# Mass Spectrometry Proteomics for the Computational Biologist

December 1, 2006

John T. Prince

"...our ability to collect large proteomic data sets already outstrips our ability to validate, to interpret and to integrate such data for the purpose of creating biological knowledge"

Patterson and Aebersold   (*Nature Genetics* **33**, 318 (2003))

MARCOTTE.LAB
THE UNIVERSITY OF TEXAS AT AUSTIN

# Mass Spectrometry (MS) Proteomics
## Needs Computational Biologists

"…our ability to collect large proteomic data sets already outstrips our ability to validate, to interpret and to integrate such data for the purpose of creating biological knowledge"

- Patterson and Aebersold  (*Nature Genetics* **33**, 318 (2003))

# MS Proteomics

- How?
- Data?
- Problems?

# Why Proteomics?
## and not just Transcriptomics

- Proteins are the actual players
- mRNA not necessarily proportional to protein level
  - translational control
  - degradation
- Post-translational modifications alter cell state
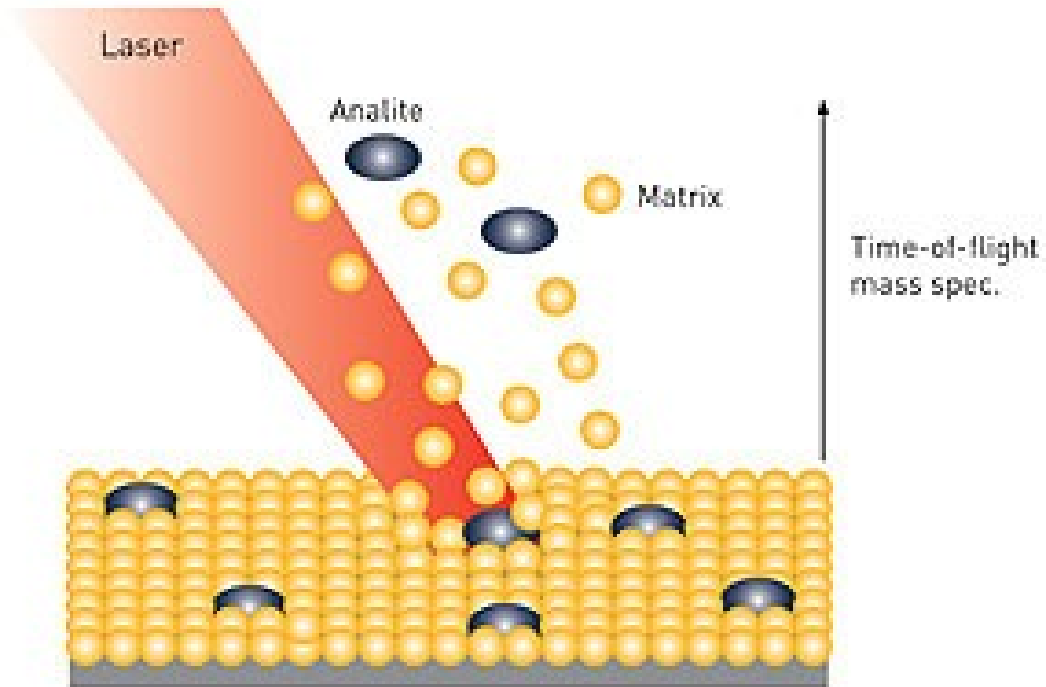- Cellular localization

# Mass Spec (Proteomics)



- Ionization
  - MALDI
  - ESI
- m/z Analysis
  - TOF
  - Quadrupole
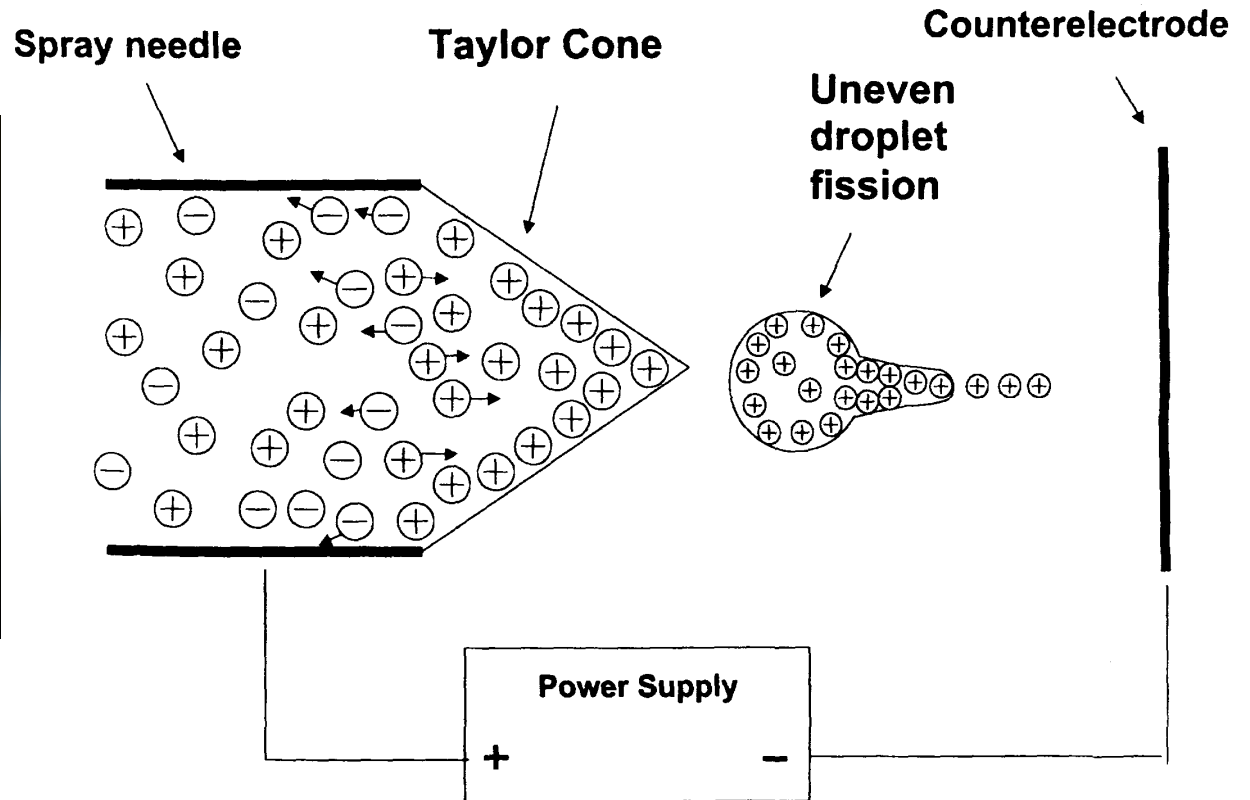  - Ion Trap
  - FTICR
  - Orbitrap

