

NANYANG
TECHNOLOGICAL
UNIVERSITY

Principles of proteomics algorithms

Wilson Wen Bin Goh

School of Biological Sciences, Nanyang Technological University

16 March 2017

Learning objectives

- Describe what an algorithm is, and how it differs from a heuristic or a computer program
- Describe the various levels of spectra data and their derivations in MS-based proteomics
- Describe the steps of a library search algorithm
- Describe the steps of a *de novo* sequencing algorithm
- Describe how peptides are assembled to proteins and associated problems
- Describe and evaluate the various levels and data representation formats in proteomics

Algorithms

- <https://www.youtube.com/watch?v=6hfOvs8pY1k>

What is an algorithm?

- An **algorithm** is an ordered series of steps for solving the problem
- It is very **exact** and **unambiguous**
- It can be **expressed** by a programming language (where it becomes a program)
- Can also be expressed in semi-human readable form (in pseudocode)

What is an algorithm? (Cont'd)

- Example: an “algorithm” to come to this lecture

1. Walk to bus stop
2. Catch bus
3. Get off near Uni
4. Walk to lecture theatre

- Is this enough as directions to give someone?

- What about: (supposing you live near Prime)

1. Walk to bus stop in Nanyang Drive, in front of Prime
2. Catch bus number 1 or 2
3. Get off at bus stop code 123456
4. Walk up the first staircase you see and turn to the right on

the second floor. Enter LT 25

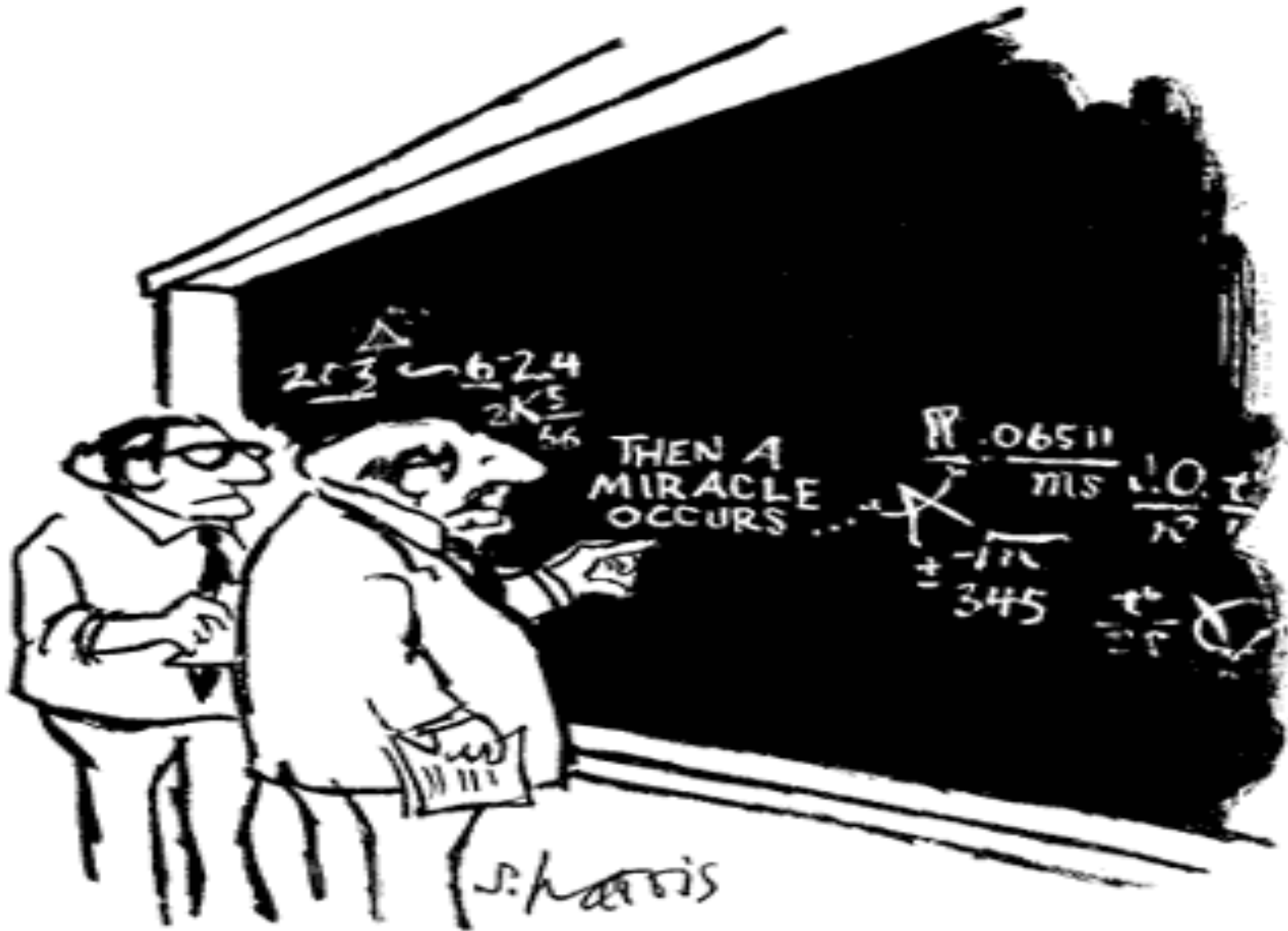
See an algorithm in action --- The bubble sort algorithm

- Objective is to sort objects using an iterative window.
- <https://www.youtube.com/watch?v=IUqFelc84XE>

How detailed should an algorithm be?

- First answer: “there should be enough detail”
 - How “enough” is “enough”?
 - Depends on who will read the algorithm (audience)
 - Depends on what the algorithm is for (problem to solve)
 - Depends on why someone will read algorithm (purpose)
- Second answer: “it should have enough detail so as to allow someone to
 - Understand each and every step
 - Follow the algorithm (manually) to work out a solution
 - Implement it, adapt it, extend it, embed it,...

How detailed should an algorithm be?



