# CS4220: Knowledge Discovery Methods for Bioinformatics
# Unit 3: Proteomic Profiling

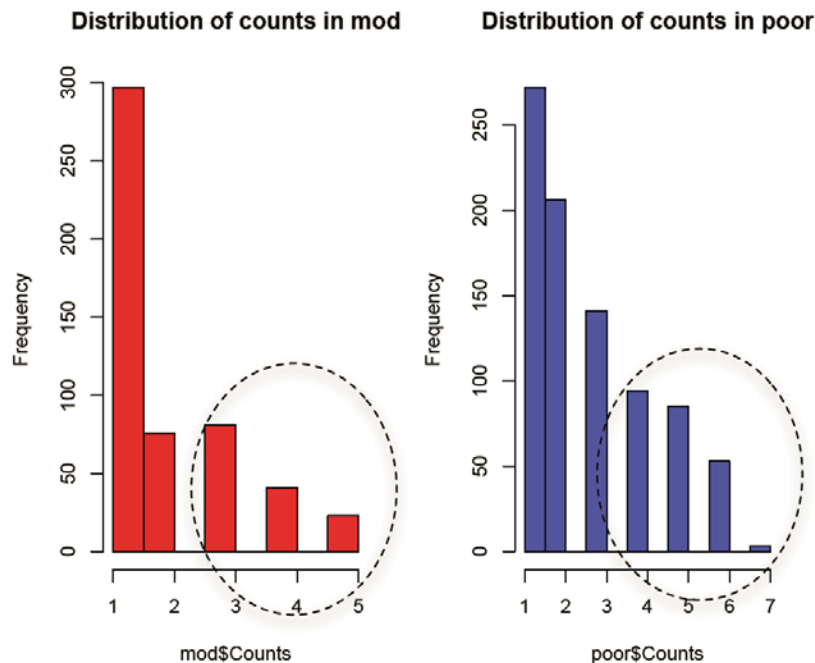**Wong Limsoon**

# Delivering more powerful proteomic profile analysis

- **Basic proteomic profile analysis**

- **Common issues**

- **Improving coverage**

- **Improving consistency**
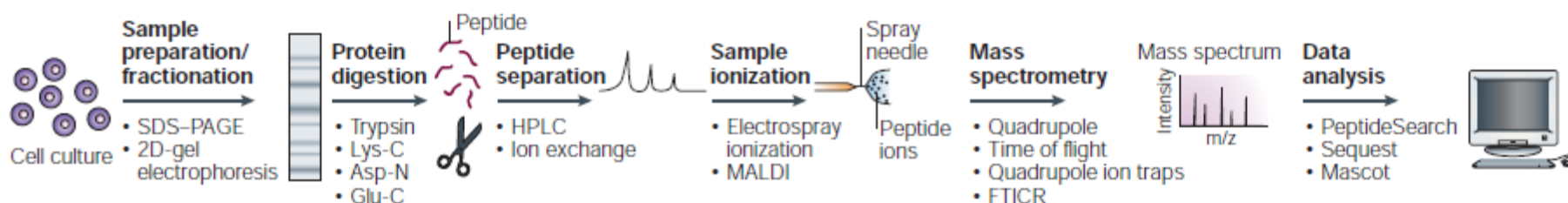
- **Finally, a quantum leap?**

# BASIC PROTEOMIC PROFILE ANALYSIS

# Typical proteomic MS expt



Figure 1 | **The mass-spectrometry/proteomic experiment.** A protein population is prepared from a biological source — for example, a cell culture — and the last step in protein purification is often SDS–PAGE. The gel lane that is obtained is cut into several slices, which are then in-gel digested. Numerous different enzymes and/or chemicals are available for this step. The generated peptide mixture is separated on- or off-line using single or multiple dimensions of peptide separation. Peptides are then ionized by electrospray ionization (depicted) or matrix-assisted laser desorption/ionization (MALDI) and can be analysed by various different mass spectrometers. Finally, the peptide-sequencing data that are obtained from the mass spectra are searched against protein databases using one of a number of database-searching programmes. Examples of the reagents or techniques that can be used at each step of this type of experiment are shown beneath each arrow. 2D, two-dimensional; FTICR, Fourier-transform ion cyclotron resonance; HPLC, high-performance liquid chromatography.

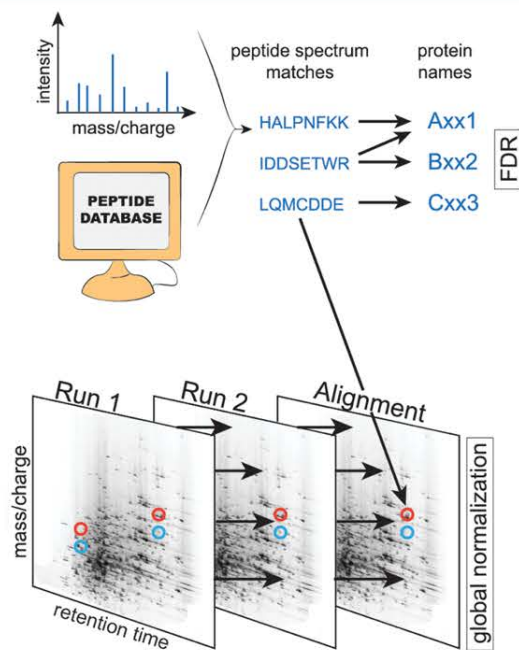See also http://www.slideshare.net/joachimjacob/bits-introduction-to-mass-spec-data-generation

Steen & Mann. The ABC's and XYZ's of peptide sequencing.
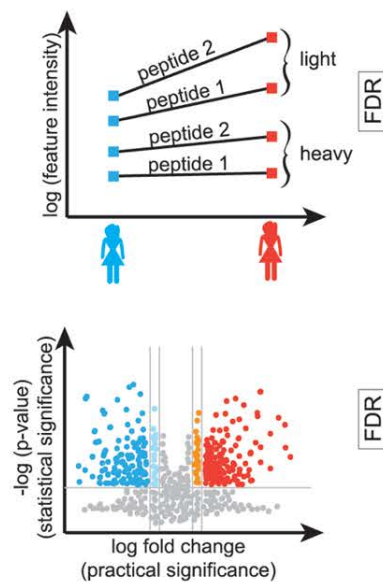*Nature Reviews Molecular Cell Biology*, 5:699-711, 2004

# Diagnosis using proteomics



Kall and Vitek, *PLoS Comput Biol ,* 7(12): e1002277, 2011
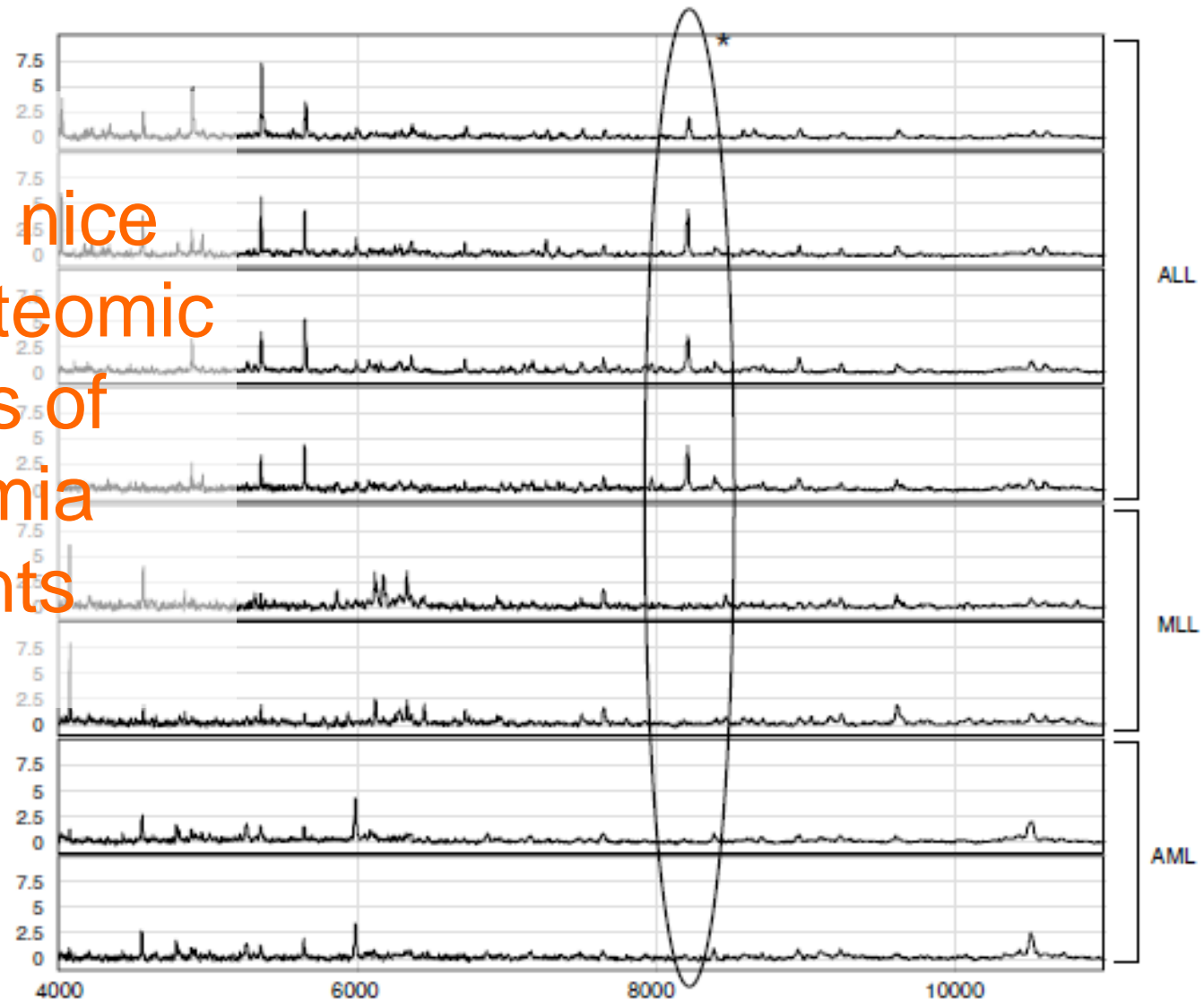
A rather nice set of proteomic profiles of leukemia patients