# MALDI
## Matrix Assisted Laser Desorption Ionization



http://www.eurogentec.com/module/images2/p23_3.jpg

# ESI
# Electrospray Ionization



Spray needle    Taylor Cone    Counterelectrode

Uneven droplet fission

http://www.phoenix-st.com/images/splash2.jpg

Power Supply

+    −

# TOF
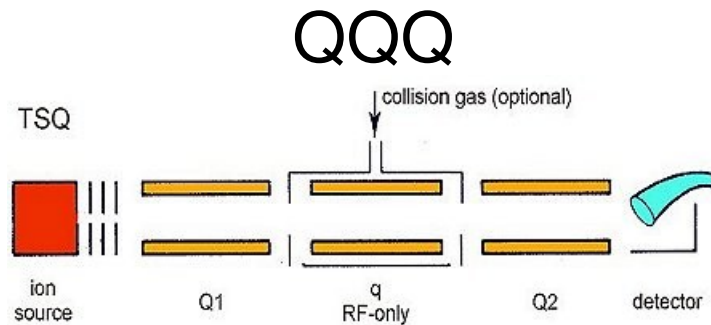## Time of Flight

detector

vacuum

## TOF (reflectron)

high velocity

low velocity

detector

vacuum

# Q
## (e.g., Q-TOF, QQQ)Quadrupole



quadrupole rods

TO DETECTOR

IONS

exit slit
(to detector)

resonant ion
(detected)

source slit

non-resonance ion
(not detected)

http://www.bris.ac.uk/nerclsmsf/images/quadrupole.gif

## QQQ



TSQ

collision gas (optional)

ion source

Q1

q RF-only

Q2

detector

http://www.rzuser.uni-heidelberg.de/~bl5/ency/pics/t_tsq1.jpg

# Quadrupole Ion Trap



Ion source

Injected ions

End-cap electrode

Ring electrode

$+V_{res}\cos(\omega_{res}t+\phi_{res})$

Resonance AC voltage

$-V_{res}\cos(\omega_{res}t+\phi_{res})$

Aperture

RF voltage $V\cos(\Omega t+\phi)$

Detector

End Cap

Ring

Ring

End Cap

Ion Detector

$a_z$

Axial (z) Stability Region

Example of Mass-selective Ejection

Ejection of Ions in Axial (z) Dimension

$q_z$

Radial (r) Stability Region

Example of Mass-selective Storage

$q_z = \dfrac{4eV}{mr_0^2\omega^2}$

$a_z = \dfrac{8eU}{mr_0^2\omega^2}$

http://www.rzuser.uni-heidelberg.de/~bl5/ency/pics/q_trap01.jpg
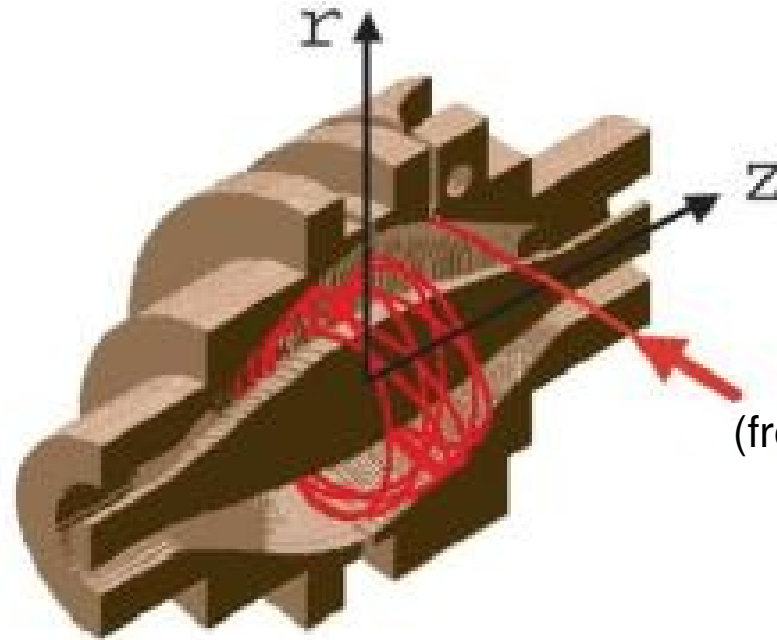K. Yoshinari, Rapid Commun. Mass Spectrom. 14, 215-223 (2000)