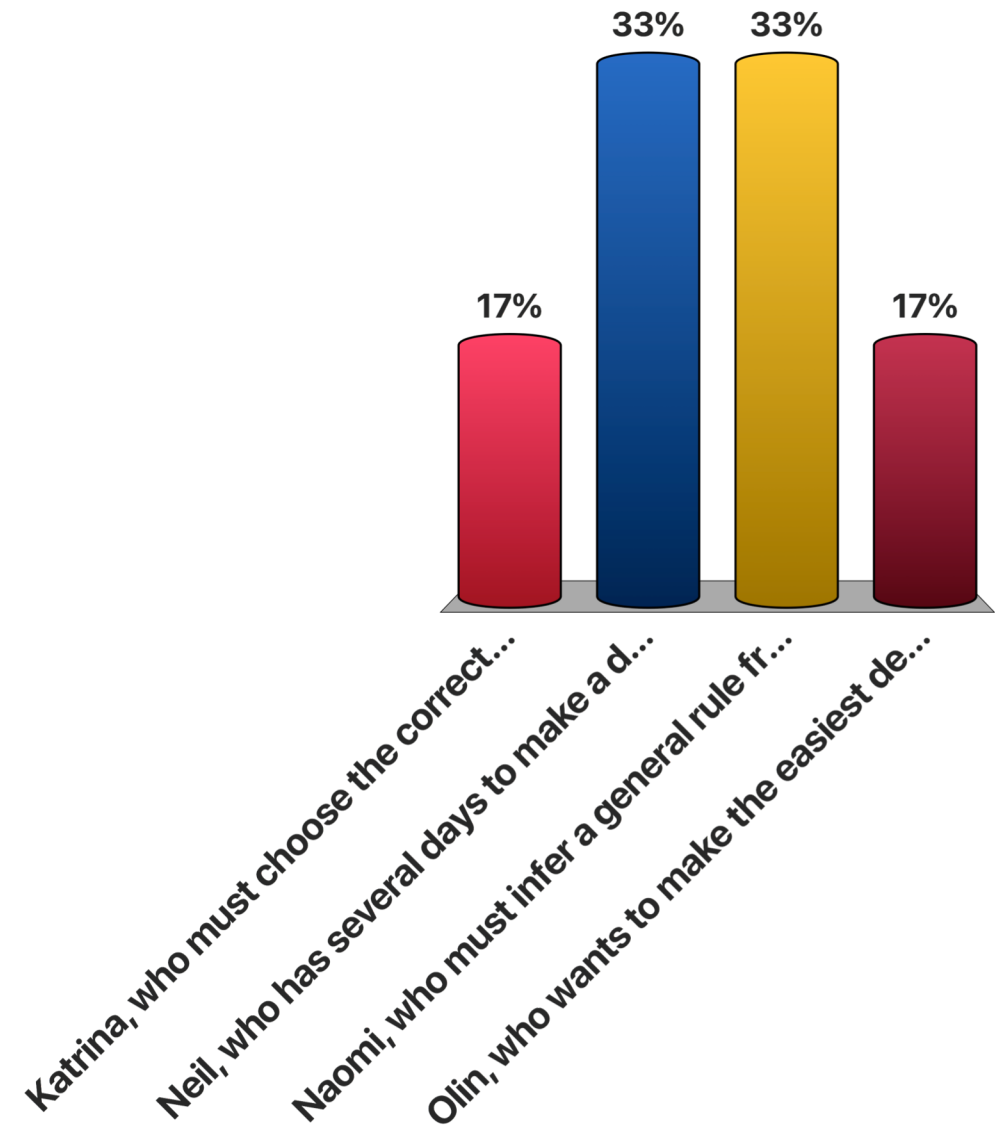
"I think you should be more explicit here in step two."

Is any process an algorithm?

- What is the difference between an *algorithm*, *heuristic* and a program?
 - Algorithm – a method to problem solving which is guaranteed to give the correct output for a given input, e.g. adding two numbers. Might not be optimally efficient (space/time).
 - Heuristic – a practical method to problem solving which is not guaranteed to give the correct output for a given input, but it is likely to be correct. Might not be optimally efficient (space/time). Predicting winner in a football match, e.g. score is 5-1, 5 minutes before end of game.
 - An algorithm is not the same as the computer program. The computer program is an implementation of the algorithm or a heuristic in a particular programming language.

Heuristics will be most useful for which of the following people?

- A. Katrina, who must choose the correct answer on a multiple-choice exam.
- B. Neil, who has several days to make a decision about which company to join.
- C. Naomi, who must infer a general rule from a set of propositions.
- D. Olin, who wants to make the easiest decision about which job to take.



A recap...Typical proteomic MS expt

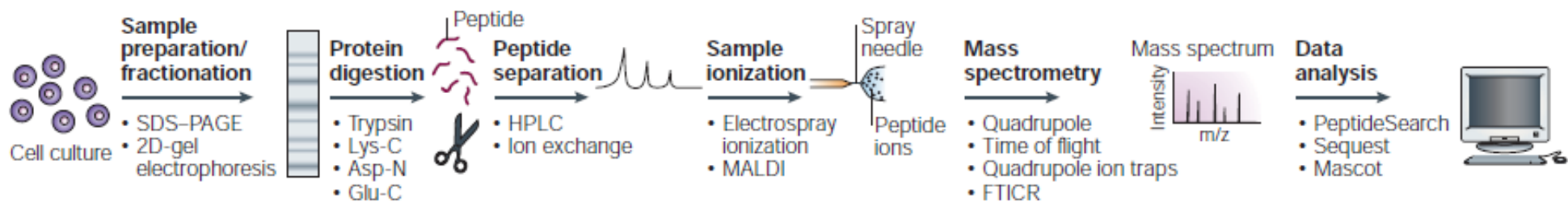
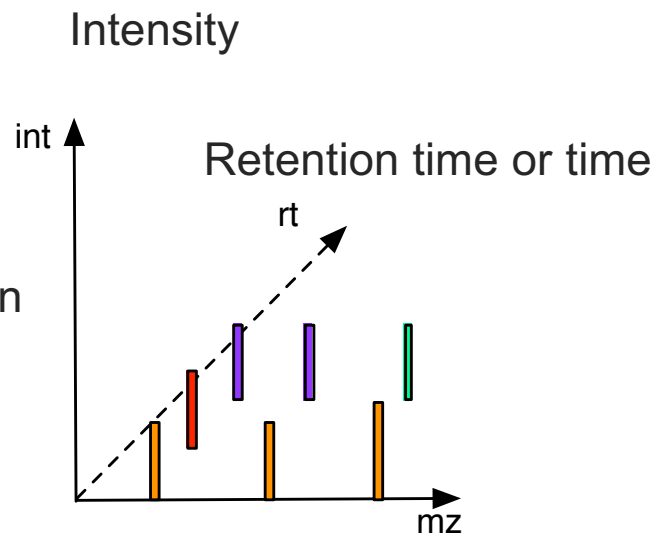


Figure 1 | **The mass-spectrometry/proteomic experiment.** A protein population is prepared from a biological source — for example, a cell culture — and the last step in protein purification is often SDS-PAGE. The gel lane that is obtained is cut into several slices, which are then in-gel digested. Numerous 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.

What we are really measuring...