**Figure 1**  Spectra from SELDI-TOF MS analysis of REH, 697, MV4;11, and Kasumi cell lines. Protein (4 μg) from each cell type was analyzed on SAX2 ProteinChip® Arrays. ALL cell lines shown are REH and 697, the MLL cell line is MV4;11, and the AML cell line is Kasumi. The asterisk indicates the differentially expressed protein at 8.3 kDa.
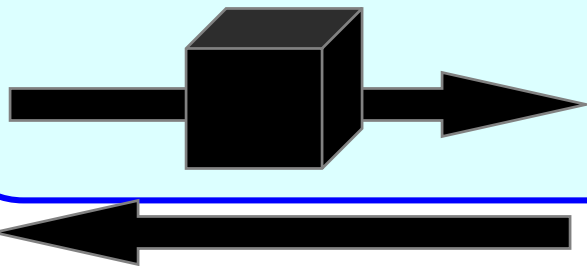
Hegedus et al. Proteomic analysis of childhood leukemia. Leukemia, 19:1713-1718, 2005

# Protein identification by mass spec



**Step 1:**

MS/MS instrument

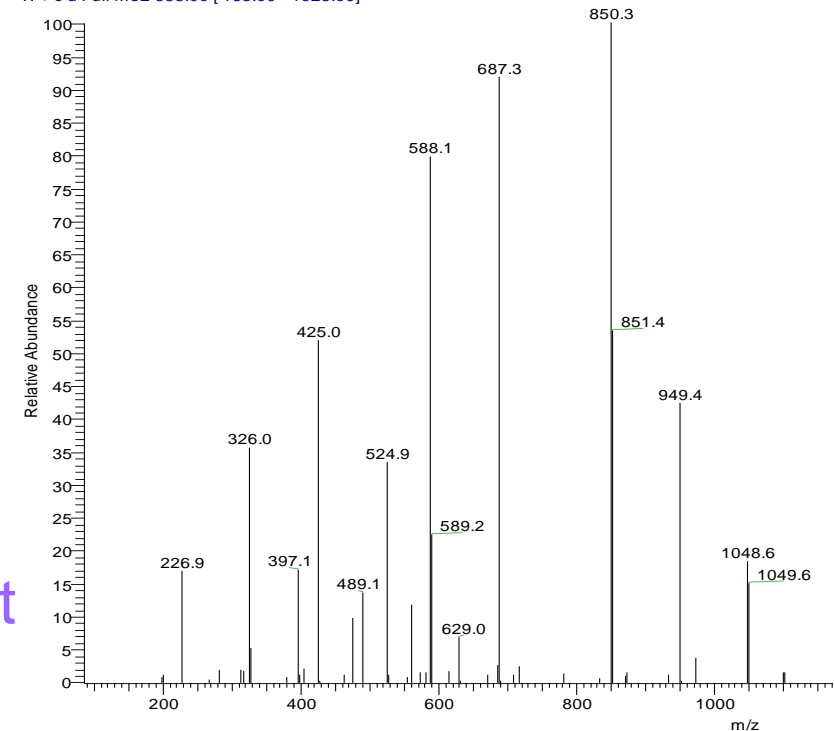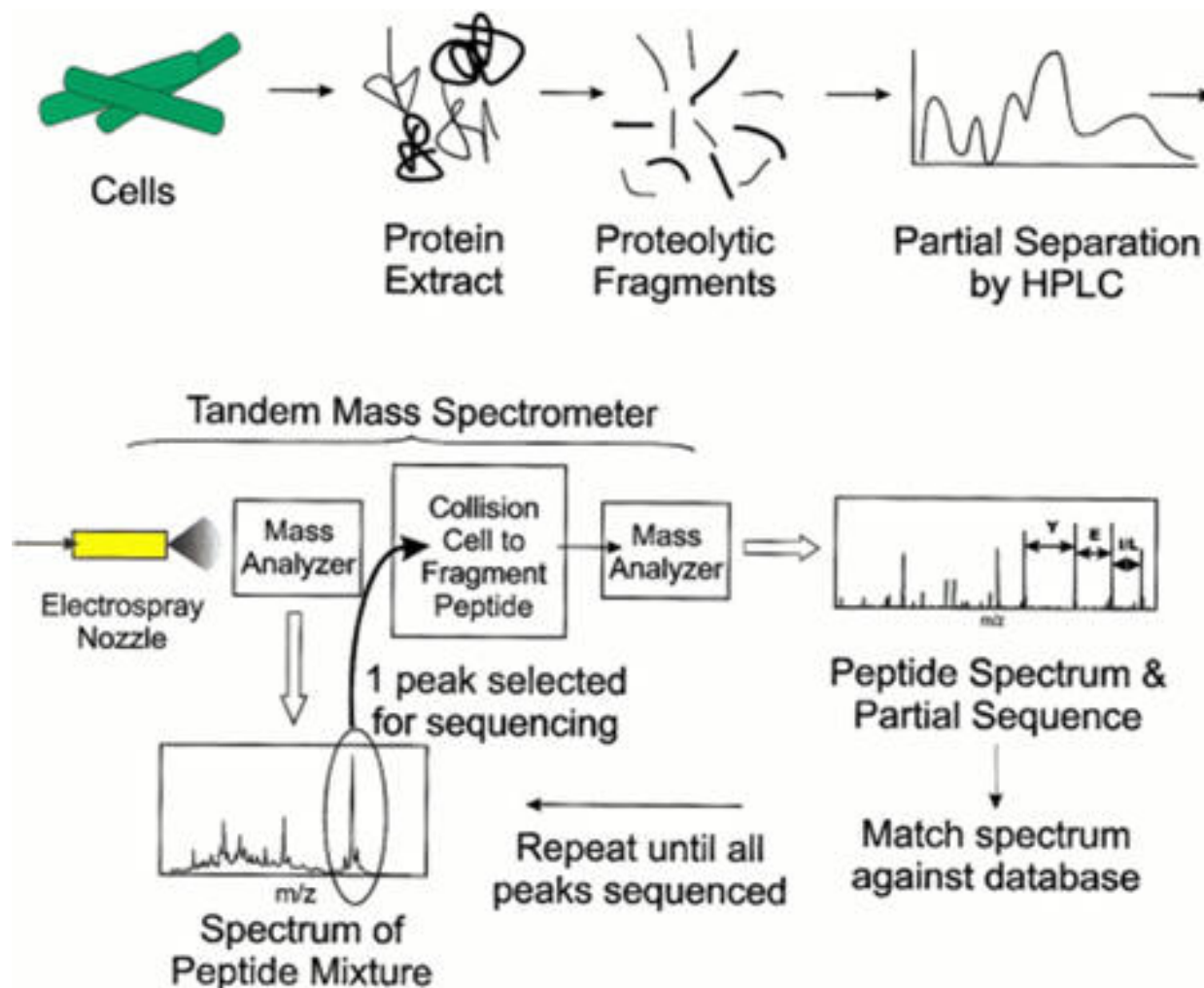**Sequence**

Database search
- Sequest, Mascot, InSpect

*de Novo* interpretation
- Lutefisk, Peaks, PepNovo

S#: 1708   RT: 54.47   AV: 1   NL: 5.27E6
T: + c d Full ms2 638.00 [ 165.00 - 1925.00]

# Tandem mass spectrometry

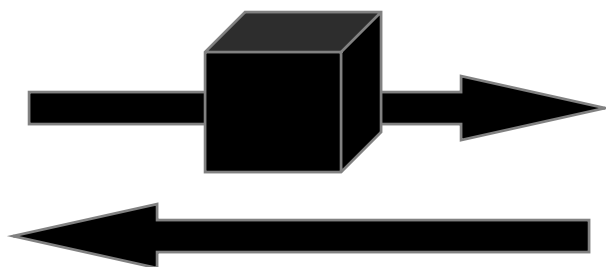# Breaking protein into peptides, and peptides into fragment ions

- **Proteases, e.g. trypsin, break a protein into peptides**

- **Tandem mass spectrometer further breaks the peptides down into fragment ions and measures the mass of each piece**

- **Mass spectrometer accelerates the fragmented ions; heavier ions accelerate slower than lighter ones**

- **Mass spectrometer measures mass/charge ratio of an ion**

# Peptide identification by mass spec

**Sequence**

MS/MS instrument

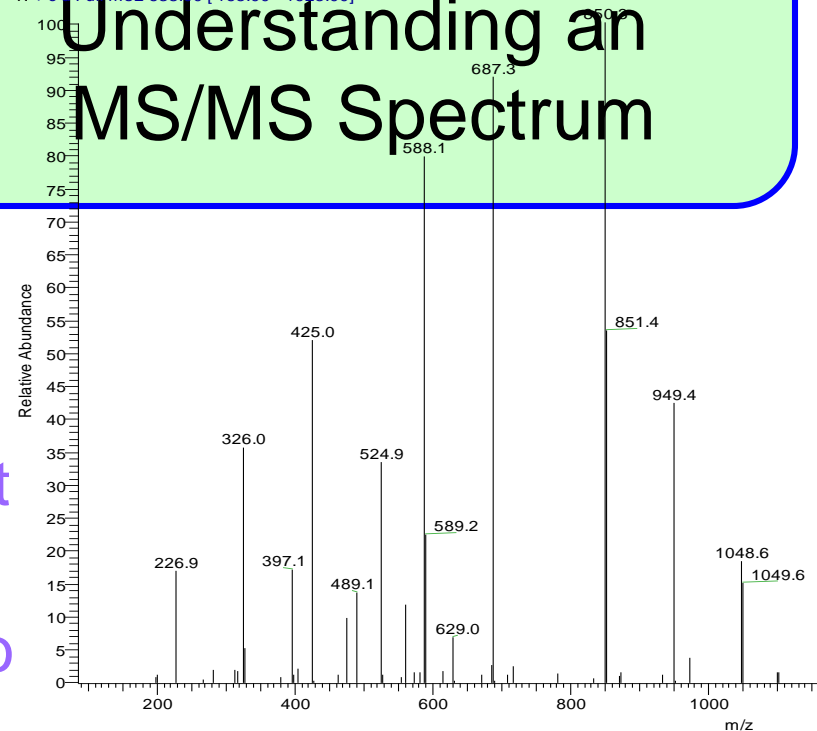Database search
- Sequest, Mascot, InSpect

