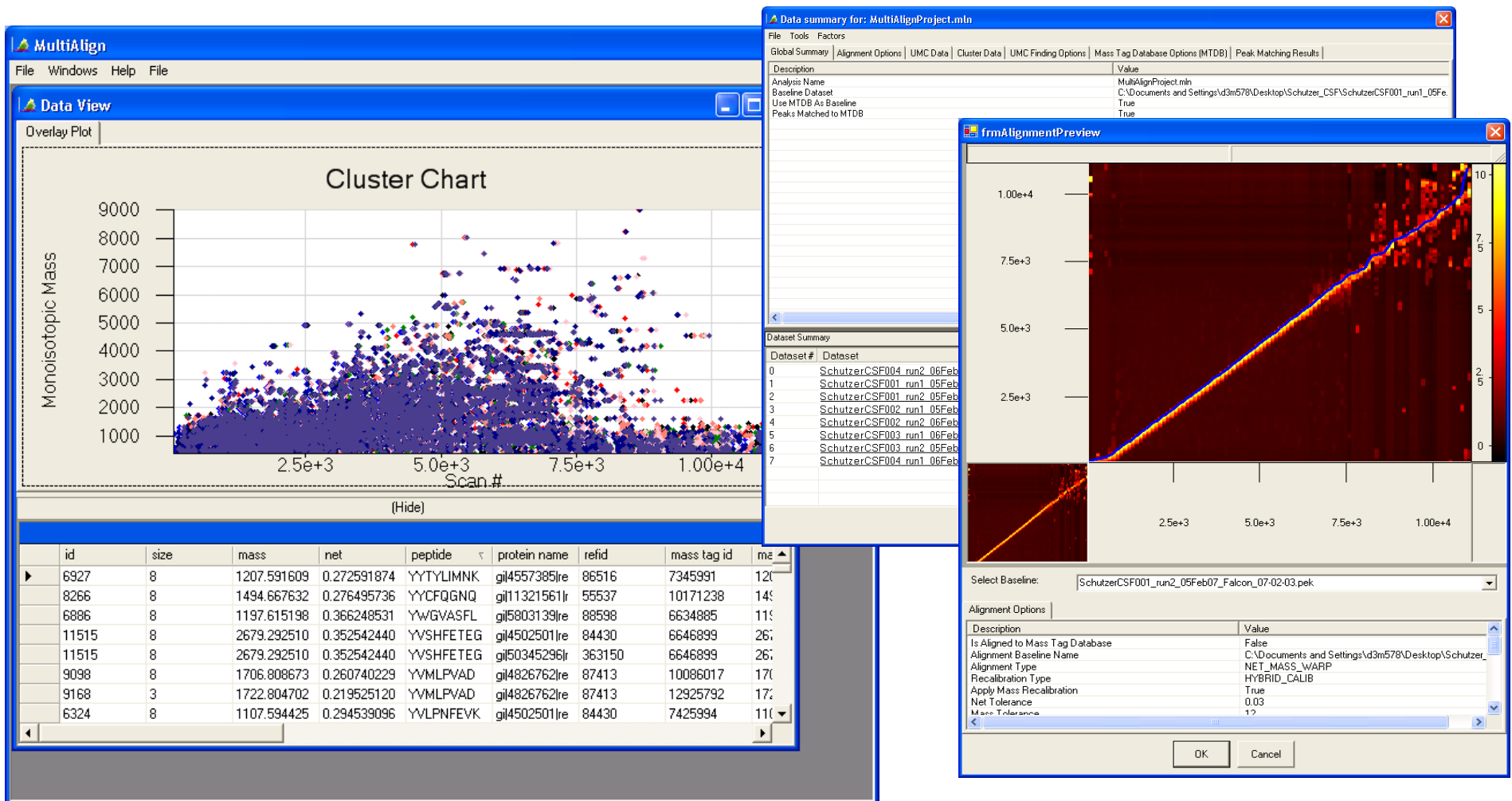
# Identifying Clustered Features

● Align mass and elution time of clusters to AMT tag database, then identify clusters by matching to AMT tags



**Comparison of Reproducible Detection**

Legend: Multialign, Qrollup, Mutialign

Y-axis (left): Frequency — 0, 1000, 2000, 3000, 4000, 5000, 6000

Y-axis (right): Percentage — 0, 10, 20, 30, 40, 50, 60, 70, 80

X-axis: # of points detected in rolled up peptides in 5 datasets — 1, 2, 3, 4, 5

Fewer missing values observed with clustered feature approach

# MultiAlign

- Represents next version of the feature identification process
- Along with MTDB Creator it represents a standalone, redistributable version of the AMT tag process

# LC-MS Feature Discovery

- Similar approaches and software tools: High Res LC-MS
  - CRAWDAD
    - G.L. Finney et al. *Analytical Chemistry* **2008**, *80*, 961-971.
  - msInspect
    - M. Bellew et. al. *Bioinformatics* **2006**, *22*, 1902-1909.
  - PEPPeR
    - J. Jaffe et.al. *Mol. Cell. Proteomics* **2006**, *5*, 1927-1941.
  - SpecArray (Pep3D, mzXML2dat, PepList, PepMatch, PepArray)
    - X.-J. Li, et. al. *Mol Cell Proteomics* **2005**, *4*, 1328-1340.
  - SuperHIRN
    - L.N. Mueller et al. *Proteomics* **2007**, *7*, 3470-3480.
  - Surromed label-free quantitation software (MassView)
    - W. Wang et al. *Analytical Chemistry* **2003**, *75*, 4818-4826.
  - XCMS (for Metabolite profiling)
    - C.A. Smith et. al. *Analytical Chemistry* **2006**, *78*, 779-787.
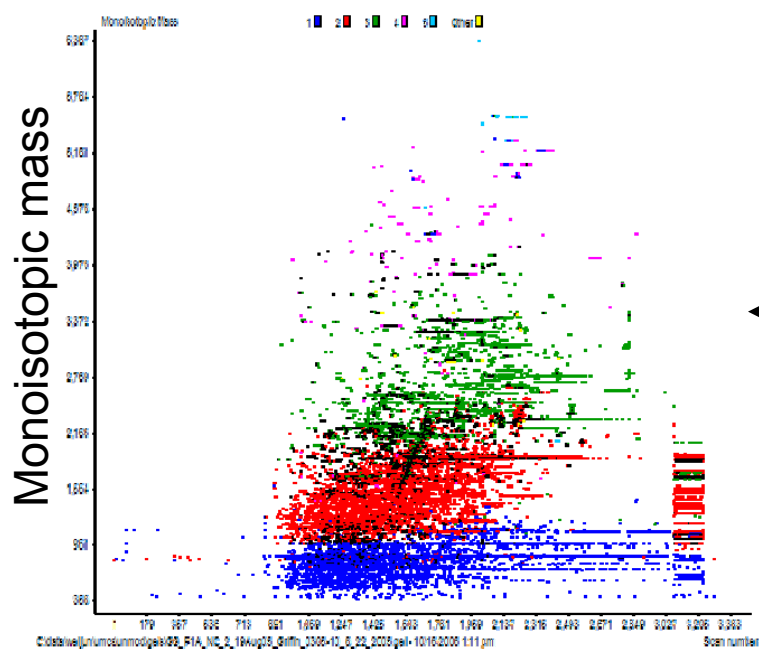
# LC-MS Feature Discovery

- **Similar approaches and software tools: Low Res LC-MS**
  - Signal maps software
    - A. Prakash et. al. *Mol. Cell Proteomics* **2006**, *5*, 423-432.
  - Informatics platform for global proteomic profiling using LC-MS
    - D. Radulovic, et al. *Mol. Cell. Proteomics* **2004**, *3*, 984-997.
  - Computational Proteomics Analysis System (CPAS)
    - A. Rauch et. al. *J. Proteome Research* **2006,** *5*, 112-121.

# Part II: LC-MS Feature Discovery

- Introduction (Adkins)
- Part I: Overview of Label-Free Quantitative Proteomics (Jaffe)
- Part II: Feature discovery in LC-MS datasets (Monroe and Jaitly)
  - ✔ Structure of LC-MS Data
  - ✔ Feature discovery in individual spectra (deisotoping)
  - ✔ Feature definition over elution time
  - ✔ Identifying LC-MS Features using an AMT tag DB
  - ✔ Extending the AMT tag approach for feature based analyses
  - Estimating confidence of identified LC-MS features
  - Downstream quantitative analysis with DAnTE
- Part III: PEPPeR, GenePattern and Real-world examples (Jaffe)
- Break
- AMT tag Pipeline Demo (general)
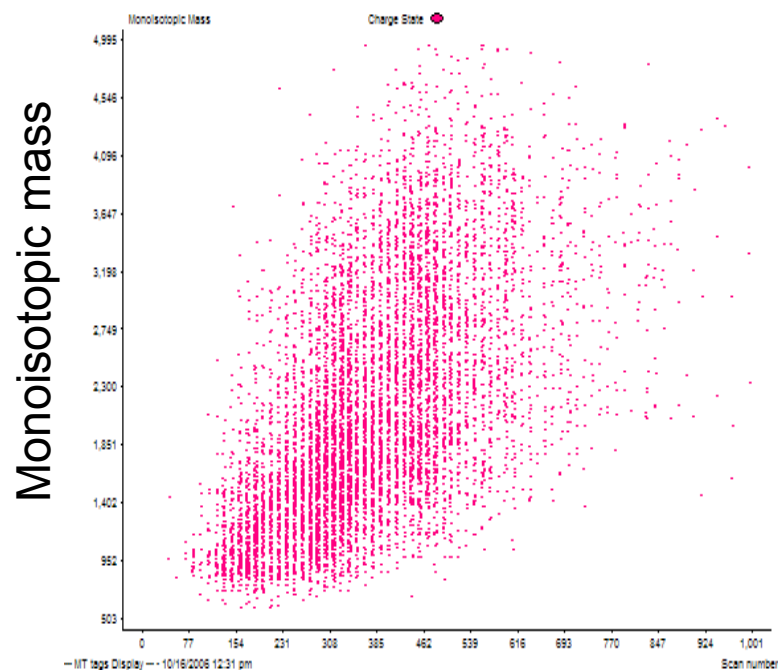- Panel Discussion