3 dimensions of information

Unlike genomics, MS does not give direct information on sequence information



The mass-to-charge ratio (also referred to as m/z)

An example of real spectra

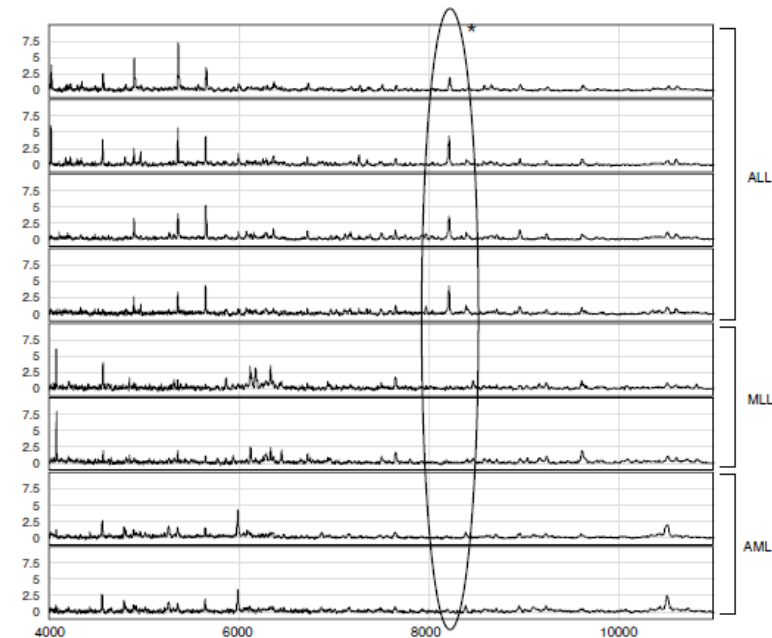


Figure 1 Spectra from SELDI-TOF MS analysis of REH, 697, MV4;11, and Kasumi cell lines. Protein (4 μ g) from each cell type was analyzed on SAX2 ProteinChip[®] Arrays. ALL cell lines shown are REH and 697, the MLL cell line is MV4;11, and the AML cell line is Kasumi. The asterisk indicates the differentially expressed protein at 8.3 kDa.

Does this look like biological information?

Note: Although we can measure 3 dimensions, normally we only use m/z and int for identification. rt is usually for effecting separation

Comparing proteomics and transcriptomics

Similarities

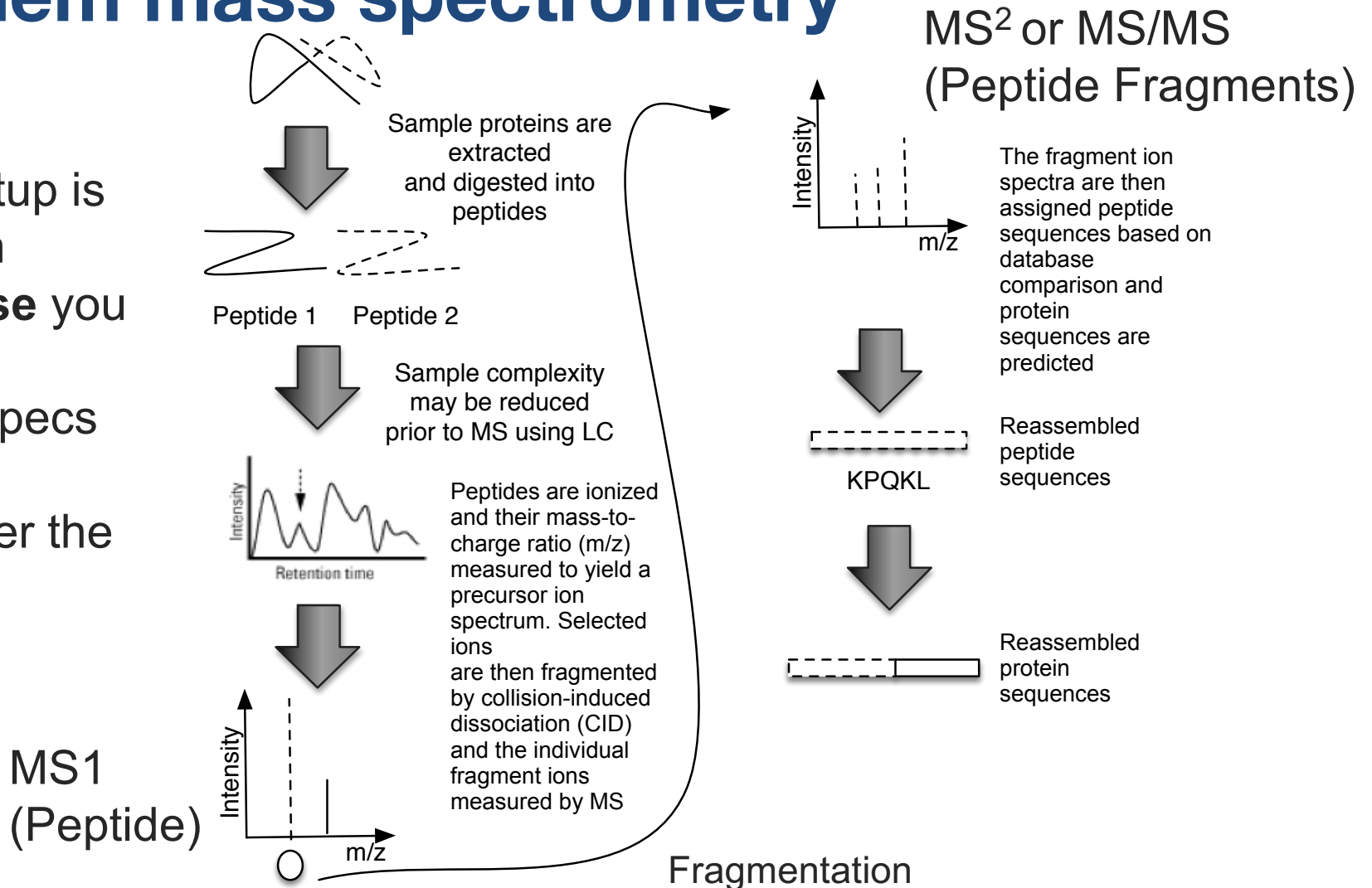
- Proteomic profile
 - Which protein is found in the sample
 - How abundant it is
- Similar to gene expression profile. So typical gene expression profile analysis methods can be applied in theory...