*de Novo* interpretation
- Lutefisk, Peaks, PepNovo

**Step 2: Understanding an MS/MS Spectrum**

S#: 1708   RT: 54.47   AV: 1   NL: 5.27E6
T: + c d Full ms2 638.00 [ 165.00 - 1925.00]

Relative Abundance

226.9
326.0
397.1
425.0
489.1
524.9
588.1
589.2
629.0
687.3
851.4
949.4
1048.6
1049.6

200   400   600   800   1000

m/z

# Peptide fragmentation

Collision-Induced Dissociation

$$H...-HN-CH-CO \quad . \quad . \quad . \quad NH-CH-CO-NH-CH-CO-...OH$$

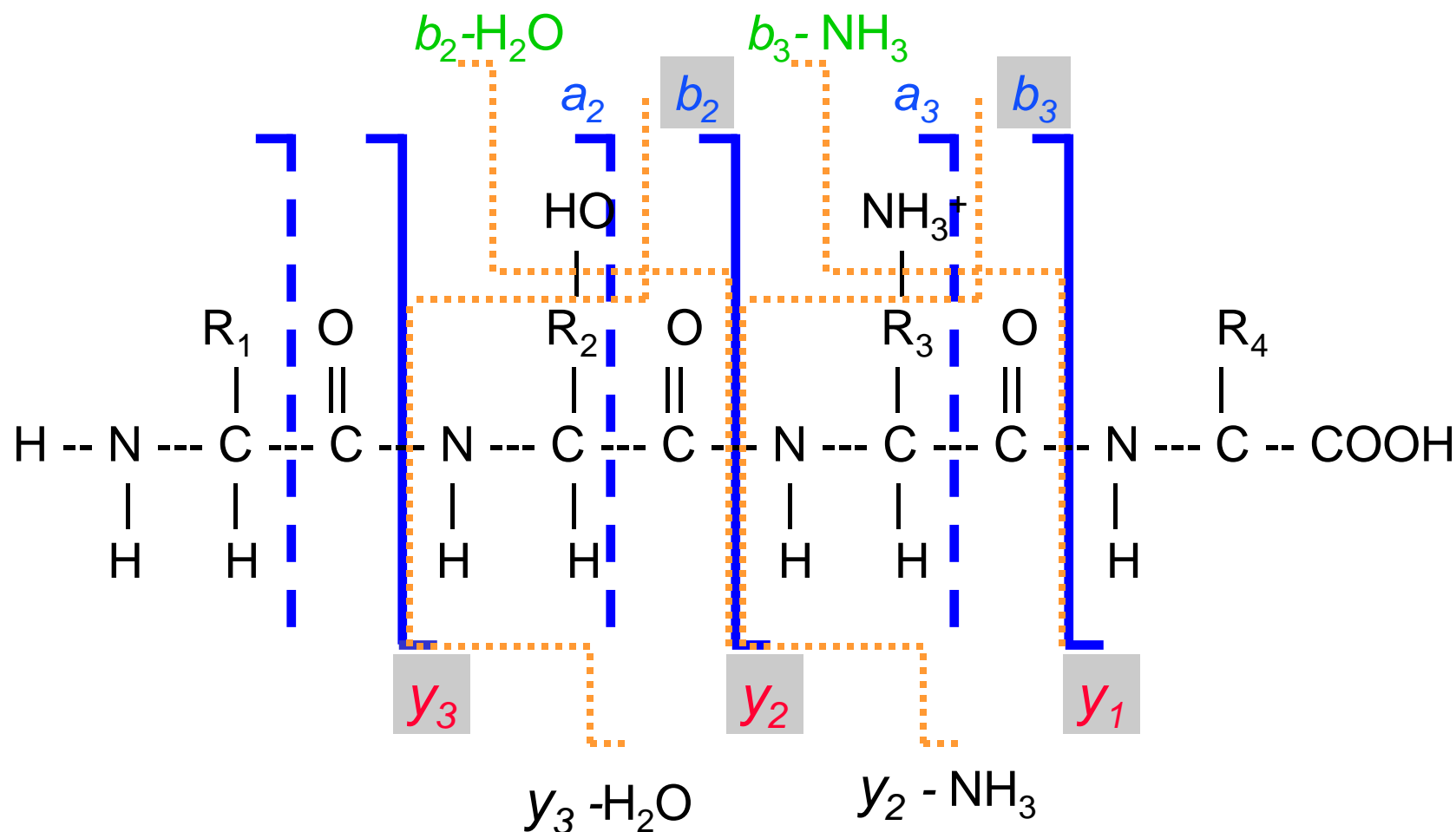$$R_{i-1} \qquad\qquad R_i \qquad\qquad R_{i+1}$$

$H^+$

Prefix Fragment                    Suffix Fragment

- **Peptides tend to fragment along the backbone**
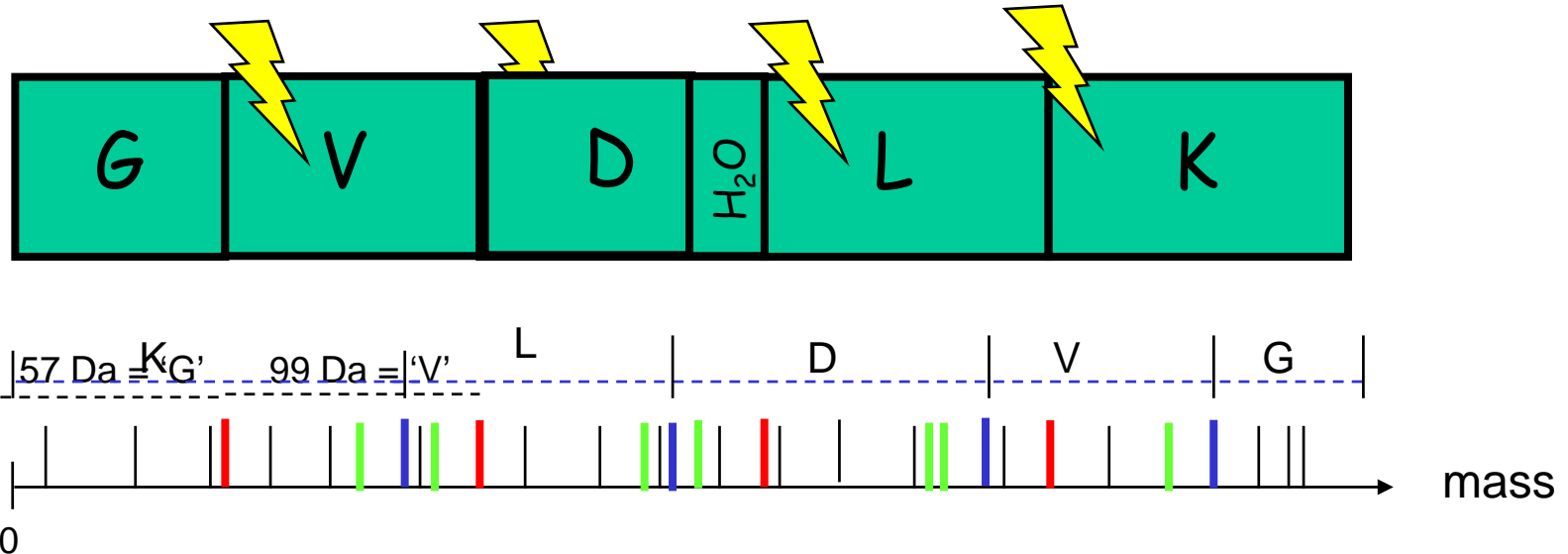- **Fragments can also loose neutral chemical groups like $NH_3$ and $H_2O$**

# … and fragments due to neutral losses

# Mass spectra



- **The peaks in the mass spectrum:**

    – Prefix and Suffix Fragments

    – Fragments with neutral losses ($-H_2O$, $-NH_3$)

    – Noise and missing peaks

Bafna & Edwards. "On de novo interpretation of tandem mass spectra for peptide identification". RECOMB 2003, pp. 9-18