# FT-ICR
## Fourier Transform Ion Cyclotron Resonance



receiver plate

Magnetic Field, B

trap plate

transmitter plate

time-domain signal

FT

mass spectrum

http://www.ivv.fraunhofer.de/ms/ms-analyzers.html

# FT-Orbi
## Orbitrap



(from Linear Ion Trap via C trap)

# Mass Spectrometry (Proteomics)

- Ionization
  - ESI (Electrospray Ionization)
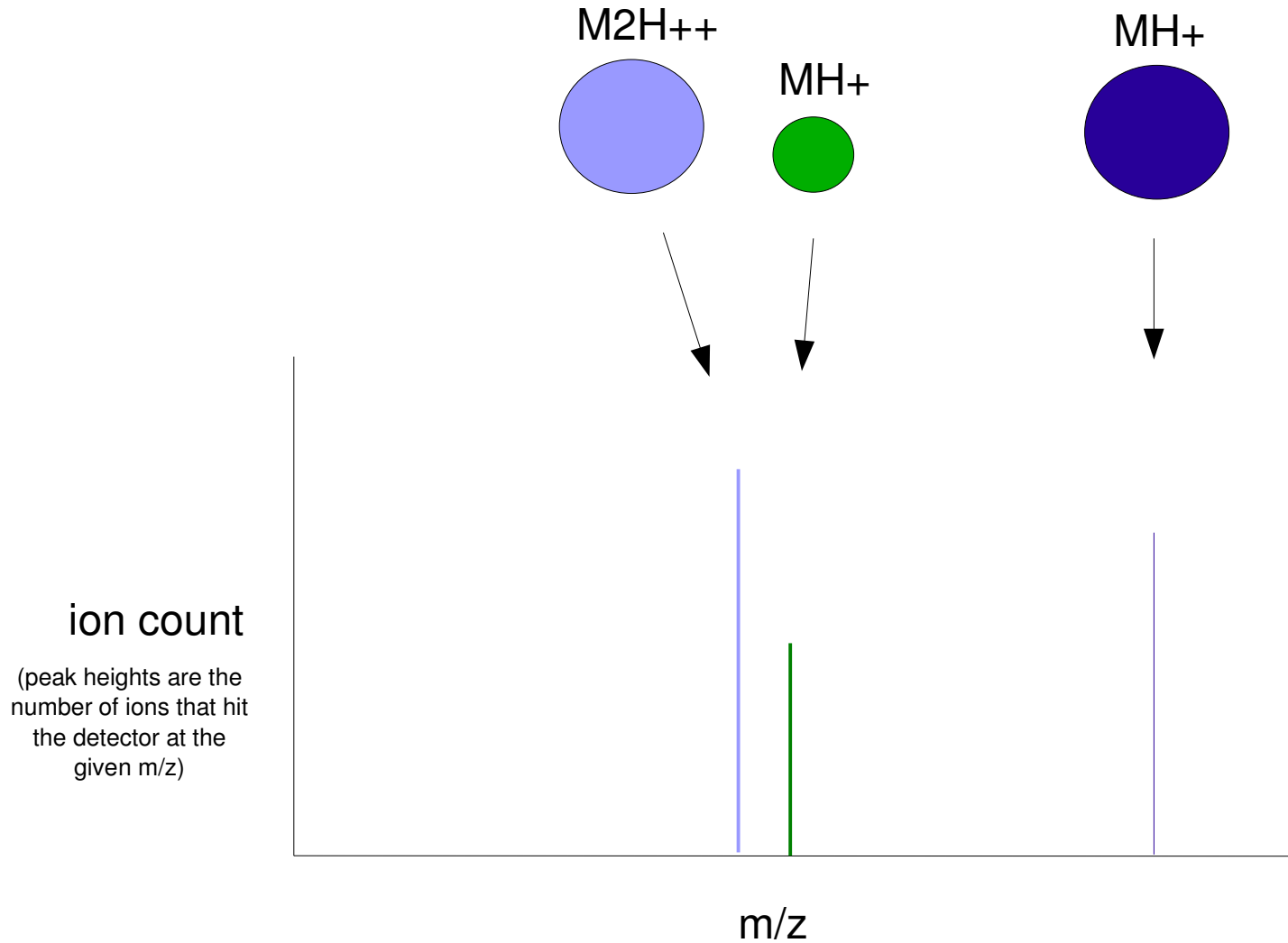  - MALDI (Matrix Assisted Laser Desorption Ionization)
- m/z Analysis
  - TOF (Time of Flight)
  - Q ([e.g. Q-TOF] Quadrupole)
  - Ion Trap
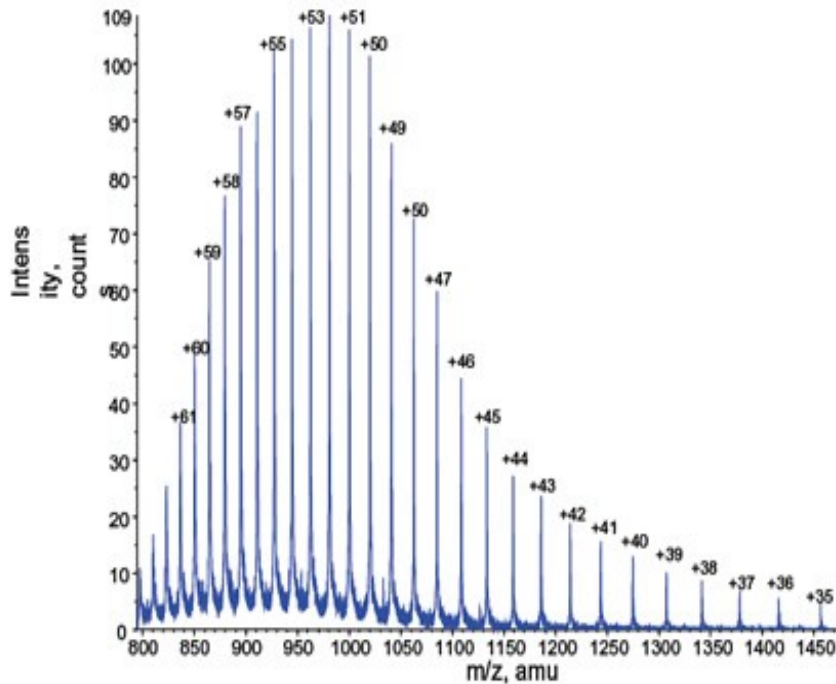  - FTICR (Fourier Transform Ion Cyclotron Resonance)
  - Orbitrap