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

A look at the important steps of tandem mass spectrometry

The setup is tandem **because** you have 2 mass specs set up one after the other



Important concepts:

- 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 signal from the mass/charge ratio of an ion

Relating the various levels of data in MS-based proteomics

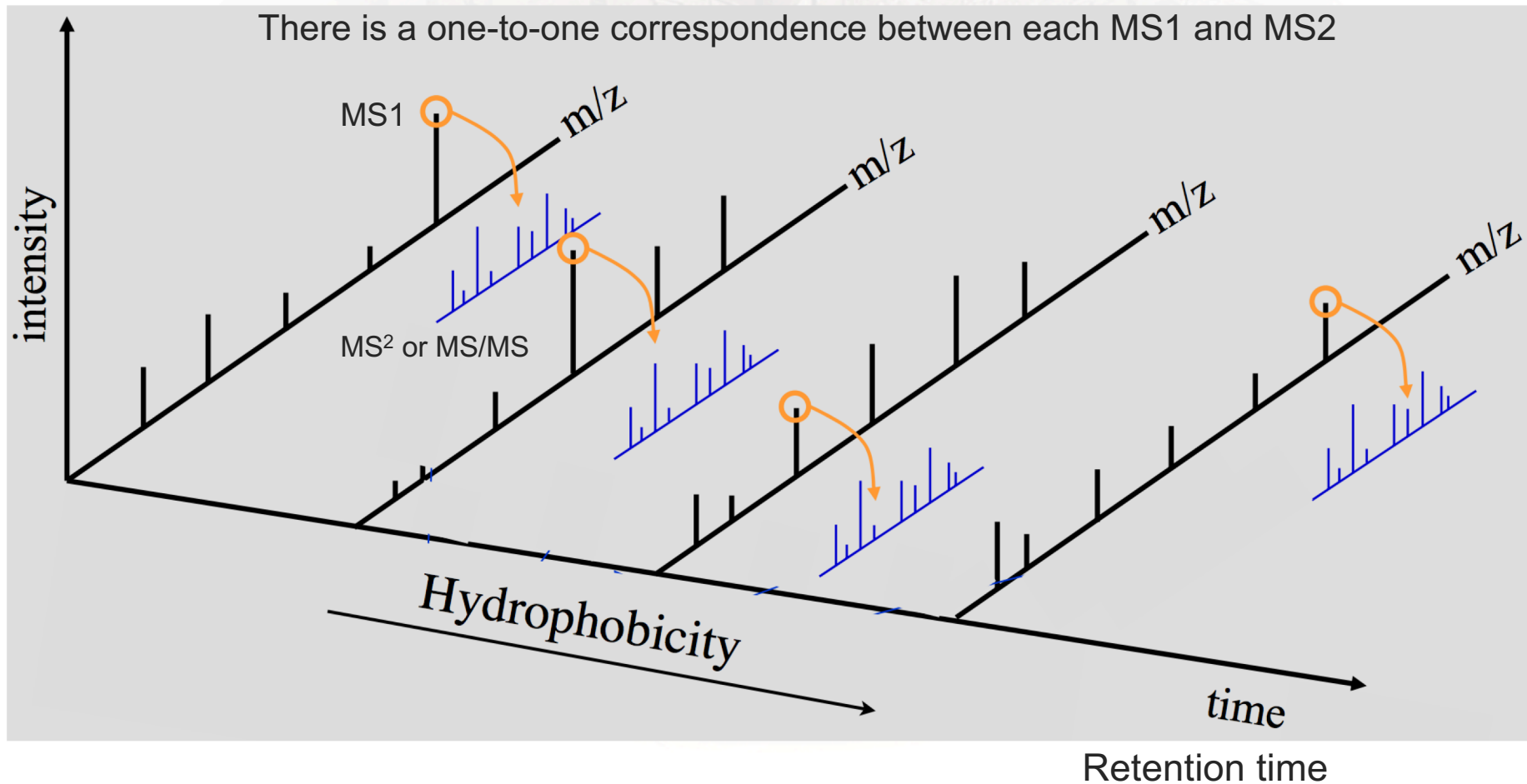


Image from Vitek O, Purdue University

The problem

- We only have a series of peaks to tell us whether a protein exists or not.
- The peaks relate to masses obtained from the fragmentation of a peptide
- Random fragmentation of peptide along any bond will result in a large number of masses

Some help from biology and chemistry

- Trypsin has defined cleavage sites (so we know the ends of each peptide)
- Collision-induced dissociation tends to break peptides along the peptide bond, resulting in partial amino acids
- We know the masses of each of the 20 amino acids
- What further information do we need?

An idea

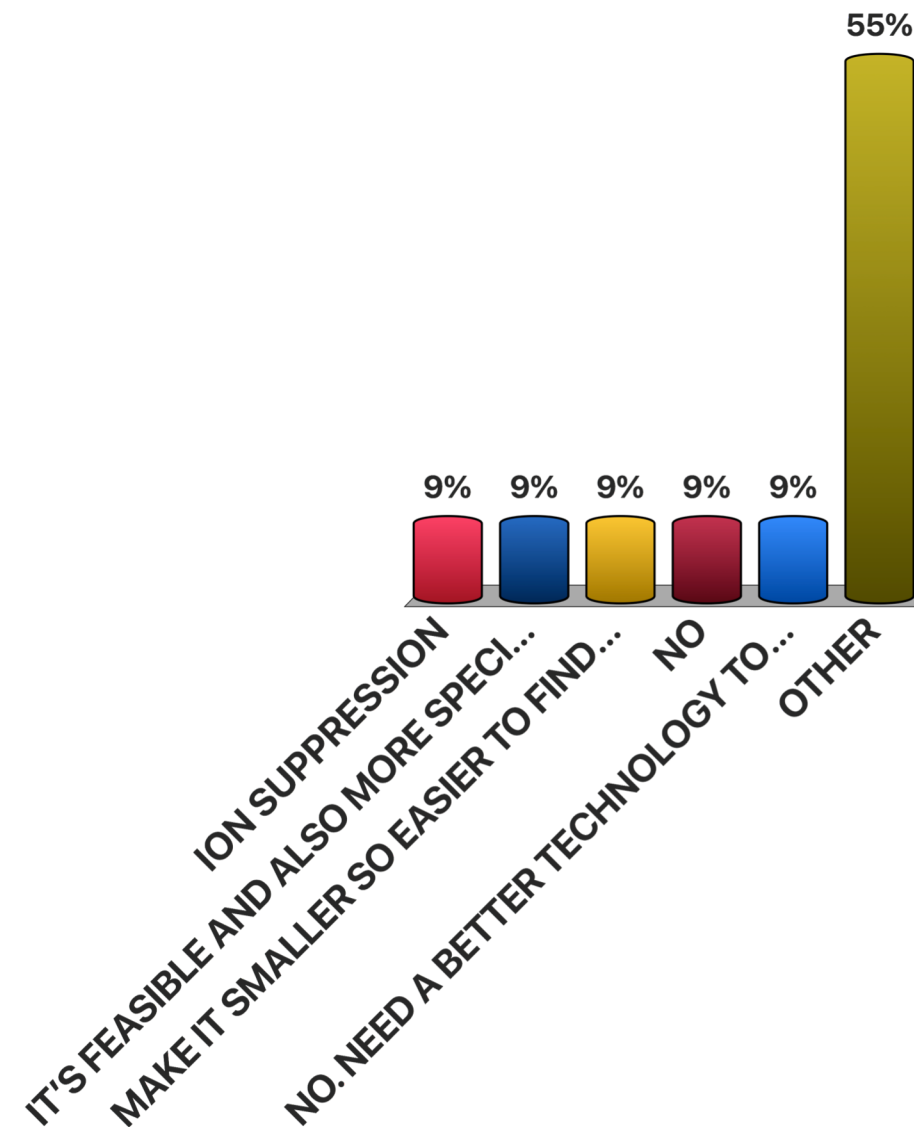
- The ends of trypsin-cleavage are “deterministic”. Bond breaking is also “deterministic”^{*}.
- We know the sequences of proteins using *in silico* translation of mRNA sequences.
- What if.. We simulate spectra from known sequences and compare it against our observed spectra?