# Example MS/MS spectrum

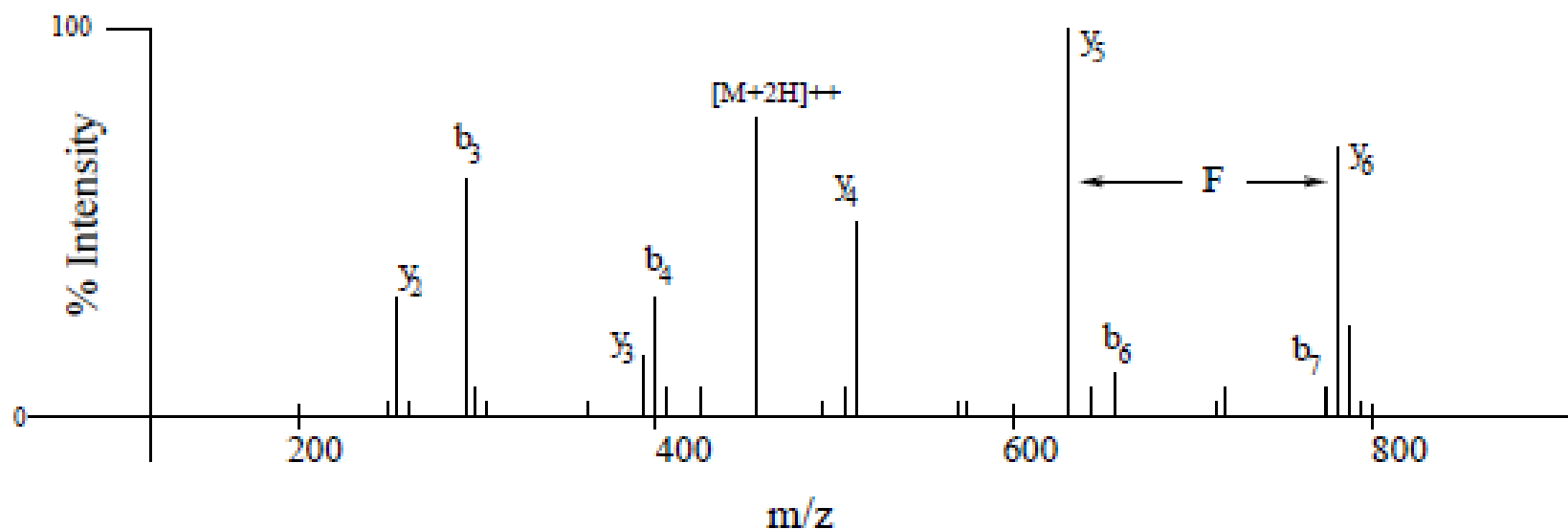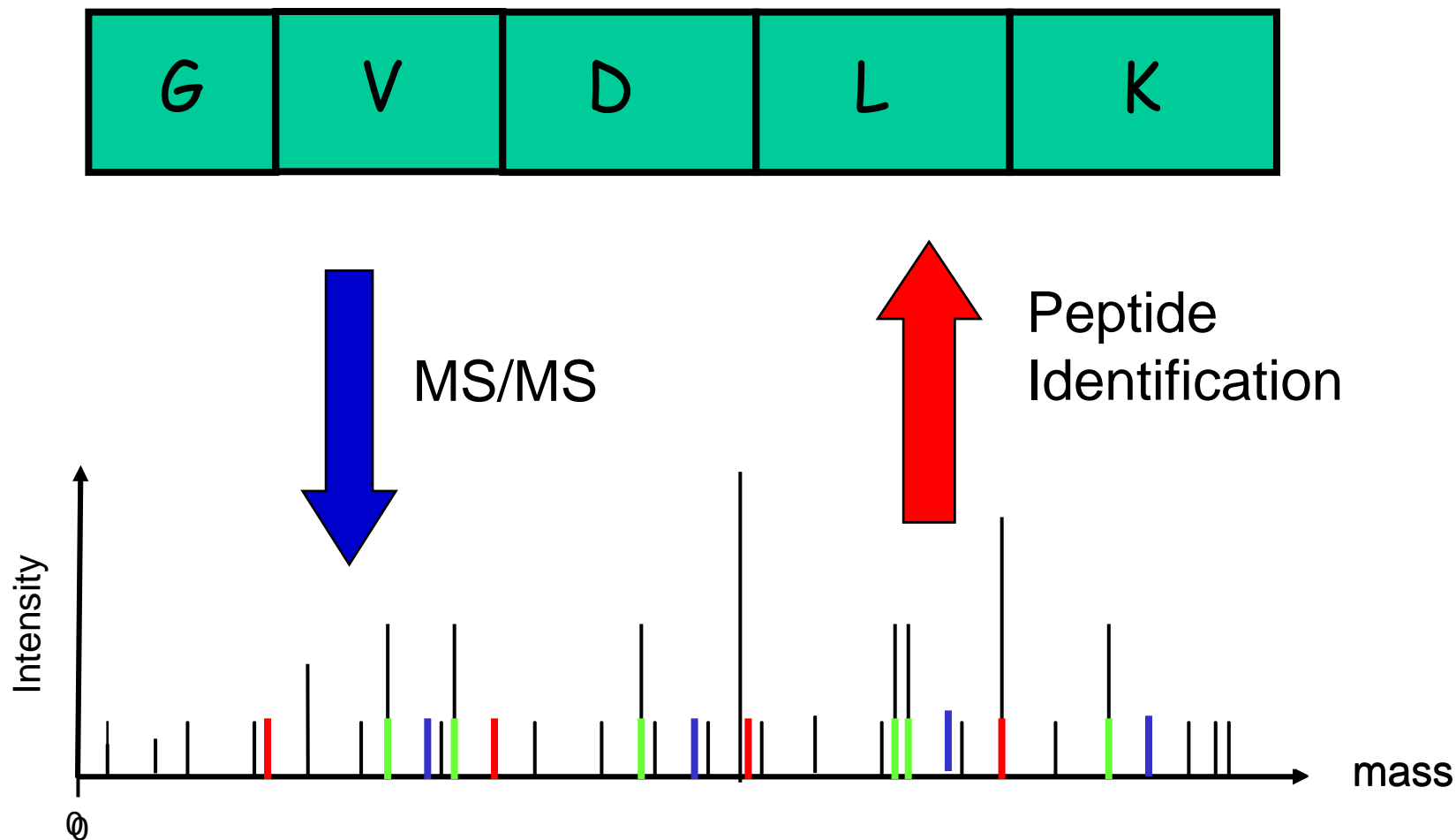| 88 | 145 | 292 | 405 | 534 | 663 | 778 | 924 | b-ions |
|----|-----|-----|-----|-----|-----|-----|-----|--------|
| S | G | F | L | E | E | D | K | |
| 924 | 837 | 780 | 633 | 520 | 391 | 262 | 141 | y-ions |



Figure 2: MS/MS spectrum for peptide SGFLEEDK.

# Protein identification with MS/MS

# Peptide identification by mass

**S e q u e n c e**

## MS/MS instrument



S#: 1708   RT: 54.47   AV: 1   NL: 5.27E6
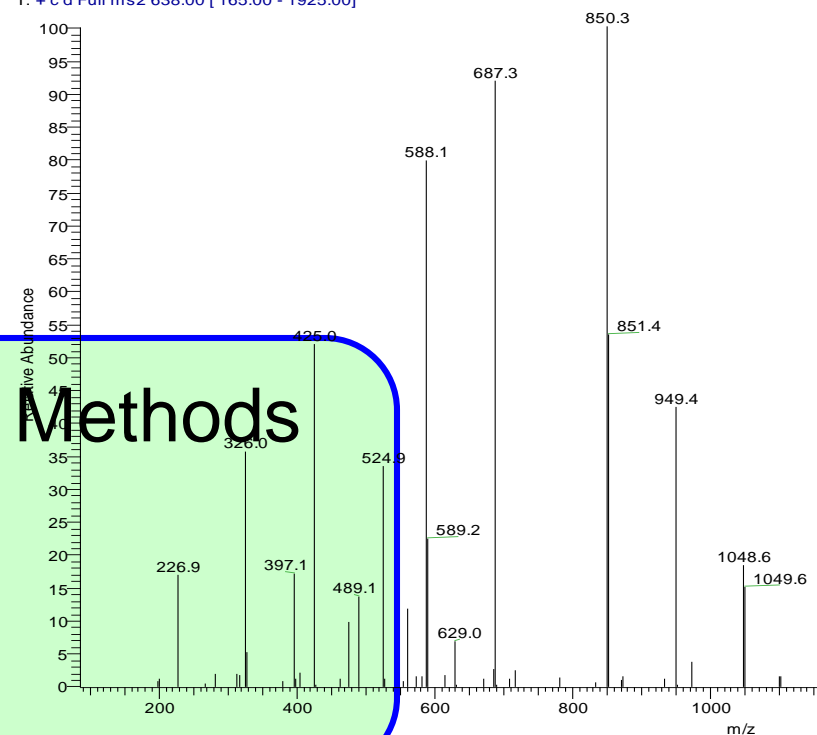T: + c d Full ms2 638.00 [ 165.00 - 1925.00]

**Step 3: Computational Methods**

Database search
  Sequest, Mascot
*de Novo* interpretation
  Lutefisk, Peaks, PepNovo

# Database search algorithms

- **Database search**
  - Used for spectrum from known peptides
  - Rely on completeness of database

- **General Approach**
  - Match given spectrum with known peptide
  - Enhanced with advanced statistical analysis and complex scoring functions

- **Methods**
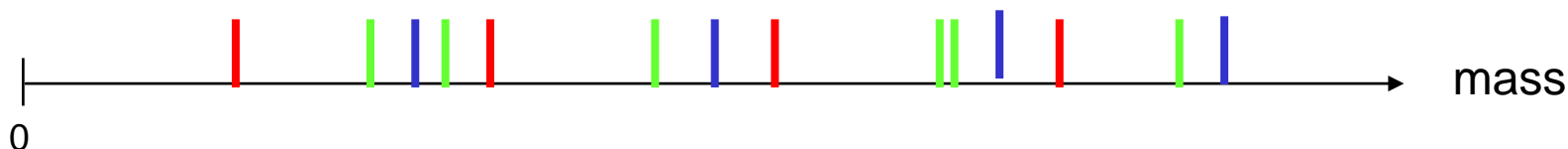  - SEQUEST, MASCOT, InsPecT, Paragon

# Theoretical spectrum for a peptide

- **Given this peptide**

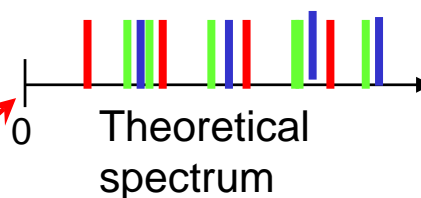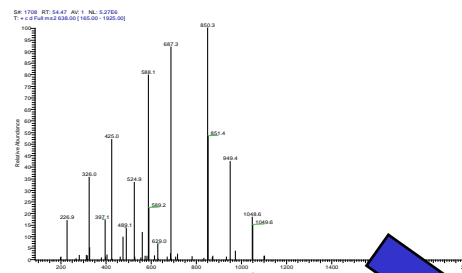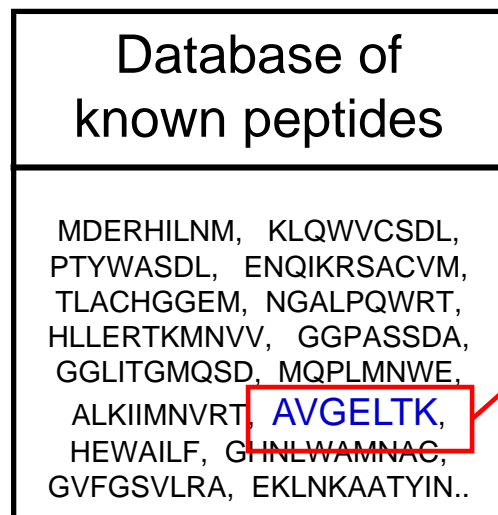| G | V | D | L | K |
|---|---|---|---|---|

- **Its theoretical spectrum is**



mass

0

- **Theoretical spectrum is dependent on**
  - Set of ion-types considered
  - Larger if multi-charge ions are considered

# Database search algorithm

**Database Search**



**Database of known peptides**

MDERHILNM, KLQWVCSDL, PTYWASDL, ENQIKRSACVM, TLACHGGEM, NGALPQWRT, HLLERTKMNVV, GGPASSDA, GGLITGMQSD, MQPLMNWE, ALKIIMNVRT, AVGELTK, HEWAILF, GHNLWAMNAC, GVFGSVLRA, EKLNKAATYIN..