# Data

- Spectrum
- ESI Protein Spectrum
- 2D MALDI imaging
- PMF (peptide mass fingerprinting)
- LC-MS
- Peptide Fragmentation
- MudPIT

# Spectrum

M2H++

MH+

MH+

ion count

(peak heights are the
number of ions that hit
the detector at the
given m/z)

m/z

# Why Not Proteins?

## Multiple Charge States (ESI)

## PTMs (Post-Translational Modifications)



http://www-methods.ch.cam.ac.uk/siteimages/sw3.jpg

Table 1. Some common and important post-translational modifications

| PTM type | ΔMass[a] (Da) | Stability[b] | Function and notes |
|---|---|---|---|
| Phosphorylation | | | Reversible, activation/inactivation of enzyme activity, modulation of molecular interactions, signaling |
| pTyr | +80 | +++ | |
| pSer, pThr | +80 | +/++ | |
| Acetylation | +42 | +++ | Protein stability, protection of N terminus. Regulation of protein–DNA interactions (histones) |
| Methylation | +14 | +++ | Regulation of gene expression |
| Acylation, fatty acid modification | | | Cellular localization and targeting signals, membrane tethering, mediator of protein–protein interactions |
| Farnesyl | +204 | +++ | |
| Myristoyl | +210 | +++ | |
| Palmitoyl | +238 | +/++ | |
| etc. | | | |
| Glycosylation | | | |
| N-linked | >800 | +/++ | Excreted proteins, cell–cell recognition/signaling |
| O-linked | 203, >800 | +/++ | O-GlcNAc, reversible, regulatory functions |
| GPI anchor | >1,000 | ++ | Glycosylphosphatidylinositol (GPI) anchor. Membrane tethering of enzymes and receptors, mainly to outer leaflet of plasma membrane |
| Hydroxyproline | +16 | +++ | Protein stability and protein–ligand interactions |
| Sulfation (sTyr) | +80 | + | Modulator of protein–protein and receptor–ligand interactions |
| Disulfide bond formation | –2 | ++ | Intra- and intermolecular crosslink, protein stability |
| Deamidation | +1 | +++ | Possible regulator of protein–ligand and protein–protein interactions, also a common chemical artifact |
| Pyroglutamic acid | –17 | +++ | Protein stability, blocked N terminus |
| Ubiquitination | >1,000 | +/++ | Destruction signal. After tryptic digestion, ubiquitination site is modified with the Gly-Gly dipeptide |
| Nitration of tyrosine | +45 | +/++ | Oxidative damage during inflammation |

[a]A more comprehensive list of PTM Δmass values can be found at: http://www.abrf.org/index.cfm/dm.home
[b]Stability: + labile in tandem mass spectrometry, ++ moderately stable; +++ stable.