# Developing Confidence Metrics

- LC-MS data is aligned against an AMT tag database
- Each LC-MS feature is matched to the closest AMT tag in mass and normalized LC elution time (NET) dimensions



Alignment & Peak Matching

scan #

LC-MS dataset

Normalized elution time (NET)
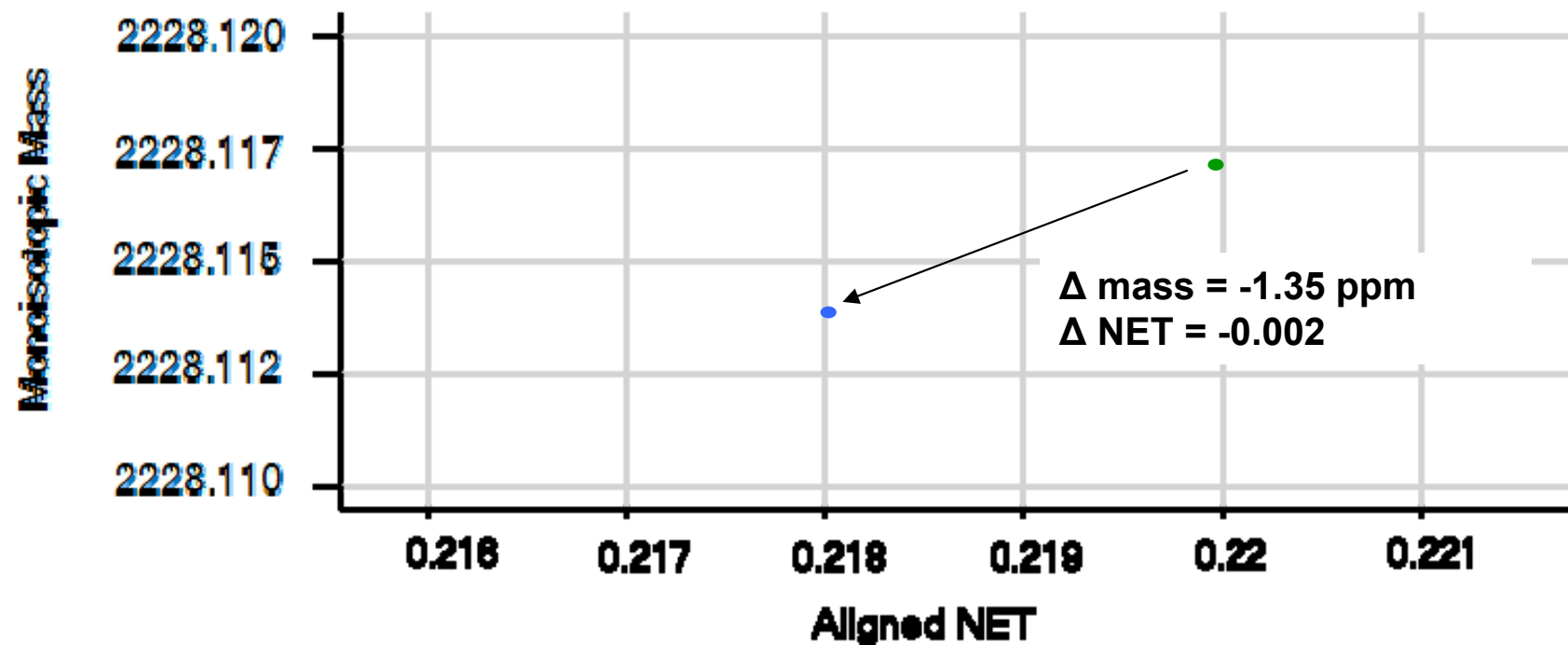
AMT tag Database

**How do we control the errors in the process?**

# Controlling Rate of Random Matches

- Size of database and degree of noise in database affects the rate of random matches

  - Building more confident AMT tag database (e.g., using strict filtering) decreases background false positives

  - But, increases false negatives

- To date, *ad hoc* rules have been used

  - *Subjectively pleasing* threshold values selected for different parameters, such as mass error tolerance, LC NET error tolerance, etc.

  - False discovery rate (FDR) was estimated using decoy methods

  - Rules were accepted if results seemed satisfactory, otherwise parameters were re-optimized

  - But, chosen parameters may not result in optimal results

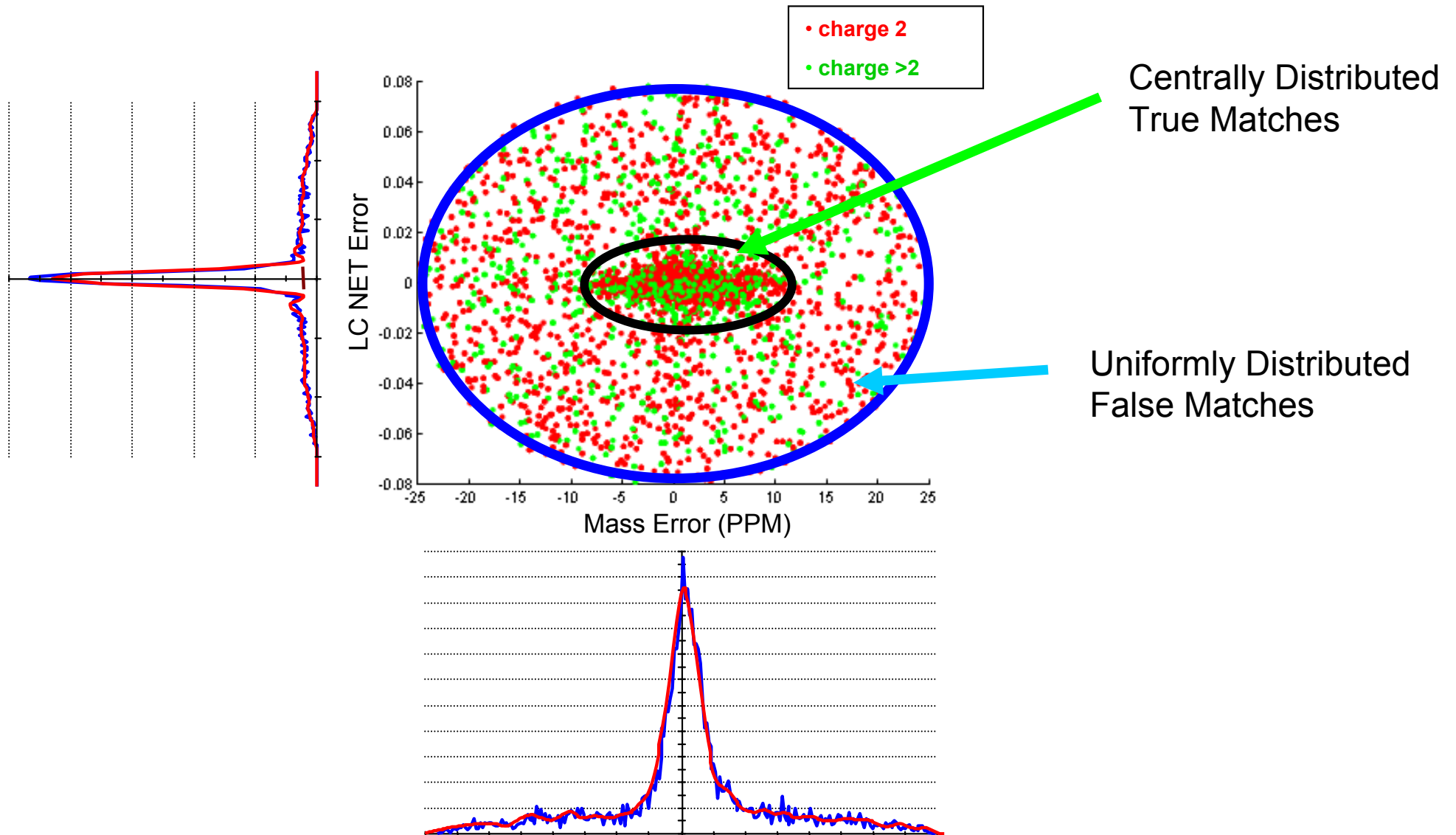# Metrics Associated with a Candidate Identification

- Each match between an LC-MS feature and a peptide AMT tag is described by a mass error and an LC NET error

| Mass | Scan | Aligned NET | Peptide | NET | Mass | ORFName |
|------|------|-------------|---------|-----|------|---------|
| 2228.114 | 1097 | 0.218 | TETQEKNPLPSKETIEQEK | 0.22 | 2228.117 | Thymosin beta-4 |