**Theoretical spectrum**

**Match**

**Matching Score for this peptide**

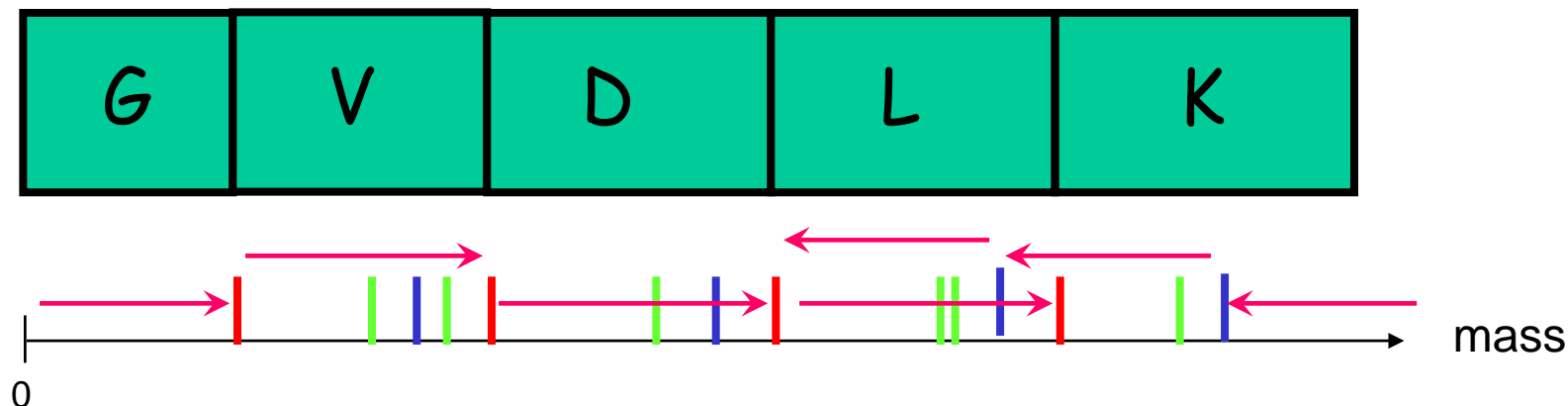**Repeat for all the peptides in the Database**

# De novo sequencing algorithms

- **Given a spectrum**
  - Build a spectrum graph
  - Peptides are paths in this graph
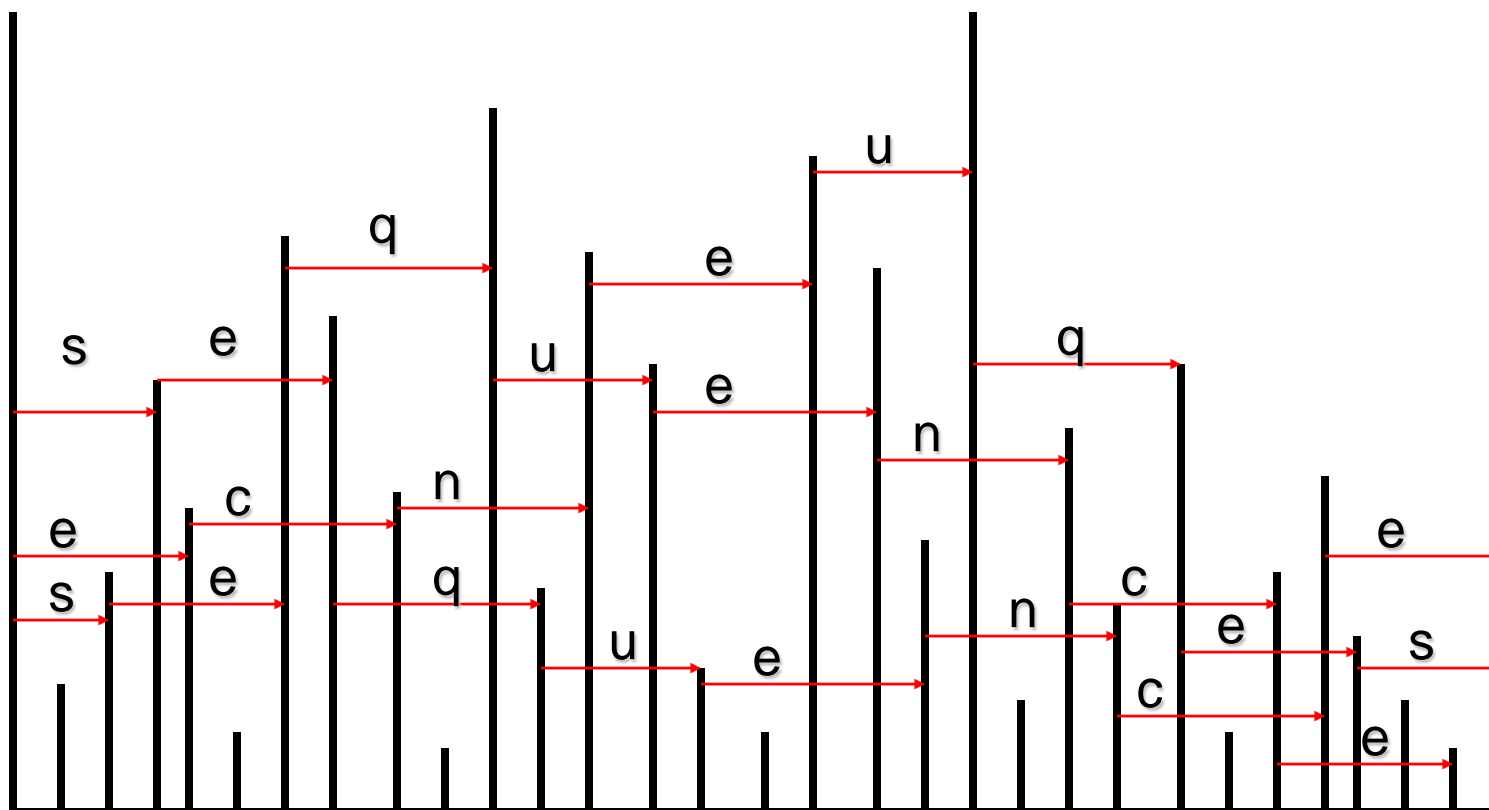  - Find the best path

# Spectrum graph for a peptide



- **Connect peaks together**

  – If their mass difference = mass of an amino acid

- **Theoretical spectrum is dependent on**
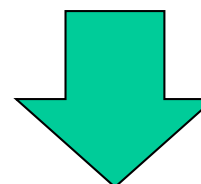
  – Set of ion-types considered

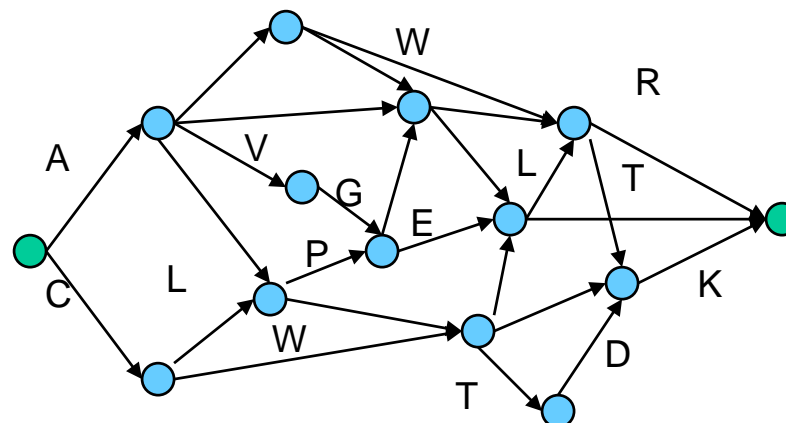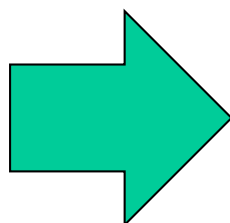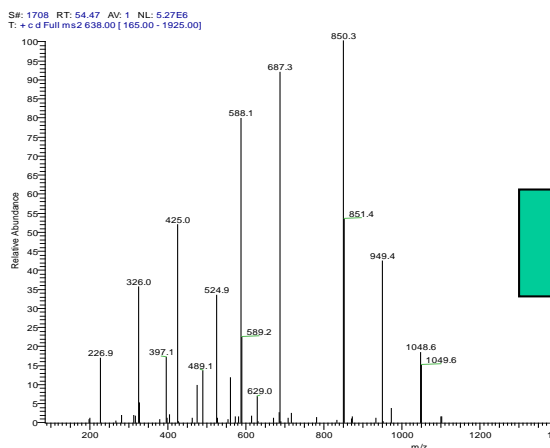  – Larger if multi-charge ions are considered

# Building a graph from a spectrum

# De novo sequencing algorithms



Find longest directed acyclic path

AVGELTK

# De novo vs. database search
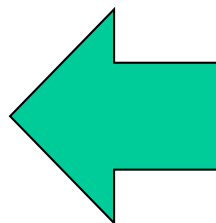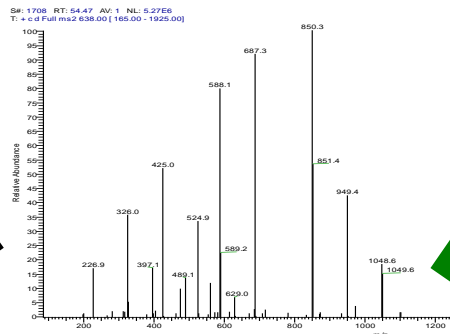


**Database Search**

**De Novo**

Database of known peptides

MDERHILNM, KLQWVCSDL, PTYWASDL, ENQIKRSACVM, TLACHGGEM, NGALPQWRT, HLLERTKMNVV, GGPASSDA, GGLITGMQSD, MQPLMNWE, ALKIIMNVRT, AVGELTK, HEWAILF, GHNLWAMNAC, GVFGSVLRA, EKLNKAATYIN..