http://www.nature.com/nbt/journal/v21/n3/images/nbt0303-255-T1.gif

# MALDI on Biological Sample



- Signal Processing
- Classification Analysis

http://www.mcb.mcgill.ca/~hallett/GEP/PLecture4/image002.gif

# MALDI-TOF Application



- organize data
- integrate data
- mine data

# PMF
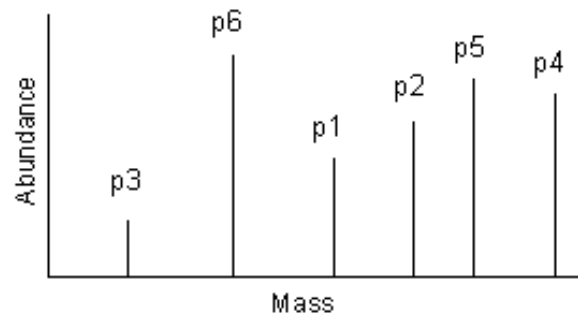## Peptide Mass Fingerprinting

**Denature**

MITGIQITKAANDLLNDSFRLLDSKGEACIVAAGYAEVVSREYPQLTIVSGQQRFNSLTPSL

**Digest**

| MITGIQITK | AANDLLNDSFR | LLDSK | GEACIVAAGYAEVVSR | EYPQLTIVSGQQR | FNSLTPSL |
|-----------|-------------|-------|------------------|---------------|----------|
| p1 | p2 | p3 | p4 | p5 | p6 |

**MS**



Abundance vs Mass spectrum showing peaks: p3, p6, p1, p2, p5, p4

– <span style="color:red">statistical validation</span>

# LC-MS
## Liquid Chromatography MS

Reverse Phase Chromatography (RPC)

ESI

MS

5 μM

http://www.ionsource.com/tutorial/chromatography/rphplc10.gif
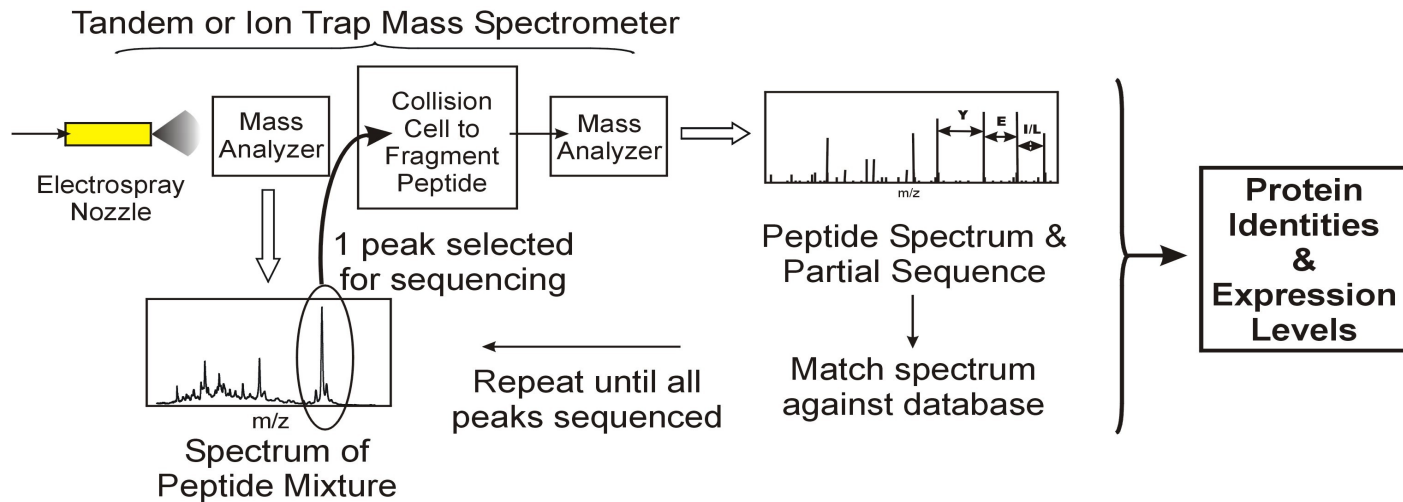
– elution prediction
– registration

time (hydrophobicity)

m/z

# MS/MS (Peptide Fragmentation)

# Peptide Fragmentation (MS/MS)

peptide

| break

y3    y2    y1

AFTG

b1    b2    b2

Ala        Phe        Thr        Gly

intensity (ion count)

71.03        147.06        101.04        57.02

b1

b2

b3

AFTG (MH+)

m/z

# Shotgun Proteomics



Cells → Protein Extract → ~50-200,000 peptides → Partial Separation by HPLC → Partial Separation by second HPLC

Protease

Tandem or Ion Trap Mass Spectrometer

Electrospray Nozzle → Mass Analyzer → Collision Cell to Fragment Peptide → Mass Analyzer → Peptide Spectrum & Partial Sequence → Protein Identities & Expression Levels

1 peak selected for sequencing

Spectrum of Peptide Mixture

Repeat until all peaks sequenced

Match spectrum against database

> 3 million data points per experiment

quantitation
peptide fragmentation prediction
spectra comparison metrics
peptides to proteins
integrating bayesian priors

# MuDPIT
## Multidimensional Protein Identification Technology



SCX
(Strong Cation Exchange)

RP
(Reverse Phase-C18)