Δ mass = -1.35 ppm
Δ NET = -0.002

# Distribution of Peak Matches

- True and false matches resulting from peak matching display different mass and LC NET error distributions
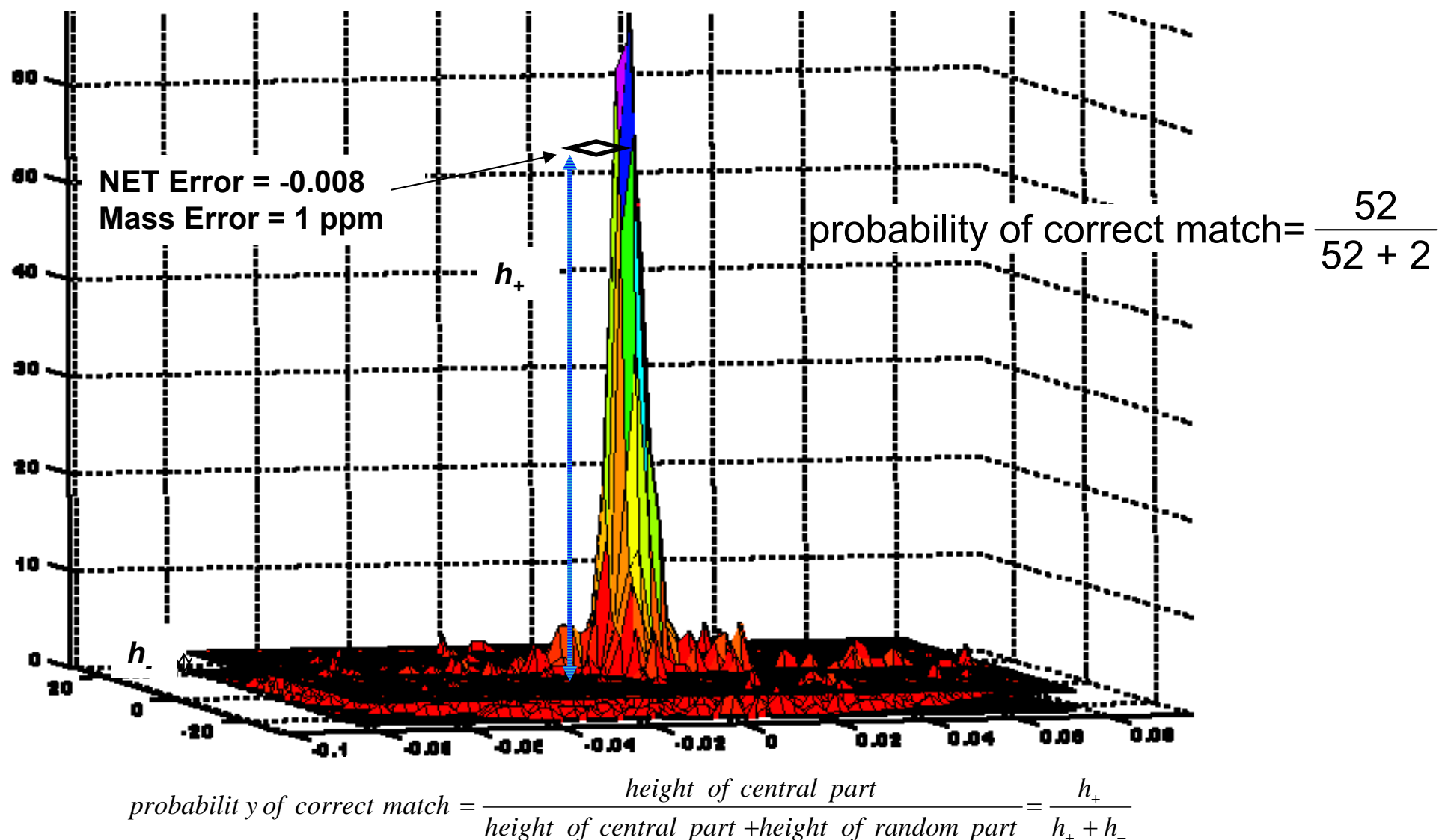
# Distribution of LC-MS Peak Matches

- Density plot of mass and LC NET error distributions is a sum of true and false components
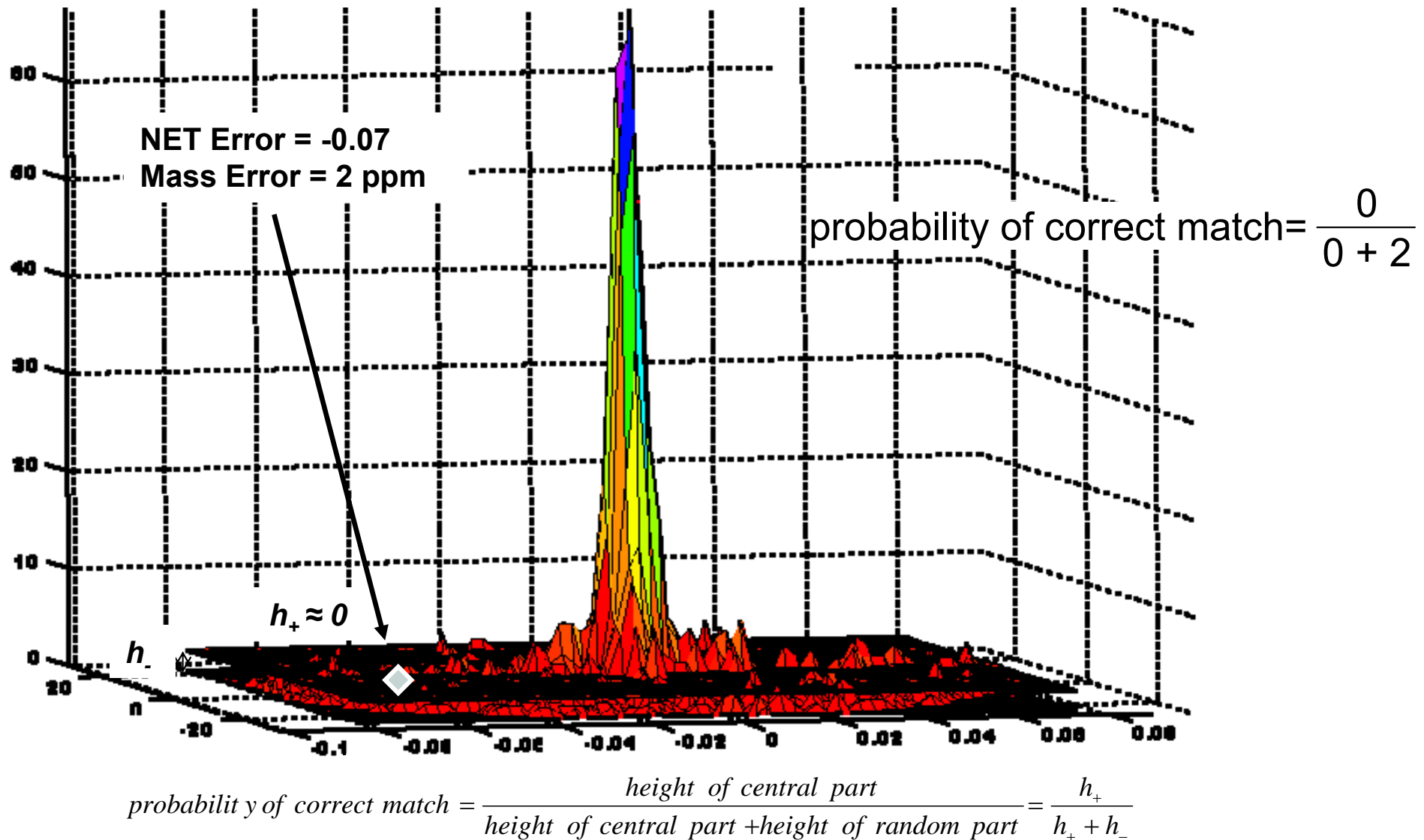
# Estimating the Probability a Match is Correct

- Approach: the probability that a peak match is correct can be estimated from where its mass and LC NET error values lie on the two-dimensional distribution



**NET Error = -0.008**
**Mass Error = 1 ppm**

$h_+$

$h_-$

$$\text{probability of correct match} = \frac{52}{52 + 2}$$

$$probabilit\,y\,of\,correct\,match = \frac{height\ of\ central\ part}{height\ of\ central\ part + height\ of\ random\ part} = \frac{h_+}{h_+ + h_-}$$

# Estimating the Probability a Match is Correct

● The probability that a peak match is correct depends on where its mass and LC NET error value lies on the two-dimensional distribution



NET Error = -0.07
Mass Error = 2 ppm

probability of correct match $= \dfrac{0}{0 + 2}$

$h_+ \approx 0$

$h_-$

$$probability \ of \ correct \ match = \frac{height \ of \ central \ part}{height \ of \ central \ part + height \ of \ random \ part} = \frac{h_+}{h_+ + h_-}$$

# Optimizing the overall matching process

- General approach; calculate confidence in peak match based on:
  - Mass and LC NET errors
  - Instrumental performance for an analysis
    - Mass error precision
    - LC-NET precision
  - LC-MS/MS ID quality (e.g., SEQUEST XCorr or X!Tandem expectation values)

- Inter-related effects of different parameters on each other complicate simple choices:
  - Lower mass and LC NET errors should allow choice of lower scores
  - Higher scores should allow somewhat wider mass and LC NET tolerances

→ **For practical value we need a single metric that calculates and combines all these factors automatically**
  - *Statistical Method for Assignment of Relative Truth (SMART)* – More details to be presented at ASMS 2008 Bioinformatics oral session

# Part II: LC-MS Feature Discovery

- Introduction (Adkins)
- Part I: Overview of Label-Free Quantitative Proteomics (Jaffe)
- Part II: Feature discovery in LC-MS datasets (Monroe and Jaitly)
  - ✓ Structure of LC-MS Data
  - ✓ Feature discovery in individual spectra (deisotoping)
  - ✓ Feature definition over elution time
  - ✓ Identifying LC-MS Features using an AMT tag DB
  - ✓ Extending the AMT tag approach for feature based analyses
  - ✓ Estimating confidence of identified LC-MS features
  - Downstream quantitative analysis with DAnTE
- Part III: PEPPeR, GenePattern and Real-world examples (Jaffe)
- Break
- AMT tag Pipeline Demo (general)
- Panel Discussion