Database of all peptides $\approx 20^n$

AAAAAAAA, AAAAAAAC, AAAAAAAD, AAAAAAAE, AAAAAAAG, AAAAAAAF, AAAAAAAH, AAAAAAAY,

AVGELTI, AVGELTK, AVGELTL, AVGELTM,

YYYYYYYS, YYYYYYYT, YYYYYYYV, YYYYYYYY

**AVGELTK**

# De novo vs. database search: A paradox

- **The database of all peptides is huge ≈ $O(20^n)$**
- **The database of all known peptides is much smaller ≈ $O(10^8)$**

- **However, de novo algorithms can be much faster, even though their search space is much larger!**
  - A database search scans all peptides in the search space to find best one
  - De novo eliminates the need to scan all peptides by modeling the problem as a graph search

# Protein identification

- **After all the peptides have been identified, they are grouped into protein identifications**

- **Peptide scores are added up to yield protein scores**

- **Confidence of a particular peptide identification increases if other peptides identify the same protein and decreases if no other peptides do so**

- **Protein identifications based on single peptides should only be allowed in exceptional cases**

Steen & Mann. The ABC's and XYZ's of peptide sequencing.
*Nature Reviews Molecular Cell Biology*, 5:699-711, 2004

# COMMON ISSUES IN PROTEOMIC PROFILE ANALYSIS

# Proteomics vs transcriptomics

- **Proteomic profile**
  - Which protein is found in the sample
  - How abundant it is

- **Similar to gene expression profile. So typical gene expression profile analysis methods can be applied in theory…**

- **Key differences**
  - Profiling
    - **Complexity: 20k genes vs 500k proteins**
    - **Dynamic range: > 10 orders of magnitude in plasma. Proteins cannot be amplified**
  - Analysis
    - **Much fewer features**
    - **Difficult to reproduce**
    - **Much fewer samples**
    - **Unstable quantitation**

# Peptide & protein identification by MS is still far from perfect

- **"… peptides with low scores are, nevertheless, often correct, so manual validation of such hits can often 'rescue' the identification of important proteins."**
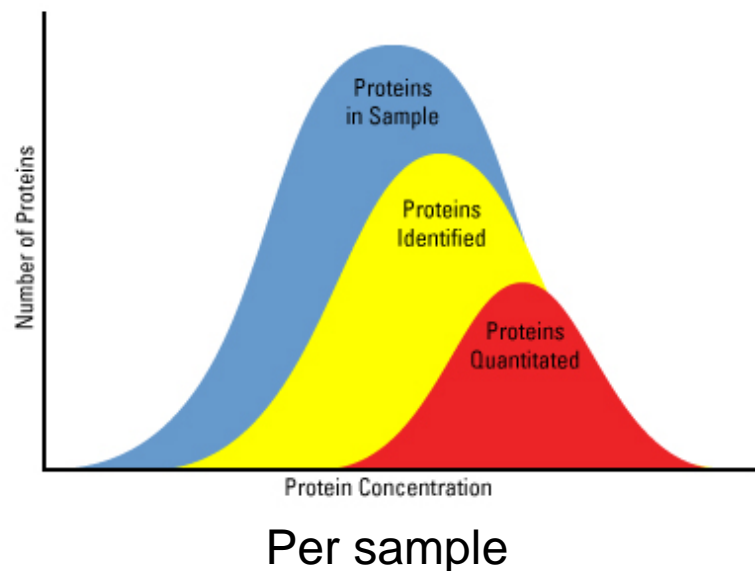
  Steen & Mann. **The ABC's and XYZ's of peptide sequencing**. *Nature Reviews Molecular Cell Biology*, 5:699-711, 2004
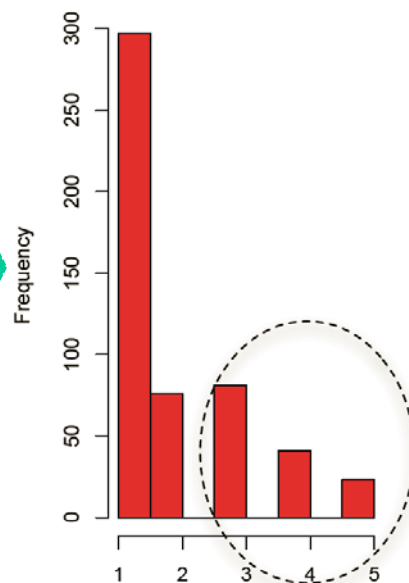
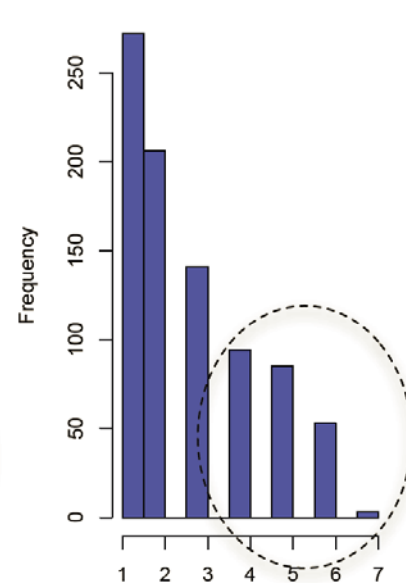# Issues in proteomics: Coverage and consistency

**Technical incompleteness**     **How it affects real data**



Per sample

Only 25 out of 800+ proteins are common to all 5 mod-stage HCC patients!
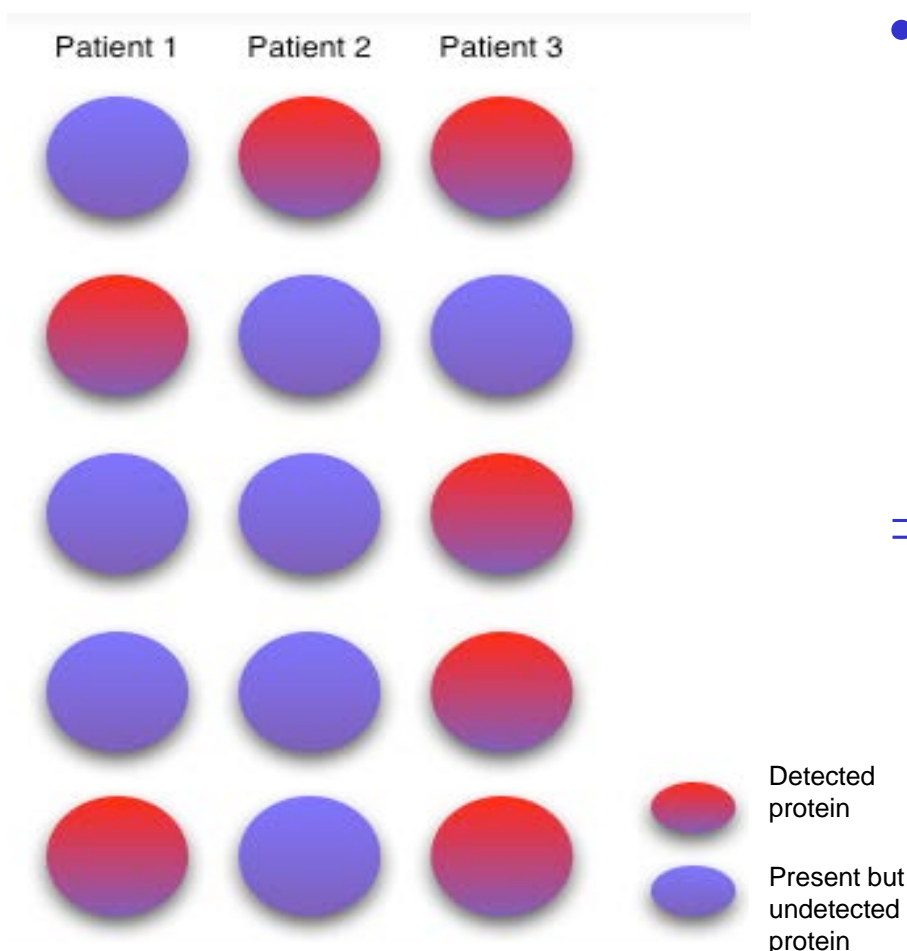
# USING PROTEIN COMPLEXES IN PROTEOMICS: BASIC IDEAS

# A postulate and some math

- **Postulate: The chance of a protein complex being present in a sample is proportional to the fraction of its constituent proteins being correctly reported in the sample**

- **Suppose proteomics screen has 75% reliability; a complex comprises proteins A, B, C, D, E; and screen reports A, B, C, D only**

$\Rightarrow$ **Complex has 60% (= 0.75 * 4 / 5) chance to be present**

$\Rightarrow$ **The unreported protein E also has $\geq$ 60% chance to be present, as presence of the complex implies presence of all its constituents**
  - $\Rightarrow$ improving coverage

$\Rightarrow$ **Each of the reported proteins (A, B, C, and D) individually has 90% (= 100% * 0.6 + 75% * 0.4) chance of being true positive, whereas a reported protein that is isolated has a lower 75% chance of being true positive**
  - $\Rightarrow$ removing noise

# An intuition



Patient 1   Patient 2   Patient 3

Detected protein

Present but undetected protein

- **Suppose the failure to form a protein complex causes a disease**
  - If any component protein is missing, the complex can't form
- ⇒ **Diff patients suffering from the disease can have a diff protein component missing**
  - Construct a profile based on complexes?

# Reference complexes

- **In this lecture, human complexes (of size at least 5) from CORUM are used as reference complexes**

- **It is possible to use subnets generated from pathway and PPI databases. However these such subnets vary significantly depending on network databases and subnet-generation algo used**

So I do not consider these…

# IMPROVING COVERAGE IN PROTEOMIC PROFILES

Guo et al. *Nature Medicine*, *21*, 407, 2015

# Lots of missing values in real proteomics datasets

Missing values are not due mostly to low-abundance proteins



**Figure 1.**

Average $\log_{10}$ intensity as measured by peptide peak area in the control group versus fraction of missing values and peptide counts associated with bins corresponding to the fraction of missing data comparing phenotypes and exposures for datasets from (A) human plasma and (B) mouse lung. The control group for the human plasma is the normal glucose tolerant (NGT) samples, and the sham group for the mouse lung is the regular weight mice with no lipopolysaccharide (LPS) exposure. The vertical red line represents median average intensity, and the horizontal red line represents the point that 50% of the values are missing.