ESI → Ion Trap

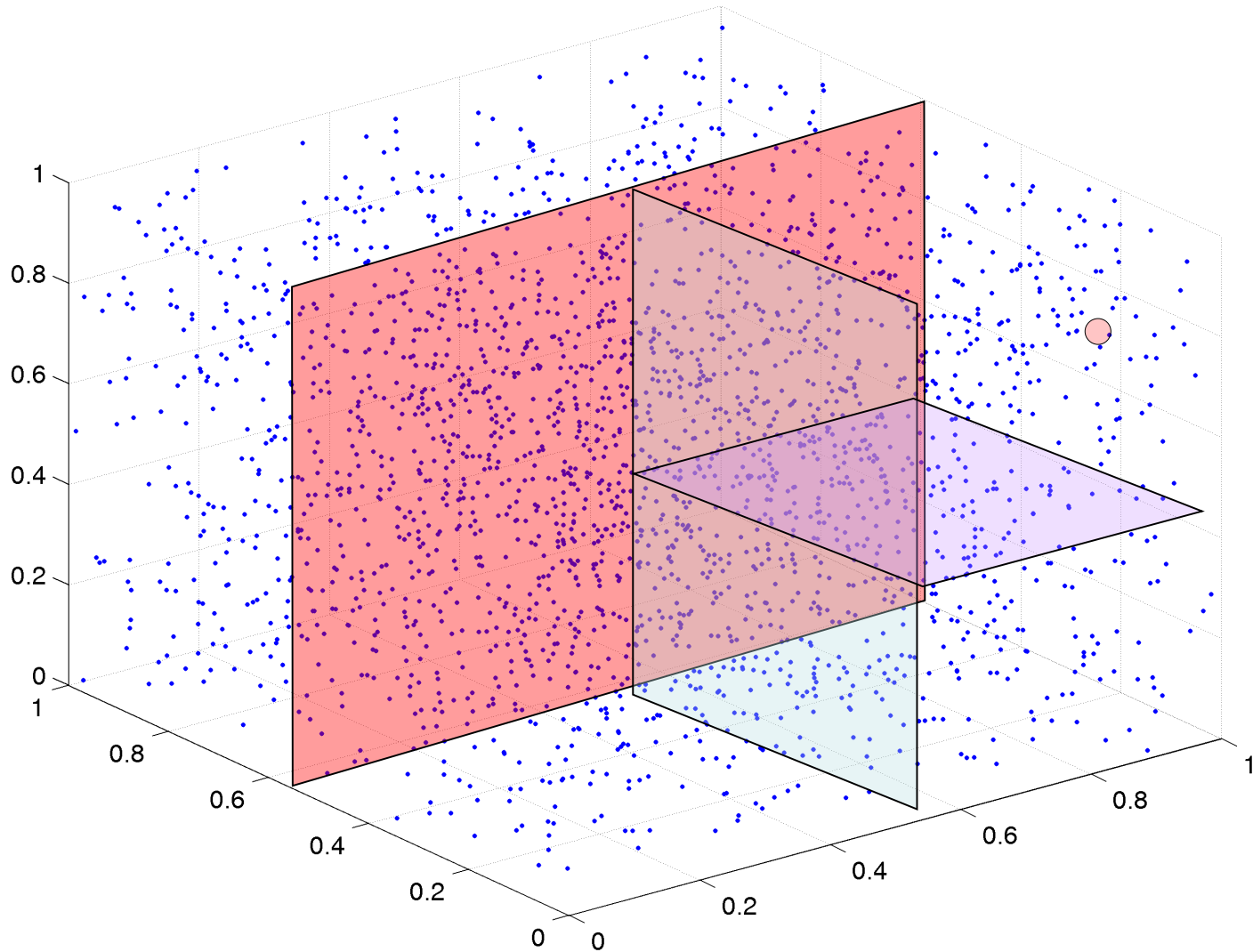– multi-dimensional dataset registration

# PTM's
## Post-translational Modifications

**Table 1. Some common and important post-translational modifications**

| PTM type | ΔMass[a] (Da) | Stability[b] | Function and notes |
|---|---|---|---|
| Phosphorylation | | | Reversible, activation/inactivation of enzyme activity, modulation of molecular interactions, signaling |
| pTyr | +80 | +++ | |
| pSer, pThr | +80 | +/++ | |
| Acetylation | +42 | +++ | Protein stability, protection of N terminus. Regulation of protein–DNA interactions (histones) |
| Methylation | +14 | +++ | Regulation of gene expression |
| Acylation, fatty acid modification | | | Cellular localization and targeting signals, membrane tethering, mediator of protein–protein interactions |
| Farnesyl | +204 | +++ | |
| Myristoyl | +210 | +++ | |
| Palmitoyl | +238 | +/++ | |
| etc. | | | |
| Glycosylation | | | |
| N-linked | >800 | +/++ | Excreted proteins, cell–cell recognition/signaling |
| O-linked | 203, >800 | +/++ | O-GlcNAc, reversible, regulatory functions |
| GPI anchor | >1,000 | ++ | Glycosylphosphatidylinositol (GPI) anchor. Membrane tethering of enzymes and receptors, mainly to outer leaflet of plasma membrane |
| Hydroxyproline | +16 | +++ | Protein stability and protein–ligand interactions |
| Sulfation (sTyr) | +80 | + | Modulator of protein–protein and receptor–ligand interactions |
| Disulfide bond formation | −2 | ++ | Intra- and intermolecular crosslink, protein stability |
| Deamidation | +1 | +++ | Possible regulator of protein–ligand and protein–protein interactions, also a common chemical artifact |
| Pyroglutamic acid | −17 | +++ | Protein stability, blocked N terminus |
| Ubiquitination | >1,000 | +/++ | Destruction signal. After tryptic digestion, ubiquitination site is modified with the Gly-Gly dipeptide |
| Nitration of tyrosine | +45 | +/++ | Oxidative damage during inflammation |

[a]A more comprehensive list of PTM Δmass values can be found at: http://www.abrf.org/index.cfm/dm.home
[b]Stability: + labile in tandem mass spectrometry, ++ moderately stable; +++ stable.