# Downstream Data Analysis

- Quantitative protein inference from peptide data
- Complications
  - Multiple, possibly inconsistent peptide measurements for same protein
  - Systematic abundance variation within and between conditions
    - How should we use information from blocking and randomization of experiments?
  - High rate of missingness in peptide measurements
- Need to combine off the shelf statistical methods and novel solutions
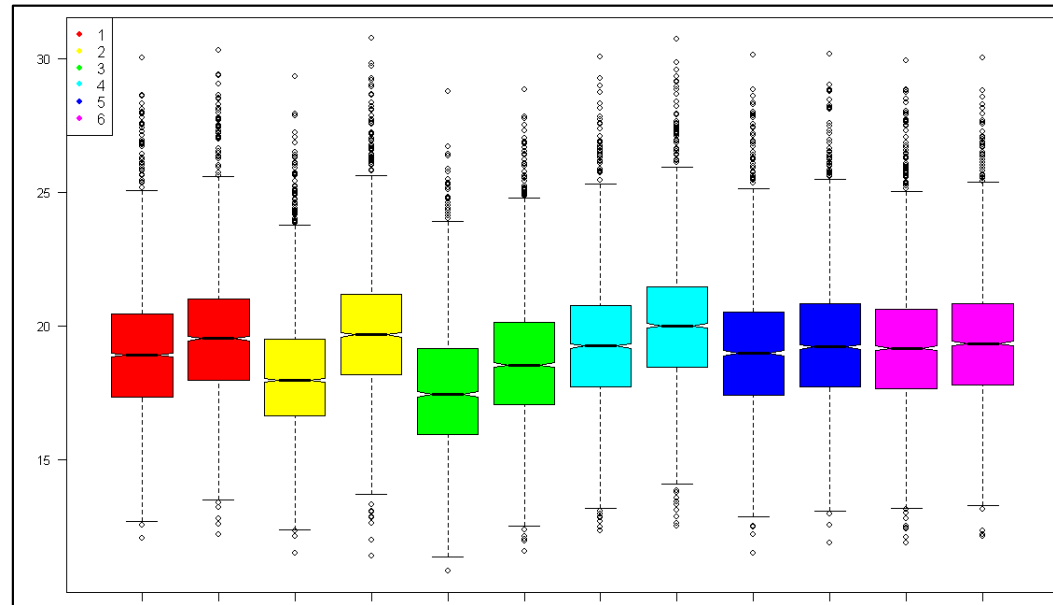  - Clustering
  - ANOVA
  - PCA

# Infer Protein Abundances from Peptide Abundances

- Multiple peptides observed for each protein
  - For example, protein with 4 peptides
    1. **SADLNVDSIISYWK**
    2. **LLLTSTGAGIIDVIK**
    3. **LIVGFPAYGHTFILSDPSK**
    4. **IPELSQSLDYIQVMTYDLHDPK**

  - Plot peptide abundance across 57 datasets (for 4 conditions)



→ Outlier detection and normalization need to be performed before meaningful abundance information can be inferred

# Outlier Detection

# Normalization

# Infer Protein Abundances from Peptide Abundances

- Scale peptide abundances to an automatically chosen "optimal" reference peptide for each protein

- Estimate relative protein abundance using scaled peptides



Raw peptide abundances vs. dataset (for 1 protein)

**Condition 1**   **Condition 2**   **Control**   **Condition 3**

Scaled peptide abundances for this protein's 4 peptides

Datasets

Median protein abundance (dark black line)

1. SADLNVDSIISYWK
2. LLLTSTGAGIIDVIK
3. LIVGFPAYGHTFILSDPSK
4. IPELSQSLDYIQVMTYDLHDPK

# <u>D</u>ata <u>A</u>nalysis <u>T</u>ool <u>E</u>xtension (DAnTE)
## A software tool for downstream quantitative protein inference

AMT tag Pipeline (LC-MS)

MultiAlign

VIPER

Tabular Data
(raw abundances, exp. ratios, spectral counts)

DAnTE

Proteins, peptides and peptide abundances

*Software by Ashoka Polpitiya*

| Ref ID | Reference | Mass Tag ID | AB_007-LP | AB_007-LP | AB_007-LP | AB-008-LP | AB-008-LP |
|--------|-----------|-------------|-----------|-----------|-----------|-----------|-----------|
| 57559 | gi\|10048460\|ref\|NP | 32275105 | | 0.0133541 | | | |
| 57559 | gi\|10048460\|ref\|NP | 47465602 | 0.214036 | 0.110191 | 0.088877 | 0.144964 | |
| 57559 | gi\|10048460\|ref\|NP | 83201030 | | | | | |
| 57559 | gi\|10048460\|ref\|NP | 83257895 | | | | | |
| 57559 | gi\|10048460\|ref\|NP | 83424583 | 0.126785 | 0.0901864 | 0.05985 | 0.0391025 | 0.0275993 |
| 57559 | gi\|10048460\|ref\|NP | 83451097 | 0.0600628 | | | 0.0392326 | 0.0274422 |
| 57559 | gi\|10048460\|ref\|NP | 83451190 | 0.0409746 | 0.0396703 | 0.0539142 | 0.0591157 | 0.0246261 |
| 57559 | gi\|10048460\|ref\|NP | 83479064 | 0.0383742 | 0.0865723 | 0.0402673 | 0.0359534 | 0.0397751 |
| 57559 | gi\|10048460\|ref\|NP | 83479231 | 0.0212362 | 0.0122634 | 0.0184865 | | |
| 57559 | gi\|10048460\|ref\|NP | 83514326 | | 0.0289326 | 0.0337594 | 0.0237968 | 0.0243036 |
| 74732 | gi\|10092608\|ref\|NP | 7824413 | | | 0.0851353 | | |
| 74732 | gi\|10092608\|ref\|NP | 8680795 | | | | | |
| 74732 | gi\|10092608\|ref\|NP | 20746162 | 0.536631 | 0.315447 | 0.231846 | 0.131434 | 0.130482 |
| 74732 | gi\|10092608\|ref\|NP | 20750908 | 0.0699752 | 0.0821782 | 0.0417777 | 0.0280187 | 0.0366826 |
| 74732 | gi\|10092608\|ref\|NP | 20750955 | 0.16129 | 0.123591 | 0.12117 | 0.126886 | 0.0970479 |
| 74732 | gi\|10092608\|ref\|NP | 20956112 | 0.0407675 | | | | |
| 74732 | gi\|10092608\|ref\|NP | 20985367 | | 0.0144197 | 0.0138798 | | 0.0125844 |
| 60259 | gi\|10181140\|ref\|NP | 7777112 | | 0.0194588 | | | |

# Interactive Analysis in DAnTE

**Data Loading**

- Peptide abundance
- Peptide-Protein relations
- Factors

**Variance Stabilization**

- log2 or log10
- Bias (additive/multiplicative)

**Investigative Plots**

- Histograms
- Boxplots
- Correlation diagrams
- MA Plots

**Replicate Normalization**

- Linear Regression
- Local regression (LOESS)
- Quantile

**Global Normalization**

- Central tendency
- Median absolute Deviation (MAD)

**Infer Proteins from Peptides**

- RRollup
- ZRollup
- QRollup
- Rollup Plots

**Statistical Tests**

- ANOVA
- Mix Models

**Visualization**

- PCA
- PLS
- Heatmaps (hierachical, kmeans)

**Impute Missing Data**

- Substitute
- Average
- KNNimpute
- SVDimpute etc.

**Other Features**

- Filter ANOVA results
- Save session

# Outline of a Typical Analysis

- Load data
- Examine diagnostic plots
- Define factors
- Normalize
  - Within a Factor
    - Linear regression
    - LOESS (LOcal regrESSion)
    - Quantile
  - Across Factors
    - MAD
    - Central tendency
- Infer protein abundances from peptide abundances
  - RRollup, QRollup, and ZRollup
- ANOVA
- Save the results to a session file (.dnt)

# Load Data

### Tabular Data File:

Proteins, peptides, and peptide abundances

| Ref ID | Reference | Mass Tag ID | AB_007-LP | AB_007-LP | AB_007-LP | AB-008-LP | AB-008-LP |
|---|---|---|---|---|---|---|---|
| 57559 | gi\|10048460\|ref\|NP | 32275105 | | 0.0133541 | | | |
| 57559 | gi\|10048460\|ref\|NP | 47465602 | 0.214036 | 0.110191 | 0.088877 | 0.144964 | |
| 57559 | gi\|10048460\|ref\|NP | 83201030 | | | | | |
| 57559 | gi\|10048460\|ref\|NP | 83257895 | | | | | |
| 57559 | gi\|10048460\|ref\|NP | 83424583 | 0.126785 | 0.0901864 | 0.05985 | 0.0391025 | 0.0275993 |
| 57559 | gi\|10048460\|ref\|NP | 83451097 | 0.0600628 | | | 0.0392326 | 0.0274422 |
| 57559 | gi\|10048460\|ref\|NP | 83451190 | 0.0409746 | 0.0396703 | 0.0539142 | 0.0591157 | 0.0246261 |
| 57559 | gi\|10048460\|ref\|NP | 83479064 | 0.0383742 | 0.0865723 | 0.0402673 | 0.0359534 | 0.0397751 |
| 57559 | gi\|10048460\|ref\|NP | 83479231 | 0.0212362 | 0.0122634 | 0.0184865 | | |
| 57559 | gi\|10048460\|ref\|NP | 83514326 | | 0.0289326 | 0.0337594 | 0.0237968 | 0.0243036 |
| 74732 | gi\|10092608\|ref\|NP | 7824413 | | | 0.0851353 | | |
| 74732 | gi\|10092608\|ref\|NP | 8680795 | | | | | |
| 74732 | gi\|10092608\|ref\|NP | 20746162 | 0.536631 | 0.315447 | 0.231846 | 0.131434 | 0.130482 |
| 74732 | gi\|10092608\|ref\|NP | 20750908 | 0.0699752 | 0.0821782 | 0.0417777 | 0.0280187 | 0.0366826 |
| 74732 | gi\|10092608\|ref\|NP | 20750955 | 0.16129 | 0.123591 | 0.12117 | 0.126886 | 0.0970479 |
| 74732 | gi\|10092608\|ref\|NP | 20956112 | 0.0407675 | | | | |
| 74732 | gi\|10092608\|ref\|NP | 20985367 | | 0.0144197 | 0.0138798 | | 0.0125844 |
| 60259 | gi\|10181140\|ref\|NP | 7777112 | | 0.0194588 | | | |