Webb-Robertson, *JPR*, 14(5):1993-2001, 2015

CV(RMSE) ~ 20%
at 75th percentile

**Current imputation methods don't work very well**



**Figure 2.**
Boxplot of the average $\log_{10}$ CV(RMSE) for the imputed dilution series datasets (Table 1) at the (A) peptide and (B) protein levels. The lower line represents the 25th percentile, the upper line of the box represents the 75th percentile, and the inner line corresponds to the median $\log_{10}$ CV(RMSE).

# FCS

- **Rescue undetected proteins from high-scoring protein complexes**

- **Goh et al. Comparative network-based recovery analysis and proteomic profiling of neurological changes in valporic acid-treated mice. *JPR*, 12(5):2116-2127, 2013**

Basic idea: Let A, B, C, D and E be the 5 proteins that function as a complex and thus are normally correlated in their expression. Suppose only A is not detected and all of B–E are detected. Suppose the screen has 50% reliability. Then, A's chance of being false negative is 50%, & the chance of B–E all being false positives is $(50\%)^4$=6%. Hence, it is almost 10x more likely that A is false negative than B–E all being false positives.

Observed

Empirical Overlap Ratios

P(Significant Enrichment)?

Observed

Identified Proteins

# Other methods for rescuing missing proteins

- **CEA**
  - Generate cliques from PPIN
  - Rescue missing proteins from cliques containing lots of high-confidence proteins
  - Li et al. Network-assisted protein identification and data interpretation in shotgun proteomics. *Mol. Syst. Biol.*, 5:303, 2009

- **MaxLink**
  - Map high-confidence proteins ("seeds") to PPIN
  - Rescue proteins that interact many seeds but few non-seeds
  - Goh et al. *Int J Bioinformatics Research and Applications*, 8(3/4):155-170, 2012

- **PEP**
  - Map high-confidence proteins to PPIN
  - Extract neighbourhood & predict protein complexes using CFinder
  - Rescue undetected proteins from high-ranking predicted complexes
  - Goh et al. A Network-based pipeline for analyzing MS data---An application towards liver cancer. *J. Proteome Research*, 10(5):2261-2272, 2011

# iTRAQ experiment

- **Valporic acid (VPA)-treated mice vs control**
  - VPA or vehicle injected every 12 hours into postnatal day-56 adult mice for 2 days
  - Role of VPA in epigenetic remodeling

- **MS was scanned against IPI rat db in round #1**
  - 291 proteins identified
- **MS was scanned against UniProtkb in round #2**
  - 498 additional proteins identified

- **All recovery methods ran on round #1 data and the recovered proteins checked against round #2**

Moderate level of agreement of reported proteins between various recovery methods



FCS (Real Complexes)

# Performance comparison

| Method | Novel Suggested Proteins | Recovered proteins | Recall | Precision |
|---|---|---|---|---|
| PEP | 1037 | 158 | 0.317 | 0.152 |
| Maxlink | 822 | 226 | 0.454 | 0.275 |
| FCS (predicted) | 638 | 224 | 0.450 | 0.351 |
| FCS (complexes) | 895 | 477 | 0.958 | 0.533 |

- **Looks like running FCS on real complexes is able to recover more proteins and more accurately**

# SWATH experiment

- **If there are technical replicates, they should have reported the same proteins. So we can run FCS on one replica, and see whether the predicted missing proteins show up in other replicas**

- **If there are multiple biological replicates (i.e. patients of the same phenotype), we can run FCS on one of them, and check on the others**

- **Proteomics data used: Renal cancer**
  - Guo et al. *Nature Medicine*, 21(4):407-413, 2015
  - 6 pairs of normal vs cancer ccRCC tissues
  - SWATH in duplicates

~20% of predicted missing proteins are supported by $\geq$1 reported peptide in the screen

# ~20% of predicted missing proteins are supported by ≥1 reported peptide in the replicate

# But ~25% of predicted missing proteins are supported by peptides in the screen or replicate

**C**                                                                              st union of self and other replicate)

# Time for Exercise #1

- **About 25% FCS-predicted missing proteins are supported by peptides in screen/replicate. Can we do better? In particular, suggest some ways to rank missing proteins, so that those which are more likely to be really present are ranked higher**

# IMPROVING CONSISTENCY IN PROTEOMIC PROFILE ANALYSIS

# Proteomic profiles generally not consistent, even for technical replicates

- **A human kidney tissue**
    - Guo et al. *Nature Medicine*, 21(4):407-413, 2015
    - Digested in quadruplicates
    - Analyzed in triplicates

- **Clustering by proteins**
    - Correlation betw replicates is not good (~0.4)
    - Technical replicates of the same biological replicate are not tightly clustered

Goh et al. **Quantitative proteomics signature profiling based on network contextualization**. *Biology Direct*, 10:71. 2015

# qPSP

- **In a sample, assign weight to proteins as follow**
  - Most-abundant 10% of proteins, wt = 1
  - Proteins at 10-12.5%, wt = 0.8
  - Proteins at 12.5-15%, wt = 0.6
  - Proteins at 15-17.5%, wt = 0.4
  - Proteins at 17.5-20%, wt = 0.2
  - All other proteins, wt = 0



- **Hit rate of a complex C wrt a sample S is sum of the wt of proteins in C in S**
  - $score(C, Si) = \Sigma_{p \in C}\, fs(p, Si)\, /\, |C|$

- **Complex C is significant if $\{score(C, S_i)\mid S_i \in A\}$ is very different by t-test from $\{score(C, S_i)\mid S_i \in B\}$**

# Why qPSP is based on the most abundant (top 10-20%) proteins

# False-positive analysis

- **12 kidney controls were randomly assigned into two groups of equal size, and qPSP analysis was performed 10000 rounds**
- **For each round, # of significant clusters (5% FDR) was determined.**
- **Histograms showed that the false-positive rates were well within the expectation levels**
  - Sig Clusters Abs E-level: 19, Observed level: 16
  - Ratio Sig Clusters E-level: 0.05, Observed level: 0.04



Sig Clusters (Abs)

Ratio Sig Clusters

# Consistency of qPSP

- **Clustering of benchmarking control data based on protein complexes (i.e. qPSP)**

  - Correlation betw replicates is >0.95
    - **Cf. 0.4 based on proteins**
  - Technical replicates are better clustered

# Application to
# renal & colorectal cancers



**Fig. 3** qPSP strongly discriminates sample classes for renal cancer (**a**) and colorectal cancer (**b**). Clustered similarity maps at the top row showed specific and consistent segregation of non-cancer and cancer samples. The trees below the heatmaps are from bootstrap analysis (PVCLUST), which demonstrates that the discrimination between sample classes based on qPSP hit-rates is highly stable

# ESSNET: A QUANTUM LEAP?

# ESSNet, adapted for proteomics

- **Let $g_i$ be a protein in a given protein complex**
- **Let $p_j$ be a patient**
- **Let $q_k$ be a normal**

- **Let $\triangle_{i,j,k}$ = Expr($g_i$,$p_j$) − Expr($g_i$,$q_k$)**

- **Test whether $\triangle_{i,j,k}$ is a distribution with mean 0**

- **Null hypothesis is "Complex C is irrelevant to the difference between patients and normals, and the proteins in C behave similarly in patients and normals"**

- **No need to restrict to most abundant proteins**
- $\Rightarrow$ **Potential to reliably detect low-abundance but differential proteins**

Goh & Wong. **Advancing clinical proteomics via analysis based on biological complexes: A tale of five paradigms**. *Journal of Proteome Research*, 15(9):3167--3179, 2016

# Five methods to compare with

- **Network-based methods**
  - Hypergeometric enrichment (HE)
  - Direct group analysis (DG), similar to GSEA
  - qPSP, Goh et al*., Biology Direct*, 10:71, 2015
  - PFSNET, Goh & Wong, *JBCB,* 14(5):16500293, 2016
  - ESSNET, Lm et al., *JBCB*, 13(4):1550018, 2015 & Goh & Wong, *JPR*, 15(9):3167-3179 2016

- **Standard t-test on individual proteins (SP)**

# Simulated data

- **Simulated datasets from Langley and Mayr**
  - D.1.2 is from study of proteomic changes resulting from addition of exogenous matrix metallopeptidase (3 control, 3 test)
  - D2.2 is from a study of hibernating arctic squirrels (4 control, 4 test)

- **Both D1.2 and D2.2 have 100 simulated datasets, each with 20% significant features**
  - Effect sizes of these differential features are sampled from one out of five possibilities (20%, 50%, 80%, 100% and 200%), increased in one class and not in the other

- **Significant artificial complexes are constructed with various level of purity (i.e. proportion of significant proteins in the complex)**
  - Equal # of non-significant complexes are constructed as well

SP shows poor performance on simulated data.

Can network-based methods do better?