^{*}When I say deterministic, I mean it is conserved behavior when it does happen. It does not mean that it always happens.

Is it feasible to use whole proteins for sequence matching? Would it make the matching more specific? Can it cause problems?

Rank	Responses
1	ION SUPPRESS...
2	IT'S FEASIBL...
3	MAKE IT SMAL...
4	NO
5	NO. NEED A B...
6	OTHER

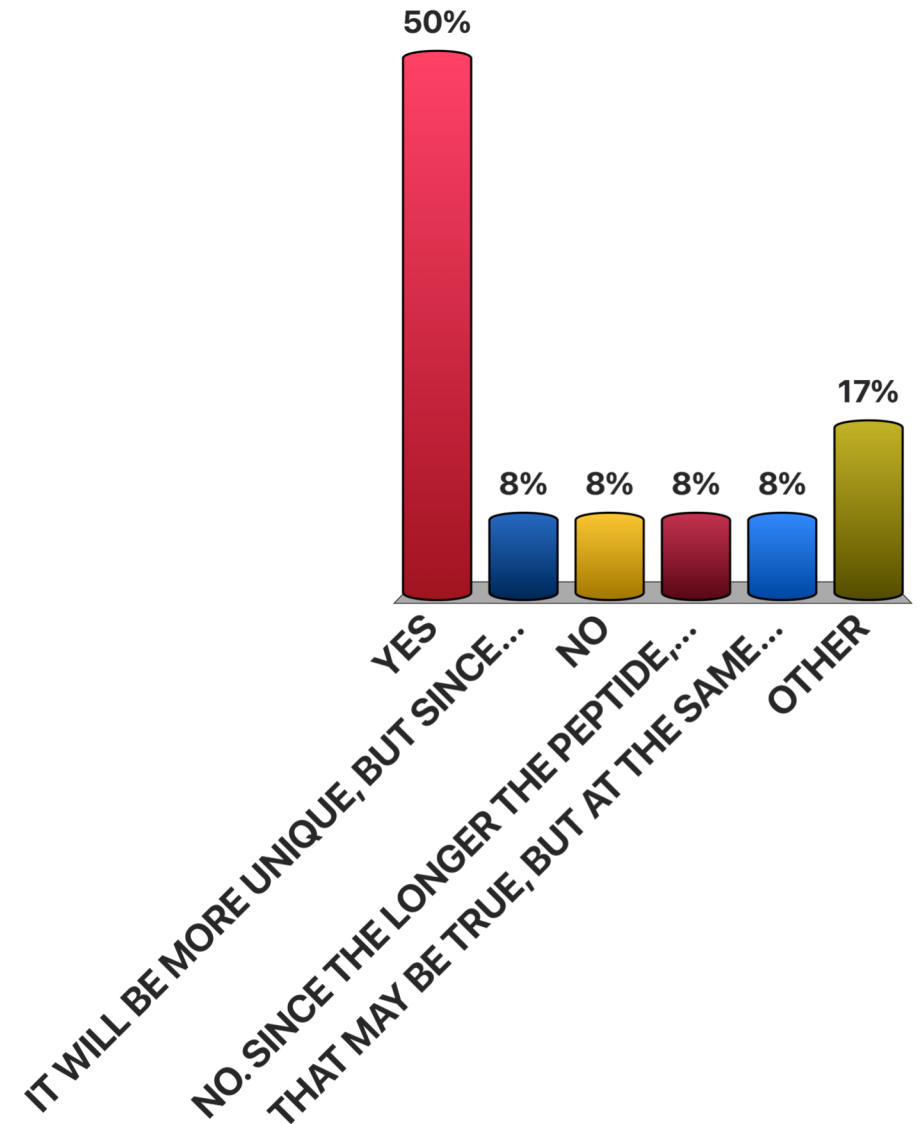
Keyword:
Keyword Matches: 0



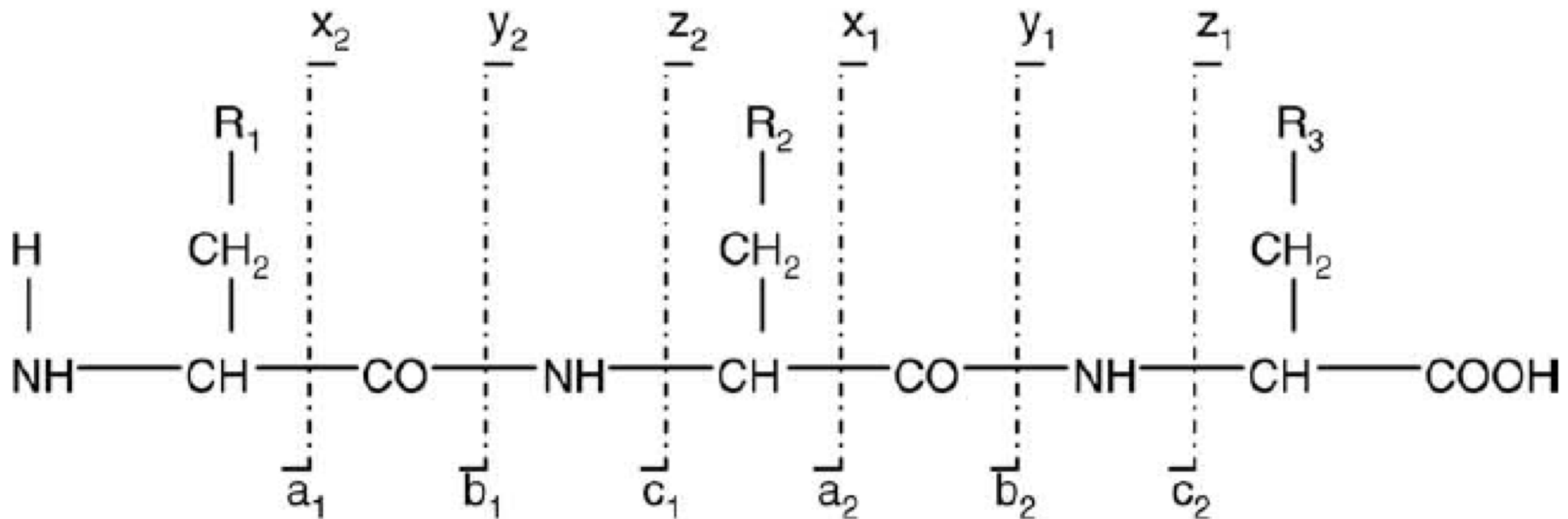
Do you agree that the longer the peptide sequence, the more unique its pattern, and therefore the more specific the match?