# Diagnostic Plots: Check Normality



**Select QQ Plot Parameters**

**QQ Plots**

Data Source: **Log Expressions**

Plot Properties

Columns on the Multi-Plot: 2
Symbol Foreground Color:
Symbol Border Color:
Line Color:
Transparent Background: ☐ (Only works with PNG format)

Reference Distribution

⦿ Normal
○ Exponential   rate: 1.0
○ Student   df: 4
○ Weibull   Shape: 2.0   Scale: 1.0

Select Datasets to Plot

☐ NIOSH_AB_007-LP       ☐ NIOSH
☐ NIOSH_AB_007-LP1      ☐ NIOSH
☐ NIOSH_AB_007-LP2      ☐ NIOSH
☐ NIOSH-AB-008-LP       ☐ NIOSH
☐ NIOSH-AB-008-LP1      ☐ NIOSH
☐ NIOSH_AB-008-LP2      ☐ NIOSH
☐ NIOSH_AB-009-LP       ☐ NIOSH
☐ NIOSH_AB-009-LP1      ☐ NIOSH
☐ NIOSH_AB-010-LP       ☐ NIOSH
☐ NIOSH_AB-010-LP1      ☐ NIOSH
☐ NIOSH_AB-011-LP       ☐ NIOSH
☐ NIOSH_AB-011-LP1      ☐ NIOSH
☐ NIOSH-AB-012-LP       ☐ NIOSH

Toggle All

**Select Histogram Plot Parameters**

**Histograms**

Data Source: **Log Expressions**

Plot Properties

Manually set the Bins: 50    ☐ Auto Binning
Columns on the Multi-Plot: 2
Foreground Color:        Border Color:
Transparent Background: ☐ (Only works with PNG format)
Add Rug: ☑

Select Datasets to Plot

☐ NIOSH_AB_007-LP       ☐ NIOSH_AB-010-LP       ☐ NIOSH_CE
☐ NIOSH_AB_007-LP1      ☐ NIOSH-AB-010-LP1      ☐ NIOSH_CE
☐ NIOSH_AB_007-LP2      ☐ NIOSH_AB-011-LP       ☐ NIOSH_CE
☐ NIOSH_AB-008-LP       ☐ NIOSH_AB-011-LP1      ☐ NIOSH_CE
☐ NIOSH_AB-008-LP1      ☐ NIOSH_AB-012-LP       ☐ NIOSH_CE
☐ NIOSH-AB-008-LP2      ☐ NIOSH_AB-012-LP1      ☐ NIOSH_CE
☐ NIOSH_AB-009-LP       ☐ NIOSH_CB_007-LP       ☐ NIOSH_CE
☐ NIOSH_AB-009-LP1      ☐ NIOSH_CB-007-LP1      ☐ NIOSH_CE

☐ Add Date/Name Stamp       Toggle All

OK    Defaults    Cancel

## Quantile-Quantile Plot



Sample Quantiles

Theoretical Quantiles (Normal)

## Histogram



Probability

Abundance

# Factors

- Capture experimental design through factors
  - For example, gender, sample type, technical replicate, and/or biological replicate

# Normalization: LOESS

# Protein Abundance Inference

- DAnTE currently has 3 different algorithms for rolling up peptide abundances to infer protein abundances
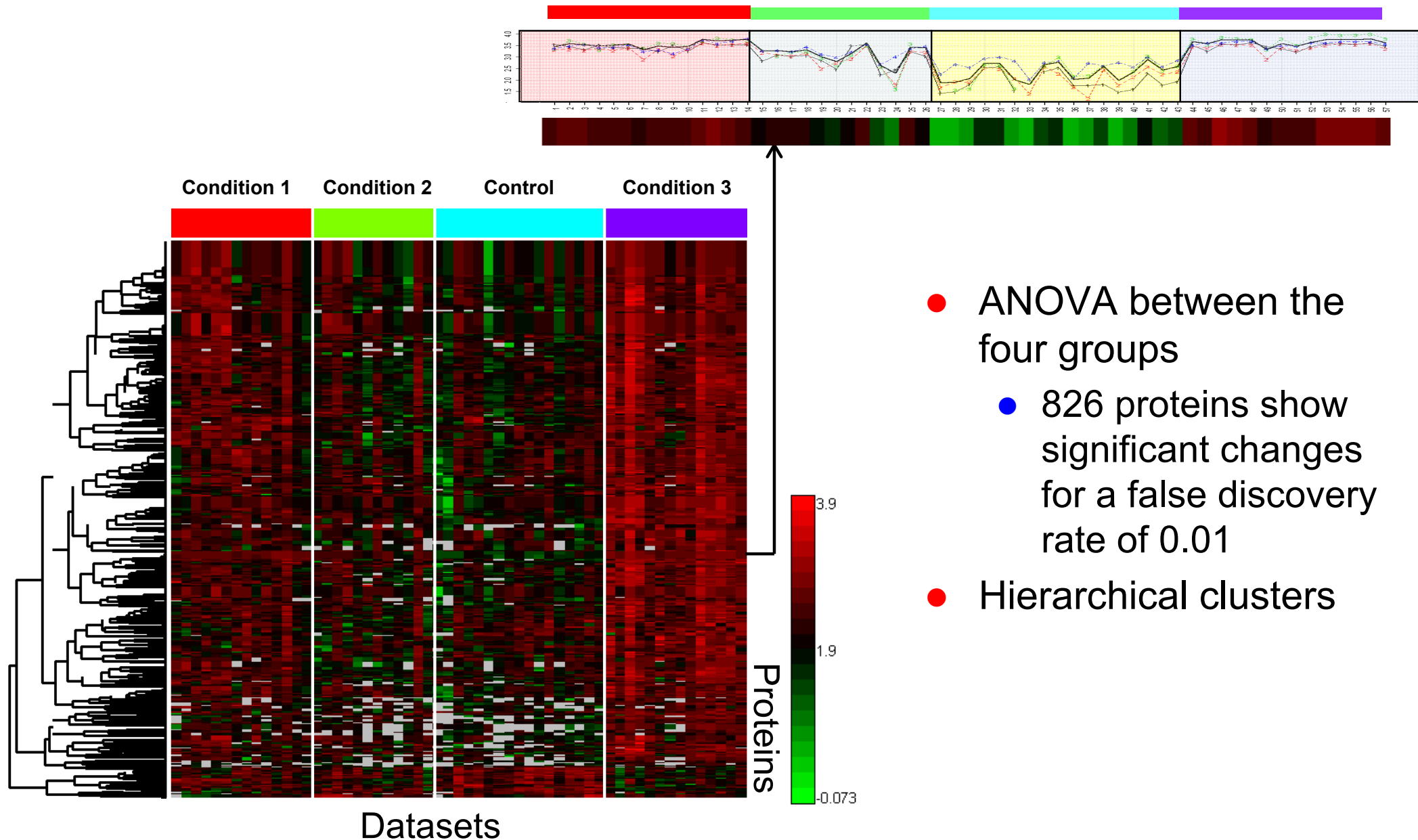- Additional algorithms can be added as needed

# Protein Heatmap

- All proteins
- Each row corresponds to a protein and each column to a dataset
- Color represents abundance (median ~2)

# Significant Proteins



- ANOVA between the four groups
  - 826 proteins show significant changes for a false discovery rate of 0.01
- Hierarchical clusters

# Complete DAnTE Feature List

- Data loading with peptide-protein group information
- Log transform
- Factor Definitions
- Normalization
  - Linear Regression
  - Loess
  - Quantile normalization
  - Median Absolute Deviation (MAD) Adj.
  - Mean Centering
- Missing Value Imputation
  - Simple
    - mean/median of the sample
    - Substitute a constant
  - Advance
    - Row mean within a factor
    - kNN method
    - SVDimpute
- Save tables / factors / session

- Plots
  - Histograms
  - QQ plots
  - Boxplots
  - Correlation plots
  - MA plots
  - PCA/PLS plots
  - Protein rollup plots
  - Heatmaps
- Rolling up to Proteins
  - Reference peptide based scaling (RRollup)
  - Z-score averaging (ZRollup)
  - QRollup
- Statistics
  - ANOVA
    - Provisions for unbalanced data
    - Random effects (multi level) models (REML)
  - Normality test (Shapiro-Wilks)
  - Non-parametric methods (Wilcoxon, Kruskal-Walis tests)
  - Q-values
  - Filters

# Course Outline

- Introduction (Adkins)
- Part I: Overview of Label-Free Quantitative Proteomics (Jaffe)
- Part II: Feature discovery in LC-MS datasets (Monroe and Jaitly)
- Part III: PEPPeR, GenePattern and Real-world examples (Jaffe)
  - PEPPeR: a self-contained web-based Biomarker Discovery pipeline
  - GenePattern: a suite of analysis and visualization tools that works with just about anything
- Break
- AMT tag Pipeline Demo (general)
- Panel Discussion
  - Questions
  - Future Directions

# Part III: PEPPeR, GenePattern and Real-world examples

## Jacob D. Jaffe

**The Broad Institute of Harvard and MIT**

**Proteomics Platform**

# Section Outline

- PEPPeR: a self-contained web-based Biomarker Discovery pipeline

- GenePattern: a suite of analysis and visualization tools that works with just about anything

- Examples of use in the real world
  - Proof of principle by accidental discovery of markers
  - In-silico defractionation
  - Breast cancer biomarker discovery

# PEPPeR:
**P**latform for **E**xperimental **P**roteomics **P**att**e**rn **R**ecognition

Jaffe JD, Mani DR, Leptos KC, Church GM, Gillette MA, Carr SA. PEPPeR, a Platform for Experimental Proteomic Pattern Recognition. *Mol Cell Proteomics*. 2006 Oct;5(10):1927-1941.