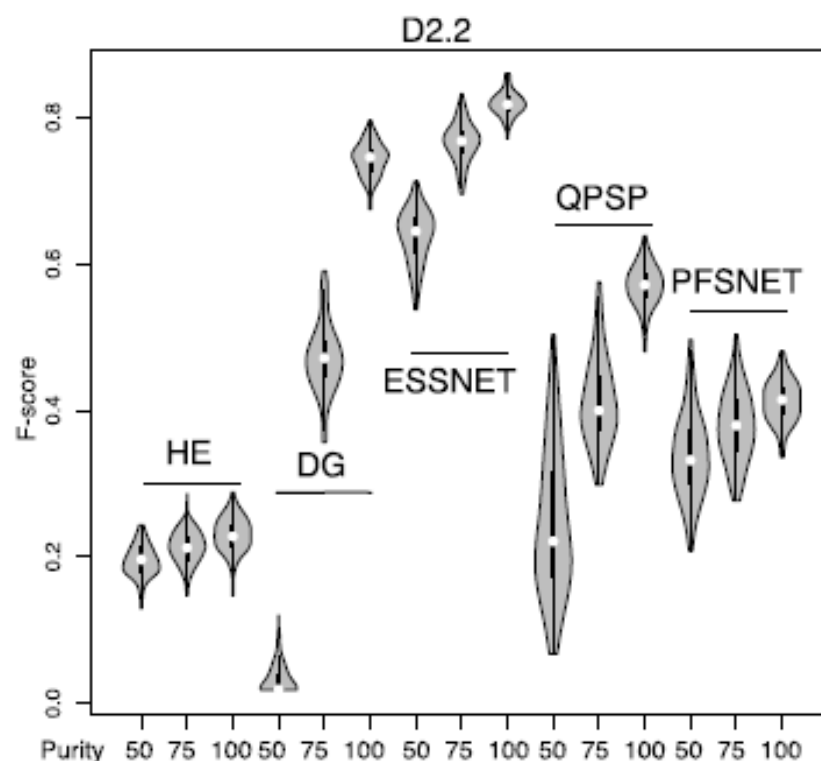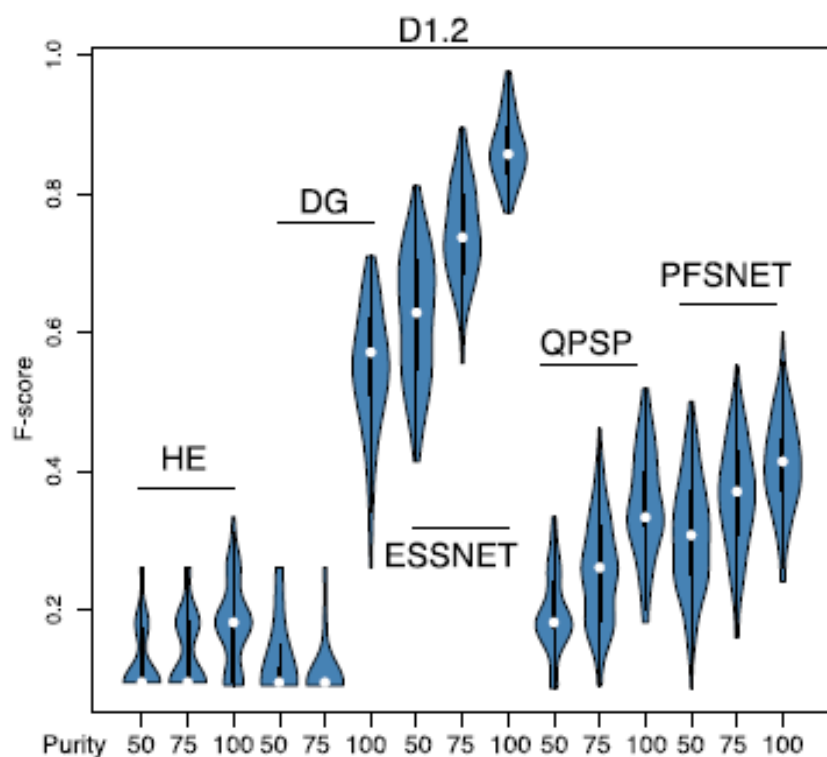**Supplementary Figure 1 Single protein (SP) precision-recall performance on D1.2.** The f-score (pink), precision (blue) and recall (green) shows that SP performs abysmally on simulated data. HE is shown next to SP as a reference.

# ESSNET shows excellent recall/precision on simulated data

Guo et al. *Nature Medicine*, 21(4):407-413, 2015

# Renal cancer control data (RCC)

- **12 runs originating from a human kidney tissue digested in quadruplicates and analyzed in triplicates**

- **Excellent for evaluating false-positive rates of feature-selection methods**
  - Randomly split the 12 runs into two groups. Report of any significant features between the groups must be false positives

Dash line corresponds to expected # of false positives at alpha 0.05 (~30 complexes)

All methods control false positives well

# Renal cancer data (RC)

- **12 samples are run twice so that we have technical replicates over 6 normal and 6 cancer tissues**

- **Excellent opportunity for testing reproducibility of feature-selection methods**
  - A good method should report similar feature sets between replicates

- **Can also test feature-selection stability**
  - Apply feature-selection method on subsamples and see whether the same features get selected

# ESSNET & PFSNET show excellent reproducibility

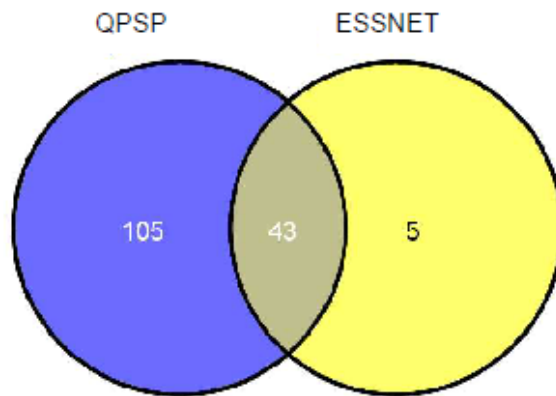| Number of terms | HE | DG | ESSNET | QPSP | PFSNET |
|---|---|---|---|---|---|
| Replicate 1 | 4 | 1 | 35 | 86 | 45 |
| Replicate 2 | 6 | 2 | 29 | 75 | 46 |
| **Overlaps** | 0.25 | 0.5 | 0.83 | 0.66 | 0.94 |

# ESSNET & PFSNET show excellent stability

# ESSNET can assay low-abundance complexes that qPSP cannot
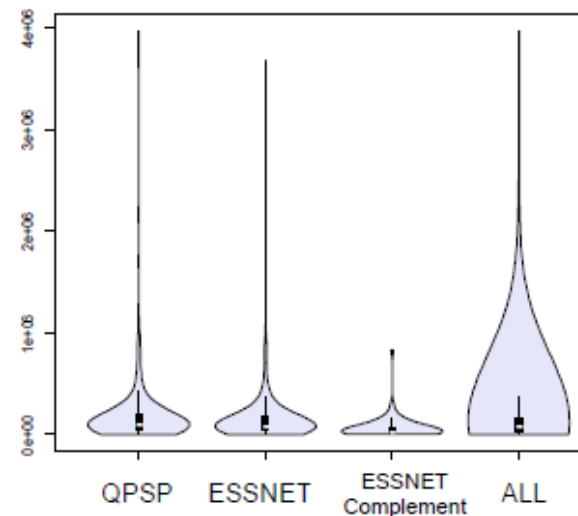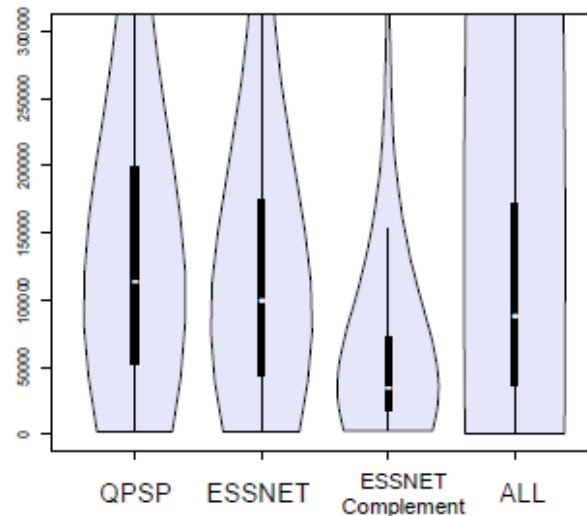


A: QPSP-ESSNET significant-complex overlaps

B: P-value distribution for overlapping and non-overlapping QPSP complexes.

C: Sampling abundance distribution. The left panel is a zoom-in of the right. The y-axis is the protein abundance while the four categories are the distribution of abundances of complexes found in QPSP, ESSNET, ESSNET unique (complement), and all proteins in RC.

# ESSNET can assay low-abundance complexes that PFSNET cannot



Of the 5 ESSNET-unique complexes, PFSNET can detect 4; the missed complex consists entirely of low-abundance proteins.

If p-value threshold is adjusted by Benjamini-Hochberg 5% FDR, PFSNET can detect only 3 of the 5 ESSNET-unique complexes while ESSNET continues to detect them all.

# CONCLUDING REMARKS

# In conclusion…

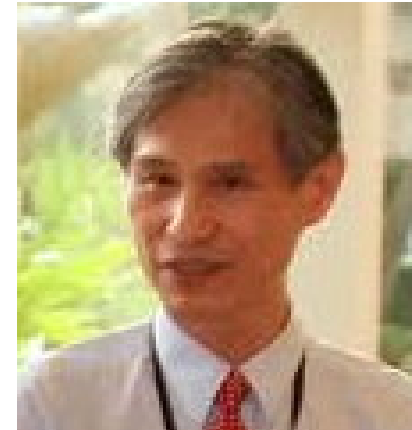## Contextualization (into complexes) can deal with coverage and consistency issues in proteomics

# Acknowledgements

- **The slides on peptide identification were adapted from those given to me by A/P Leong Hon Wai**

- **A lot of the slides on PSP, qPSP, and ESSNet came from the work of Wilson Goh**

**Leong Hon Wai**

**Wilson Goh**

# Must read

- Steen & Mann. **The ABC's and XYZ's of peptide sequencing.** *Nature Reviews Molecular Cell Biology*, 5:699-711, 2004

- Cottrell. **Protein identification using MS/MS data**. *Journal of Proteomics*, 74:1842-1851, 2011

- [FCS] Goh et al. **Comparative network-based recovery analysis and proteomic profiling of neurological changes in valporic acid-treated mice.** *Journal of Proteome Research*, 12(5):2116-2127, 2013

- [qPSP] Goh et al. **Quantitative proteomics signature profiling based on network contextualization**. *Biology Direct*, 10:71, 2015

- [ESSNET] Goh & Wong. **Advancing clinical proteomics via analysis based on biological complexes: A tale of five paradigms**. *Journal of Proteome Research*, 15(9):3167-3179, 2016

# Good to Read

- Frank, et al. **De Novo Peptide Sequencing and Identification with Precision Mass Spectrometry.** *J. Proteome Res*. 6:114-123, 2007

- Sung. **Chap. 12: Peptide sequencing**. *Algorithms in Bioinformatics: A Practical Introduction*. CRC Press, 2010

- Käll & Vitek. **Computational mass spectrometry–based proteomics**. *PLoS Comput Biol ,* 7(12): e1002277, 2011