Rank	Responses
1	YES
2	IT WILL BE M...
3	NO
4	NO. SINCE TH...
5	THAT MAY BE ...
6	OTHER

Keyword:
Keyword Matches: 0



Revisiting peptide fragmentation



Generation of beta and gamma ions are more common so they may be more informative during analysis (more completeness)

Fragments can also lose neutral chemical groups like NH_3 and H_2O

Peptide fragmentation (mass calculations)

Ion type	K	V	P	Q	V	S	T	P	T	L	R
b+		228.2	325.2	453.3	552.4	639.4	740.4	837.5	938.5	1051.6	
b++		114.6	163.1	227.1	276.7	320.2	370.7	419.3	469.8	526.3	
y+		1097.6	998.6	901.5	773.5	674.4	587.4	486.3	389.3	288.2	175.1
y++		549.3	499.8	451.3	387.2	337.7	294.2	243.7	195.1	144.6	88.1

$$M_V = 99.1$$

$$M_K = 128.1$$

$$M_P = 97$$

Peptide fragmentation (mass shifts)

Ion type	K	V	P	Q	V	S	T	P	T	L	R
b →		228.2	325.2	453.3	552.4	639.4	740.4	837.5	938.5	1051.6	
b++		114.6	163.1	227.1	276.7	320.2	370.7	419.3	469.8	526.3	
y		1097.6	998.6	901.5	773.5	674.4	587.4	486.3	389.3	288.2	175.1
y++		549.3	499.8	451.3	387.2	337.7	294.2	243.7	195.1	144.6	88.1
							+80	+80	+160	+160	

Ion type	K	V	P	Q	V	S	T*	P	T*	L	R
b →		228.2	325.2	453.3	552.4	639.4	820.4	917.5	1098.5	1211.6	
y		1257.6	1158.5	1061.4	933.4	834.3	747.3	566.3	469.2	288.2	175.1

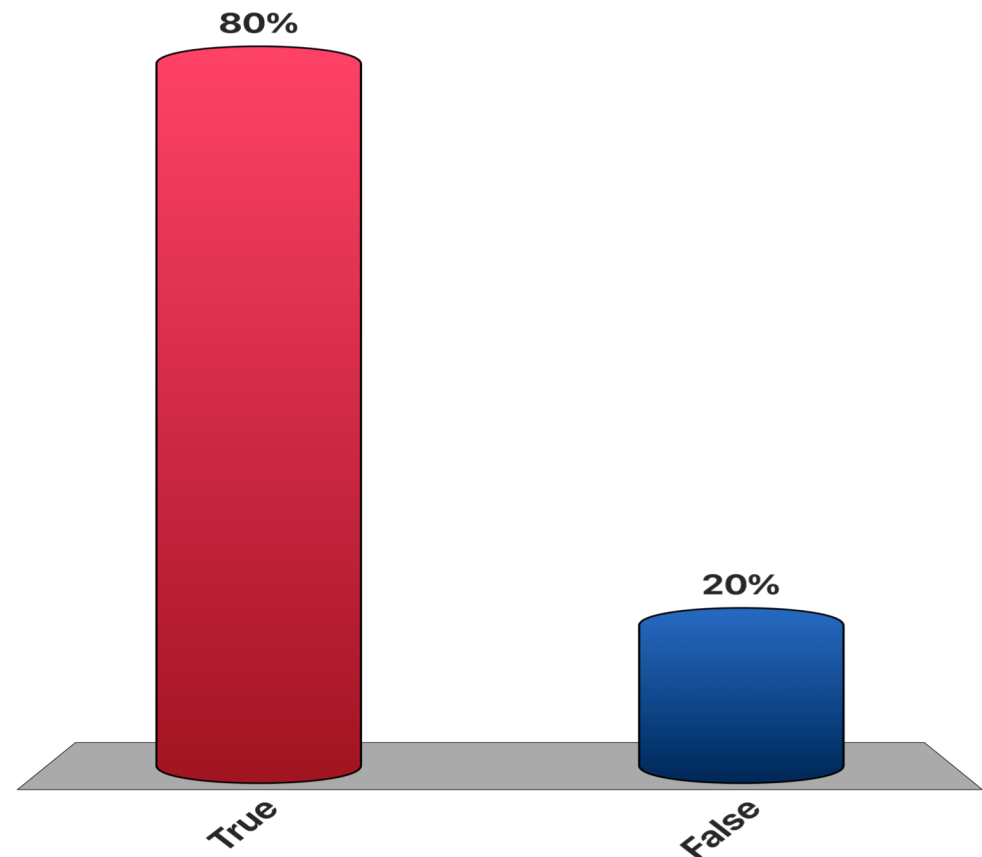
Phosphorylated threonine residues (.79.9663 Da). Note that all fragment ions including the ion with one or two threonine residues are shifted in mass once or twice, respectively.

Food for thought --- where are these PTM sites supposed to be anyway?