# Multiple LCMS Experiments: Good with the Bad

- **There is a lot of information in there**
  - Peptide/protein IDs
  - Quantitative data
  - Statistical assessment

- **The information may be noisy**
  - Retention time drift
  - Instrument response noise

- **Are there methods to leverage this information?**
  - Without 'perfect' chromatography?
  - Without strict alignment?

# PEPPeR Concepts – Samples and Data Acquisition

# PEPPeR Concepts – Data Processing

# PEPPeR Concepts – Processing Continued…

# PEPPeR Concepts – Analysis and Follow up



Landmark Matching

Peak Matching

Marker Discovery

Targeted MS/MS Identifications

# Landmark Matching: Identity Propagation

- Use accurate mass, relative retention order comparison to identify peaks

# Landmark Matching: Identity Propagation

- Use accurate mass, relative retention order comparison to identify peaks

**Current Experiment**

m/z=999.4991

X

A

B

m/z=999.4996

Y

C

**Comparison Experiment**

M

B

A

N

C

APEPTIDEK
m/z=999.4993

APDITEPEK
m/z=999.4993

# Landmark Matching: Identity Propagation

# Nuts and bolts: How it works

- Match features to sequenced peptides in a single LCMS run

- Refine/recalibrate m/z tolerance

- Re-match features to sequenced peptides in a single LCMS run

- Now compare list of all features to Basis Set for mass, relative elution order matches given landmarks as reference points – *propagation of identified features across multiple experiments*

# Landmark Scoring and Confidence

$$S = \sum_{i=1}^{w} \left[ \xi(\Lambda_{-i}, \Lambda_0) + \xi(\Lambda_0, \Lambda_i) \right]$$

$$\xi(m,n) = \begin{cases} 1 \text{ if } \tau(m) < \tau(n) \\ \text{if } \tau(m) > \tau(n) \begin{cases} 0.5 \text{ if } \tau(n) - \tau(m) < \delta \text{ and } \mu(m) + \sigma(m) > \mu(n) - \sigma(n) \\ -1 \text{ if } else \end{cases} \\ 0 \text{ if } else \end{cases}$$

$$P_{overall} = P_{m/z} P_{landmark}$$

$$P_{landmark} = P(landmark \mid m/z) = $$

$$\frac{P(m/z \mid landmark)P(landmark)}{P(m/z \mid landmark)P(landmark) + (1 - P(m/z \mid landmark))(1 - P(landmark))}$$

Let:

$\Lambda$ be a list of peptides observed in the comparison experiment ordered by elution time. Here, elution time is defined by the centroid of all MS/MS scans leading to the identification of the peptide.

$\Lambda_0$ is defined as the position of the putative assignment in $\Lambda$

$\mu(x)$ be the centroid of elution time of peptide $x$ in the comparison experiment (in scans)

$\sigma(x)$ be the standard deviation of elution time of peptide $x$ in the comparison experiment (in scans)

$\tau(x)$ be the centroid of elution time of peptide $x$ in the current experiment (in seconds)

$\delta$ be the average retention time peak width, such that peptides eluting within $\delta$ sec are considered to be co-eluting (typically $\delta$ = 30 s)

$w$ the number of peptides to consider before and after the putative assignment on the landmark list (typically $w$ = 3)

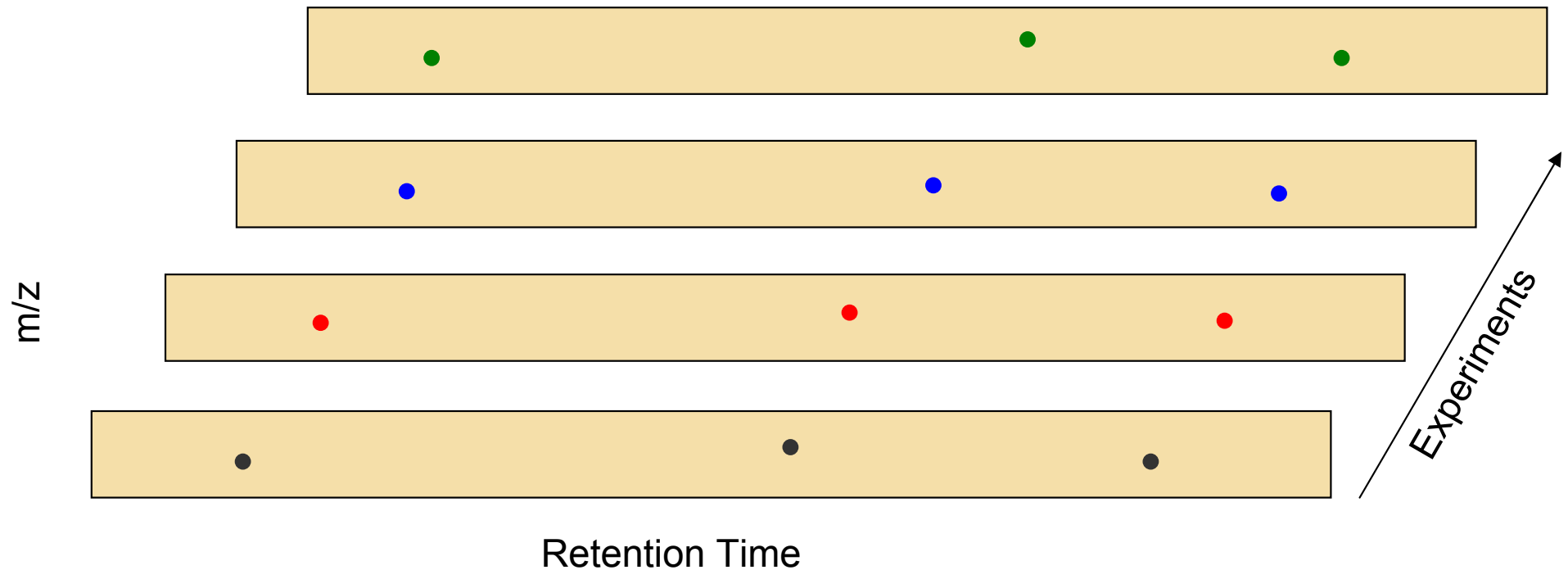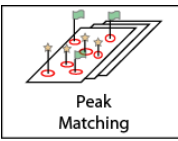# Peak Matching: Recognizing Identical Features

- Use landmarks to derive corrections and tolerances for clustering of features across LCMS experiments

  - Break down the problem to make it parallelizable

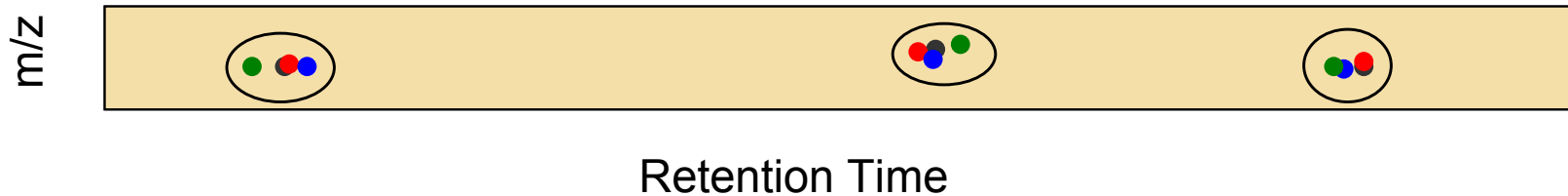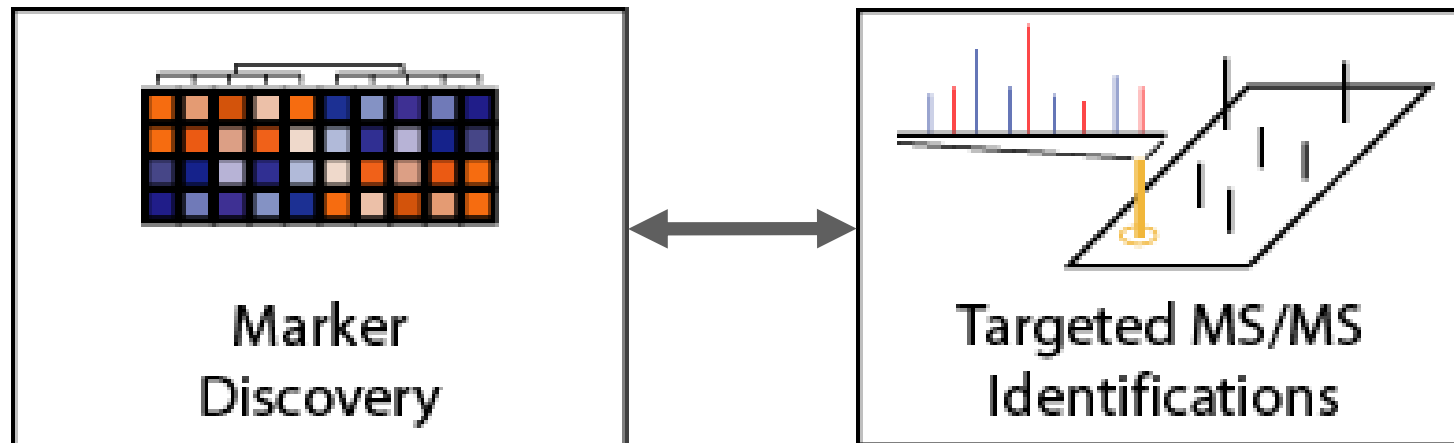# Peak Matching: Recognizing Identical Features