# Spectra In Metric-Space



2300 points (at random) in 3D space

# Data Format/Storage/Sharing

- Object Models still being worked out
- Huge Datasets
  - how much to save?
  - how much is it worth?
- Sharing
  - OPD
  - Peptide Atlas
  - PRIDE
  - GPM

# Biological Integration

# mRNA vs. Protein

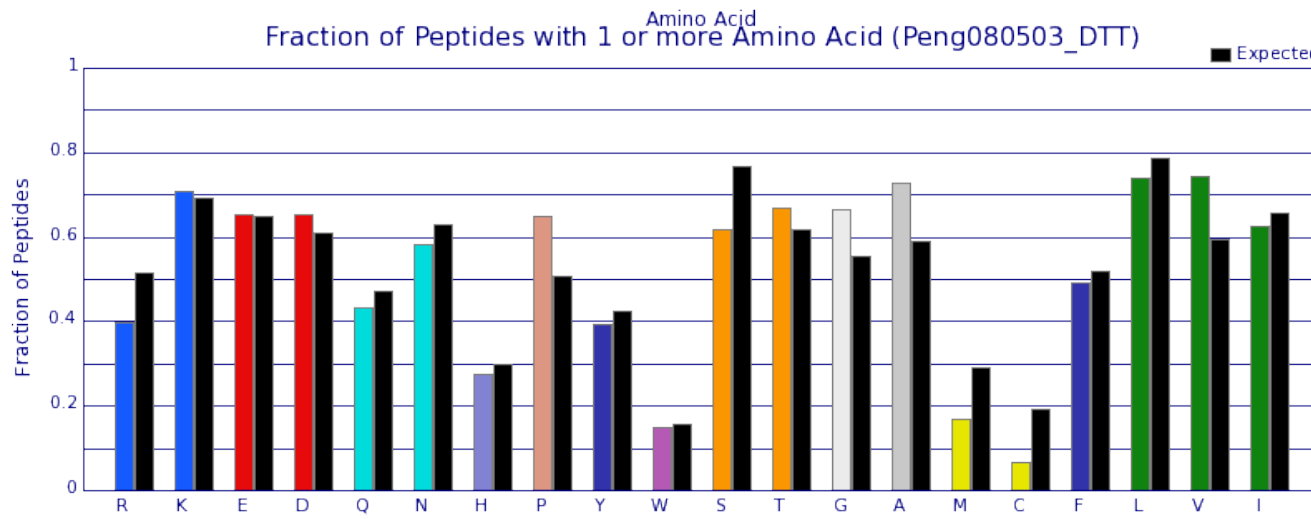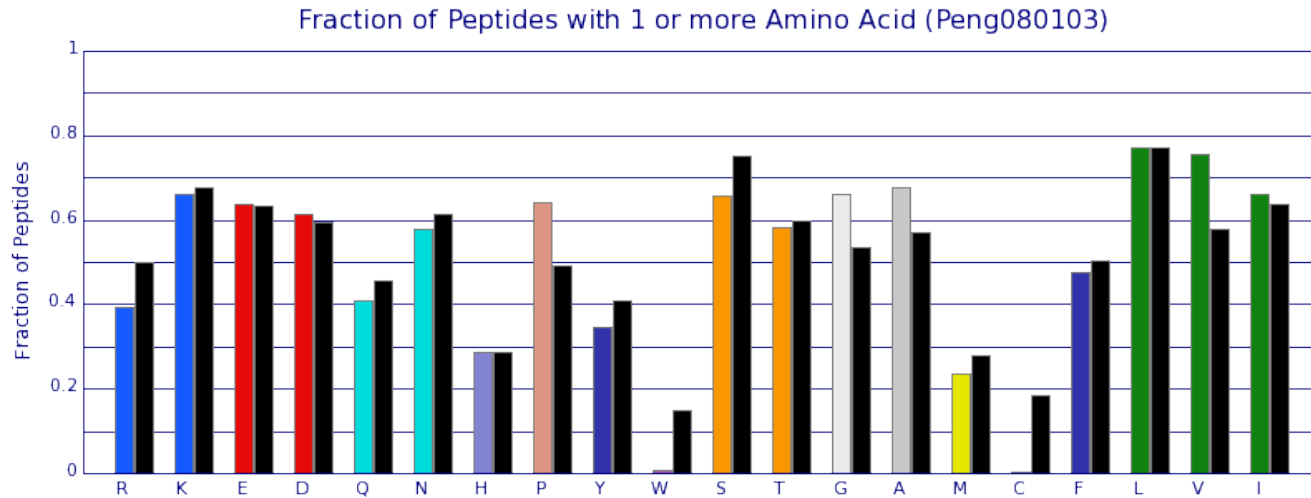| Source | Subject | Perturbation (or sample) | Num Genes | Correlation |
|---|---|---|---|---|
| Ideker et. al. | Yeast | +/- gal (gal inducing media) | 289 | $r_p = 0.61$ |
| Futcher et. al. | Yeast | 2% ethanol/ 2% glucose | 148 | $r_s = 0.74$ $r_p = 0.76$ [a] |
| Washburn et. al. | Yeast | rich/minimal | 678 | $r_s = 0.45$ |
| Griffin et. al. | Yeast | 2% ethanol/ 2% galactose | 245 | $r_s = 0.21$ |
| Gygi et. al. | Yeast | mid-log | 106 | $r_p = 0.94$ $r_s = 0.59$ [b] $r_p = 0.356$ [c] |
| Chen et. al. | Lung adenocarcinomas | 57 stage I, 19 stage III, 9 non-neoplastic | 98 (165 prots) | $r_p = -0.025$ [d] |