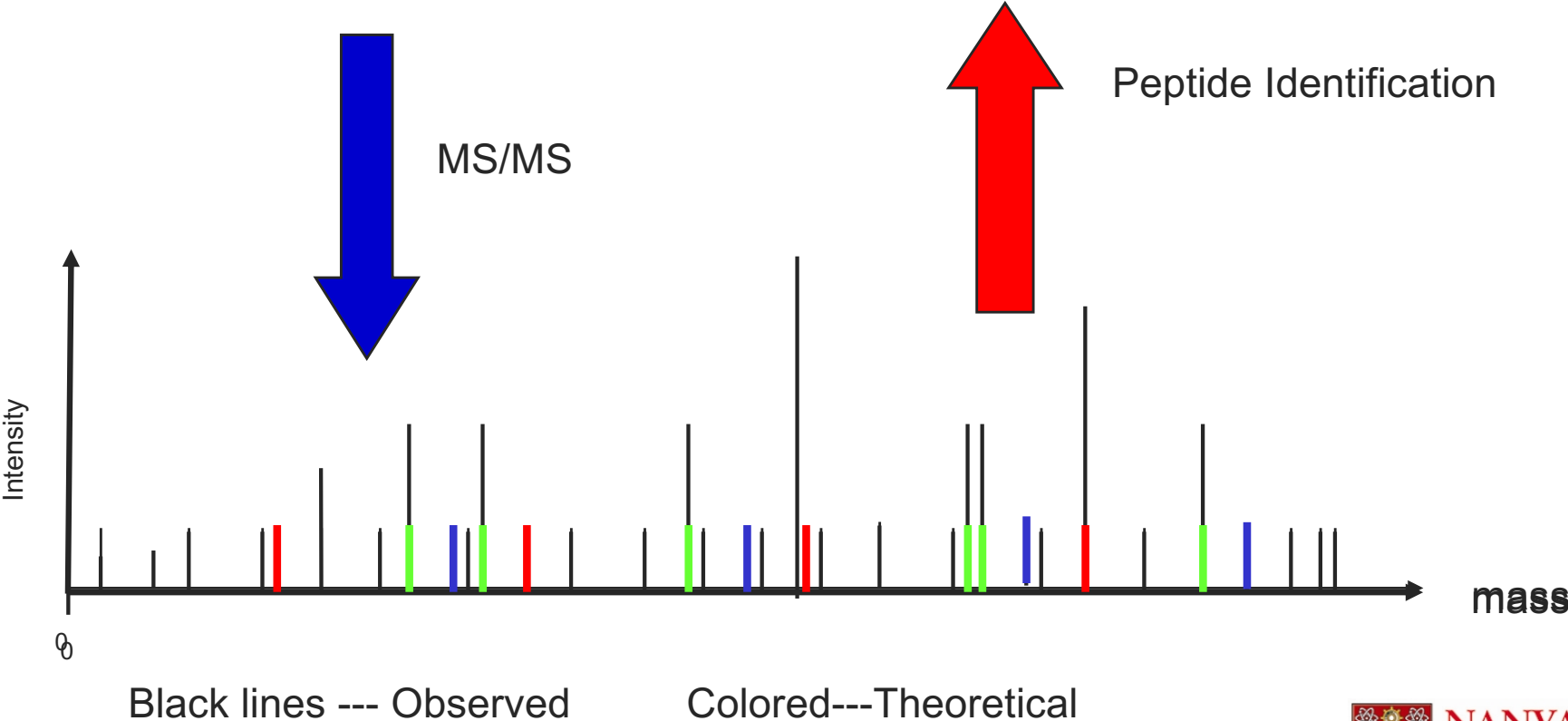
It is important to use bioinformatics tools to predict sites for posttranslational modifications based on specific protein sequences. However, prediction of such modifications can often be difficult because the short lengths of the sequence motifs associated with certain modifications.

A. True

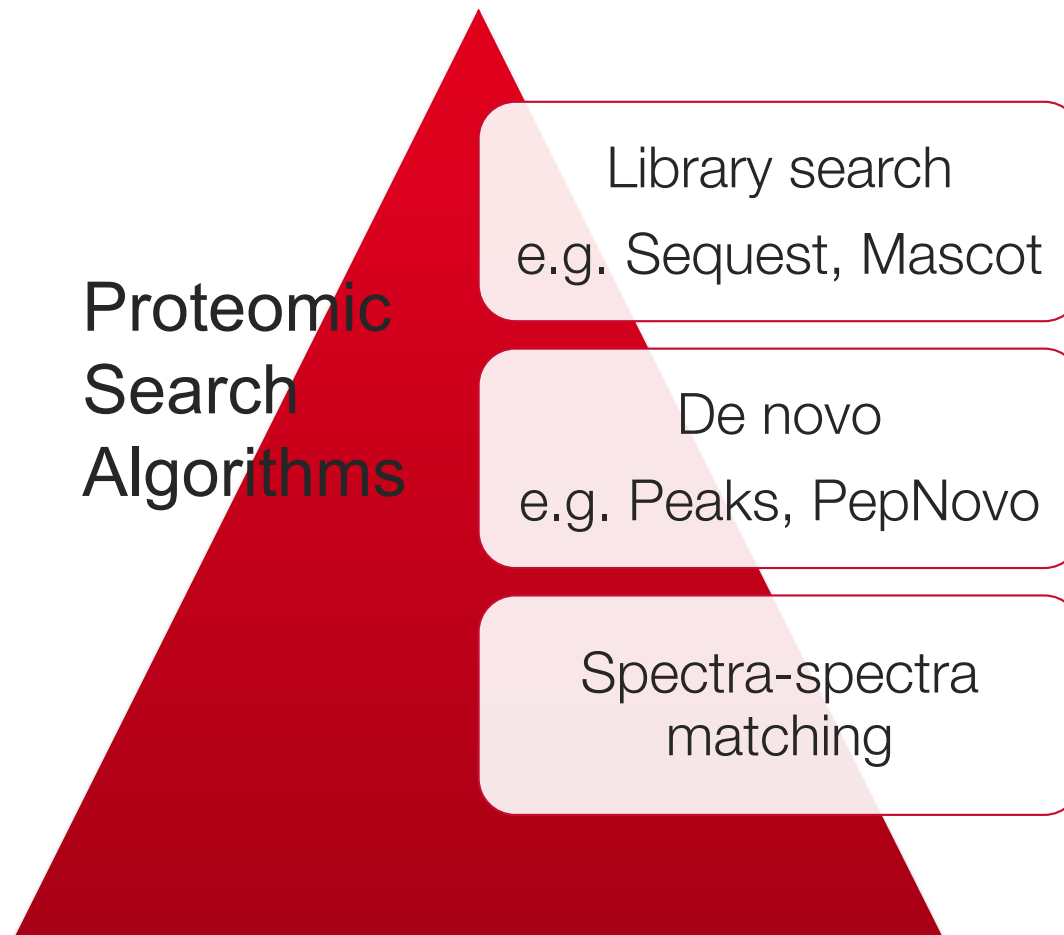
B. False



Protein identification with MS/MS



3 main approaches for peptide identification



Library search algorithms

- Library (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

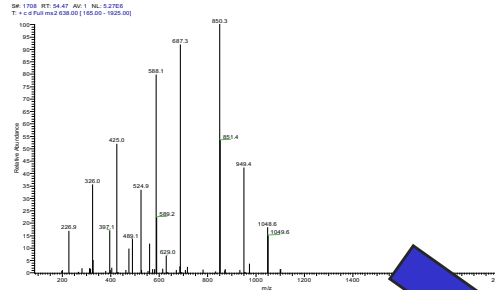
- Given this peptide



- Its theoretical spectrum is
- Theoretical spectrum is dependent on
 - Set of ion-types considered (beta-gamma)
 - Larger if multi-charge ions are considered
 - (+2, +3, +4)

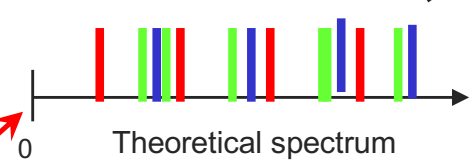
The steps of a library search algorithm

Database Search



Database of known peptides

MDERHILNM, KLQWVCS DL, PTYWASDL,
ENQIKRSACVM, TLACHGGEM, NGALPQWRT,
HLLERTKMNVV, GGPASSDA, GGLITGMQSD,
MQPLMNWE,
ALKIIMNVRT, **AVGELTK**, HEWALIF,
GHNLWAMNAC, GVFGSVLRA, EKLNKATYIN..