- Use landmarks to derive corrections and tolerances for clustering of features across LCMS experiments

# Peak Matching: Recognizing Identical Features

- ## Use landmarks to derive corrections and tolerances for clustering of features across LCMS experiments

  - Gaussian mixture model (GMM) with parameters determined by maximizing Likelihood ratio using Expectation Maximization (EM)
  - Number of clusters determined using Bayesian Information Criterion (BIC)
  - Coalesce clusters if M/Z and RT variation is within tolerance

# Parameterized Peaks

| Peak ID | m/z | R.T. | z | Run 1 | Run 2 | Run 3 | Run … | Identity |
|---------|-----|------|---|-------|-------|-------|-------|----------|
| 1 | 490.3144 | 62.0 | 3 | 607.6 | 544.2 | 581.0 | … | |
| 2 | 743.3549 | 56.2 | 3 | 694.4 | 682.6 | 691.4 | … | |
| 3 | 999.4991 | 22.5 | 2 | 209.6 | 247.6 | 232.6 | … | APEPTIDEK |
| 4 | 396.7187 | 20.5 | 3 | 321.7 | 344.9 | 318.5 | … | |
| 5 | 934.6045 | 31.7 | 2 | 722.7 | 753.0 | 701.3 | … | |
| 6 | 678.1993 | 32.4 | 3 | 371.2 | 387.2 | 441.4 | … | |
| 7 | 999.4994 | 56.8 | 2 | 857.1 | 811.0 | 750.5 | … | APDITEPEK |
| 8 | 526.6502 | 46.0 | 3 | 183.6 | 169.0 | 155.2 | … | |
| 9 | 1105.3597 | 69.4 | 3 | 1130.1 | 1075.7 | 1075.1 | … | |
| 10 | 1292.0880 | 34.5 | 2 | 709.7 | 614.0 | 656.0 | … | |



Marker Discovery  ⟷  Targeted MS/MS Identifications

# Calibration and Landmark Performance

*Scale Mixture*

| | A | B | C | D | E | F | G | H | I |
|---|---|---|---|---|---|---|---|---|---|
| **Aprotinin** | 1 | 2 | 3 | 10 | 20 | 30 | 100 | 200 | 300 |
| **Ribonuclease A** | 300 | 1 | 2 | 3 | 10 | 20 | 30 | 100 | 200 |
| **Myoglobin** | 200 | 300 | 1 | 2 | 3 | 10 | 20 | 30 | 100 |
| **beta-Lactoglobulin** | 100 | 200 | 300 | 1 | 2 | 3 | 10 | 20 | 30 |
| **alpha Casein** | 30 | 100 | 200 | 300 | 1 | 2 | 3 | 10 | 20 |
| **Carbonic anhydrase** | 20 | 30 | 100 | 200 | 300 | 1 | 2 | 3 | 10 |
| **Ovalbumin** | 10 | 20 | 30 | 100 | 200 | 300 | 1 | 2 | 3 |
| **Fibrinogen** | 3 | 10 | 20 | 30 | 100 | 200 | 300 | 1 | 2 |
| **BSA** | 2 | 3 | 10 | 20 | 30 | 100 | 200 | 300 | 1 |
| **Transferrin** | 100 | 100 | 100 | 100 | 100 | 100 | 100 | 100 | 100 |
| **Plasminogen** | 30 | 30 | 30 | 30 | 30 | 30 | 30 | 30 | 30 |
| **beta-Galactosidase** | 10 | 10 | 10 | 10 | 10 | 10 | 10 | 10 | 10 |

All concentrations in fmol/ul (nM)
Inject 1 ul x 5 replicates each

Peaks with IDs (avg. per run):
165 $\Rightarrow$ 281        +70%

False positive rate:
93%    $p < 0.005$
100%  $p < 0.05$

False negative rate:
~2%

# Measurement of Ratios with Variability

**Variability Mixture**

| | Person A | | Person B | | Person C | | Person D | | Person E | |
|---|---|---|---|---|---|---|---|---|---|---|
| | α | β | α | β | α | β | α | β | α | β |
| Aprotinin | 100 | 5 | 100 | 5 | 100 | 5 | 100 | 5 | 100 | 5 |
| Ribonuclease A | 100 | 100 | 100 | 100 | 100 | 100 | 100 | 100 | 100 | 100 |
| Myoglobin | 100 | 100 | 100 | 100 | 100 | 100 | 100 | 100 | 100 | 100 |
| beta-Lactoglobulin | 50 | 1 | 50 | 1 | 50 | 1 | 50 | 1 | 50 | 1 |
| alpha Casein | 100 | 10 | 100 | 10 | 100 | 10 | 100 | 10 | 100 | 10 |
| Carbonic anhydrase | 100 | 100 | 100 | 100 | 100 | 100 | 100 | 100 | 100 | 100 |
| Ovalbumin | 5 | 10 | 5 | 10 | 5 | 10 | 5 | 10 | 5 | 10 |
| Fibrinogen | 25 | 25 | 25 | 25 | 25 | 25 | 25 | 25 | 25 | 25 |
| BSA | 200 | 200 | 200 | 200 | 200 | 200 | 200 | 200 | 200 | 200 |
| Transferrin | 10 | 5 | 10 | 5 | 10 | 5 | 10 | 5 | 10 | 5 |
| Plasminogen | 2.5 | 25 | 2.5 | 25 | 2.5 | 25 | 2.5 | 25 | 2.5 | 25 |
| beta-Galactosidase | 1 | 10 | 1 | 10 | 1 | 10 | 1 | 10 | 1 | 10 |

All concentrations in fmol/ul (nM)
Inject 1 ul x 5 replicates each

**Complex Variability Mixture:**

**Mix α + Mitochondrial Protein from 2 wk. mouse liver**

**Mix β + Mitochondrial Protein from 6 wk. mouse liver**

**1 prep each sample, 6 injections each**

# PEPPeR and GenePattern



- GenePattern is a suite of tools originally developed for microarray analysis
    - AIM: reproducible research through well-defined processing pipelines

- Many analysis modules available
    - PEPPeR: Landmark Matching and Peak Matching
    - Daisy-chainable into pipelines
    - Feed into statistical tools

# PEPPeR in GenePattern

# Insert your favorite stuff here…

- Landmark Matching is platform agnostic
  - Need to get your data into a few simple flat-file formats and then zip them up together
  - Search engines i.e. SEQUEST, SpectrumMill, Mascot, etc.
  - Peak Pickers: MAPQUANT, msInspect, Decon2LS, etc.
  - Some helper apps can be found with the PEPPeR bundle on the GenePattern website

- All works via web-client interface
  - Just press go (but beware of this!)

# Landmark Matching Output



The main output is a zipped directory of all the processed files. This can be used as input into the PeakMatch module.

It is a good idea to check the error log to make sure that everything was processed correctly.

# Peak Matching Interface

# GenePattern Downstream Tools

- **Differential analysis/marker selection**
  - Gene/Class neighbors
  - Comparative marker selection
  - Gene Set Enrichment Analysis

- **Class Prediction – supervised learning – with cross-validation**
  - Regression trees
  - K-nearest neighbors
  - Neural networks
  - Support Vector Machine

- **Class Discovery – unsupervised learning**
  - Hierarchical clustering
  - Self-organizing maps
  - Principal Component Analysis

- **Data Visualization**
  - Heat Maps, etc.

Note: Data analysis on subsequent slides done using GenePattern

# Discovery of Novel Markers with PEPPeR



β
α

peptides / m/z features

- ● Designed accurate mass 'inclusion lists' to hit these targets

- ▪ Confident IDs of previously identified peptides agree 100% of the time (59/59)

- ▪ 60 novel confident peptide IDs
  - ▪ 25 belong to proteins in the mix
    - ▪ 24/25 are changing
  - ▪ 35 are from proteins not designed to be in the mixture