a = after normalizing the data
b = calculated by Futcher et. al.
c = 73 genes with lower abundance transcripts
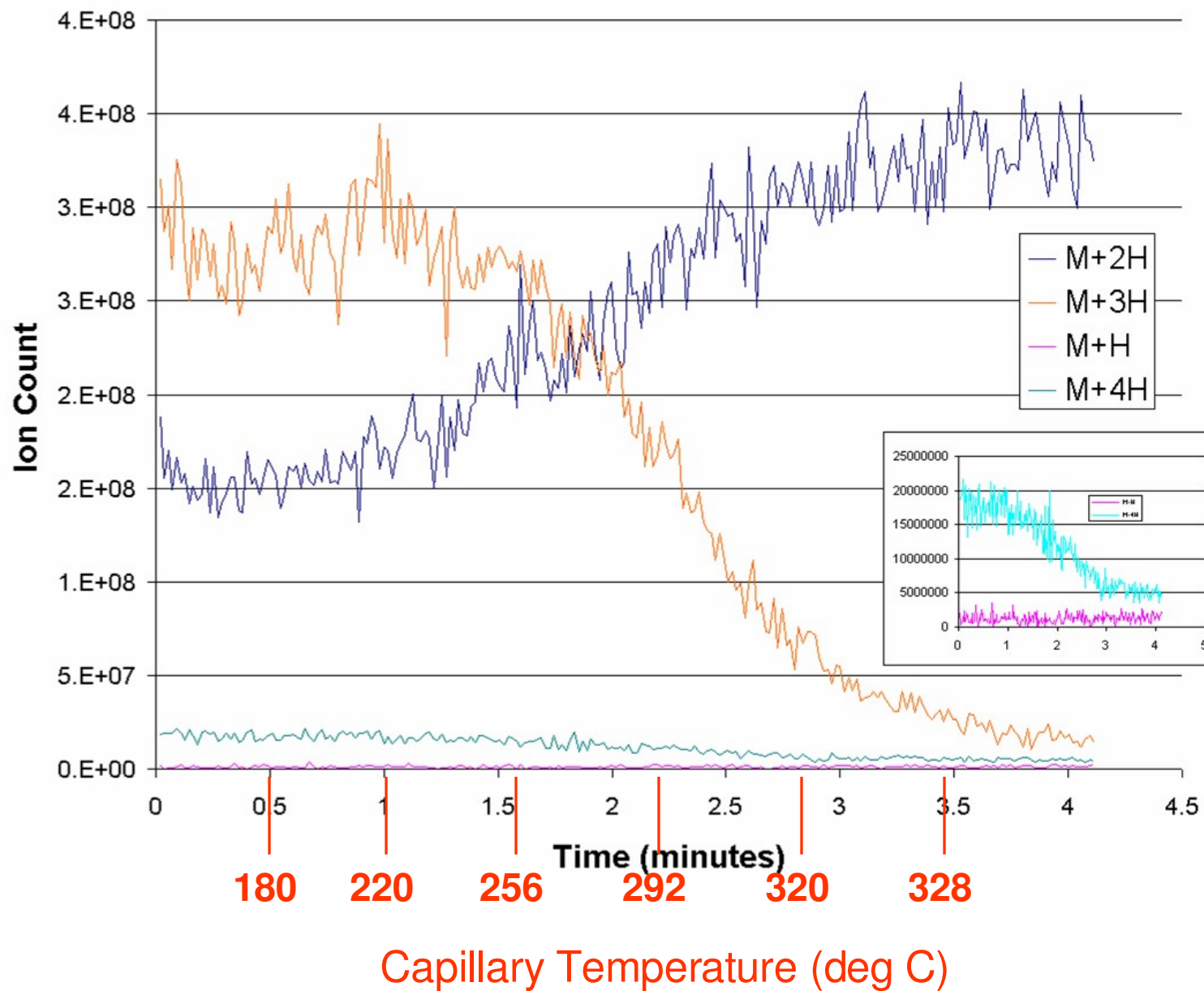d = after detailed statistical analysis

# Disulfide Bonds

Expected fraction: 1 - (1-freq)^n
[Rolling one "6" in **n** rolls is $1 - (5/6)^n$]



Fraction of Peptides with 1 or more Amino Acid (Peng080103)



Fraction of Peptides with 1 or more Amino Acid (Peng080503_DTT)

Using RasMol amino acid color scheme

# Charge State vs. Capillary Temperature

# Acknowledgments

- Dr. Edward Marcotte
- Dr. Klaus Linse
- Dr. Maria Person
- Dr. Aleksey Nakorshevskiy
- Dr. Rong Wang
- Dr. Peng Lu
- Zhihua Li