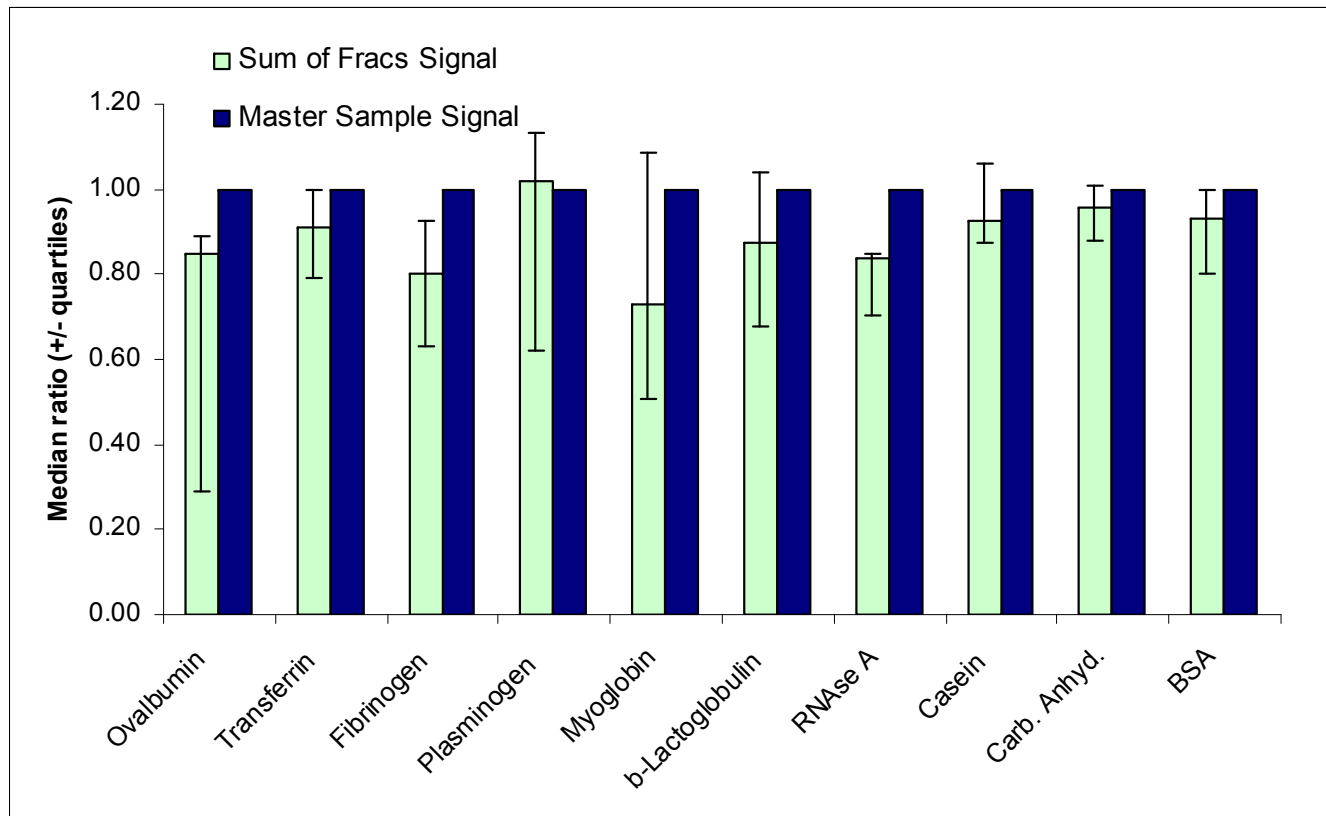| gi Number | Species | Name |
|---|---|---|
| 223424 | *E. coli* | RNA polymerase β' |
| 38491462 | *E. coli* | GroEL |
| 42144 | *E. coli* | NusA |
| 42818 | *E. coli* | RNA polymerase β |
| 42900 | *E. coli* | Ribosomal protein S1 |
| 26249756 | *E. coli* | Argininosuccinate synthase |
| 8099322 | *B. taurus* | κ-casein |

B-Galactosidase had 1:10 ratio!
Casein had 10:1 ratio!

# In-silico defractionation of 2D-LC

- Wanted to mimic SCX fractionation scheme

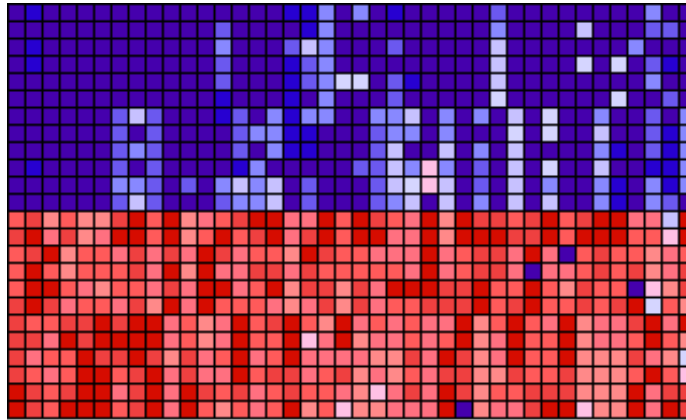| | Frac. 1 | Frac. 2 | Frac. 3 | Frac. 4 | Frac. 5 | Master |
|---|---|---|---|---|---|---|
| *Ovalbumin* | 0 | 0 | 0 | 5 | 0 | 5 |
| *Transferrin* | 0 | 0 | 0 | 0 | 5 | 5 |
| *Fibrinogen* | 0 | 12.5 | 0 | 0 | 0 | 12.5 |
| *Plasminogen* | 5 | 5 | 5 | 5 | 5 | 25 |
| *Myoglobin* | 0 | 0 | 50 | 0 | 0 | 50 |
| *$\beta$-Lactoglobulin* | 0 | 0 | 12.5 | 25 | 12.5 | 50 |
| *RNAse A* | 50 | 10 | 0 | 0 | 0 | 60 |
| *Casein* | 5 | 50 | 5 | 0 | 0 | 60 |
| *Carb. Anhyd.* | 0 | 0 | 0 | 50 | 50 | 100 |
| *BSA* | 100 | 0 | 0 | 0 | 0 | 100 |

All values in fmol injected

# Breast Cancer Biomarker Discovery

- Sample source: nipple aspirate fluid (NAF) from malignancy affected breast
  - Unaffected contra-lateral breast used for control
  - Pools of several patient samples made <- low starting material

- Samples depleted of abundant proteins by affinity chromatography

- Separate ID-centric (fractionation) and Pattern Centric runs conducted for PEPPeR analysis

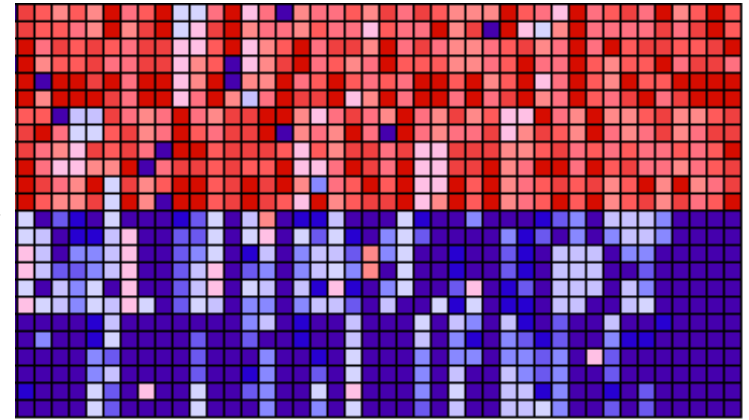- Performed marker selection with allowed FDR of 5%

# Breast Cancer Marker Selection
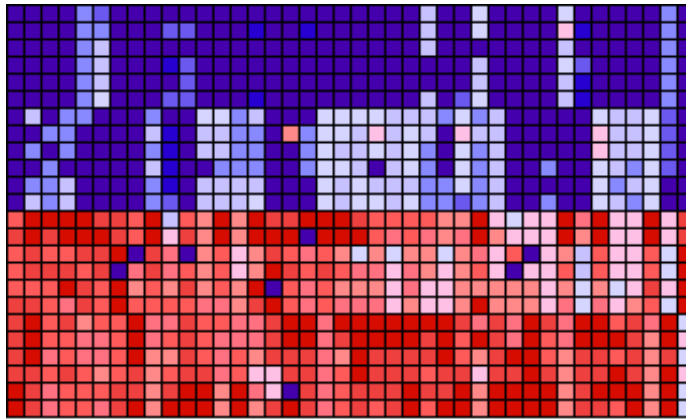


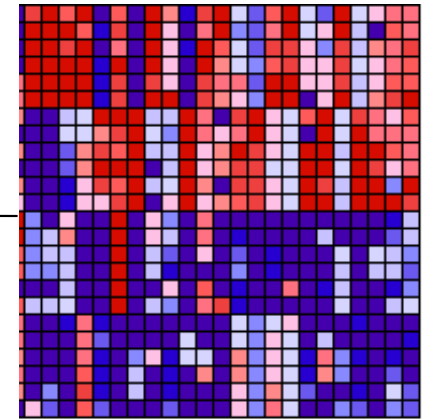n=1520

Features down in cancer

Features up in cancer

n=264

Assigned peptides down in cancer

Assigned peptides up in cancer

# Features vs. Assignments

- There's more out there than we can catalog
  - Low intensity features never trigger MS$^2$ in complex samples
  - Unidentified features may be better classifiers

- Direct follow-up easily achieved
  - We know exactly where and when to look
  - Targeted accurate mass methods can be employed

- Hopefully increase coverage and confidence in certain **proteins** as markers, rather than just peptides or features

# Summary – what I hope you learned

- PEPPeR: Landmark Matching and Peak Matching
  - Keep track of all of those pesky peaks that you picked!

- GenePattern: A web-based tool to coordinate reproducible research

- An entrée into downstream discovery methods in an automated pipeline (more GenePattern)

- Some real world examples of its application

# Acknowledgements

- Broad Proteomics:
  - Steve Carr
  - **D.R. Mani**
  - Vincent Fusaro
  - Mike Gillette

- Broad GenePattern Team:
  - **Michael Reich**
  - Josh Gould

- Church Lab, Harvard Medical School
  - George Church
  - **Kyriacos Leptos**

# URLs:

- PEPPeR / GenePattern:
    - http://www.broad.mit.edu/cancer/software/genepattern/
    - http://www.broad.mit.edu/cancer/software/genepattern/desc/proteomics.html

- MAPQUANT:
    - http://arep.med.harvard.edu/MapQuant/

# Live DEMO Time

- Thanks to the many developers, beta testers, and users



Note: PNNL is always looking for good and knowledgeable informatics staff and post-docs (see us afterwards for more information, or visit http://jobs.pnl.gov/)

# Funding for Tool Development

- DOE Office of Biological and Environmental Research
  - http://ober-proteomics.pnl.gov/
- NIH
  - National Center for Research Resources
    - http://ncrr.pnl.gov/
  - National Institute of Allergy and Infectious Diseases
    - http://proteomicsresource.org/
  - National Cancer Institute
  - National Institute of General Medical Sciences
  - National Institute of Diabetes & Digestive & Kidney Diseases