Match

Matching Score for this peptide

Repeat for all the peptides in the Database

What is the database?

Advances in high-throughput and advanced technologies allow researchers to routinely perform whole genome and proteome analysis. For this purpose, they need high-quality resources providing comprehensive gene and protein sets for their organisms of interest. Using the example of the human proteome, we will describe the content of a complete proteome in the UniProt Knowledgebase (UniProtKB). We will show how manual expert curation of UniProtKB/Swiss-Prot is complemented by expert-driven automatic annotation to build a comprehensive, high-quality and traceable resource. We will also illustrate how the complexity of the human proteome is captured and structured in UniProtKB.

Database URL : <http://www.uniprot.org/>

The expert curation of all functional protein isoforms produced by **alternative splicing** will require more resources and time. However, a complementary pipeline for the import of predicted human protein sequences in UniProtKB/TrEMBL has been developed in collaboration with Ensembl to complete the set of human isoform sequences. >49 000 additional predicted alternative products are currently available in UniProtKB/TrEMBL.

Read: <https://academic.oup.com/database/article/doi/10.1093/database/bav120/2630095>

The steps of a library search algorithm

Pseudocode example for peptide mass finger printing

```
Get the experimental mass list L
For each sequence s in the database do
  digest s and obtain a set of peptides P
  for each peptide p in P do
    compute mass(p)
    push mass(p) in M
  x <- score(M, L)
  store score x for protein s
compute p-values for each score
return the n best proteins /* highest score or lowest p-value */
```

Is this suitable for helping us do library search as we have looked at earlier?

Is this suitable for helping us do library search as we have looked at earlier? What is missing?

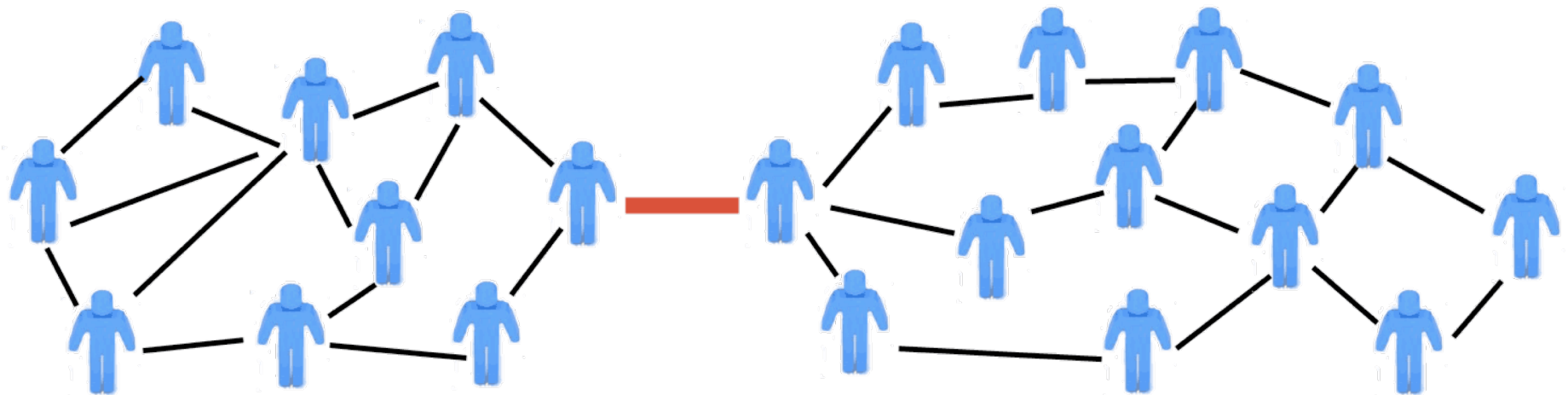
Rank	Responses
1	
2	
3	
4	
5	
6	Other

The steps of a *de novo* sequencing algorithms

- Given a spectrum
 - Build a spectrum graph (A spectrum graph is simply a graph that is built from your Mass spectra.)
 - Peptides are paths in this graph (A path is simply a directed walk from point A to point B or C)
 - Find the best path

A graph

- A generalization of the simple concept of a set of items and the relationships between them
- Representation: Graph $G = (V, E)$ consists set of vertices denoted by V , or by $V(G)$ and set of edges E , or $E(G)$



Connectivity

- Basic Idea: In a Graph Reachability among vertices by traversing the edges

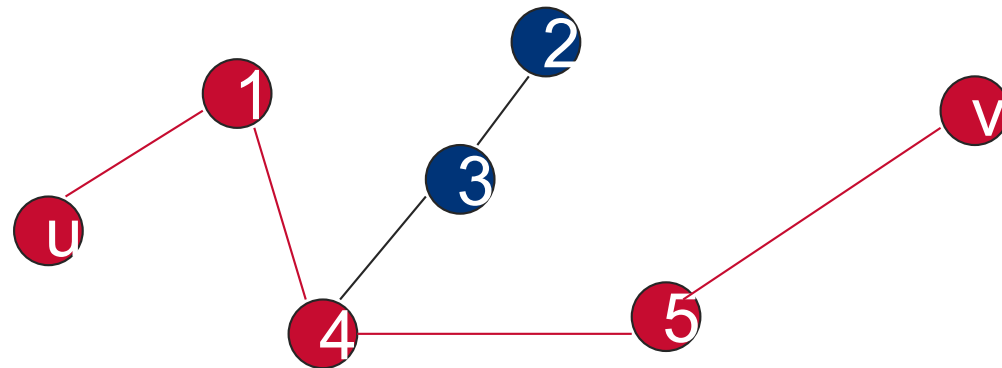
Application Example:

- In a city to city road-network, if one city can be reached from another city.
- Problems if determining whether a message can be sent between two computer using intermediate links
- Efficiently planning routes for data delivery in the Internet

Connectivity and paths

A **Path** is a sequence of edges that begins at a vertex of a graph and travels along edges of the graph, always connecting pairs of adjacent vertices.

Representation example: $G = (V, E)$, Path P represented, from u to v is $\{u, 1\}, \{1, 4\}, \{4, 5\}, \{5, v\}$



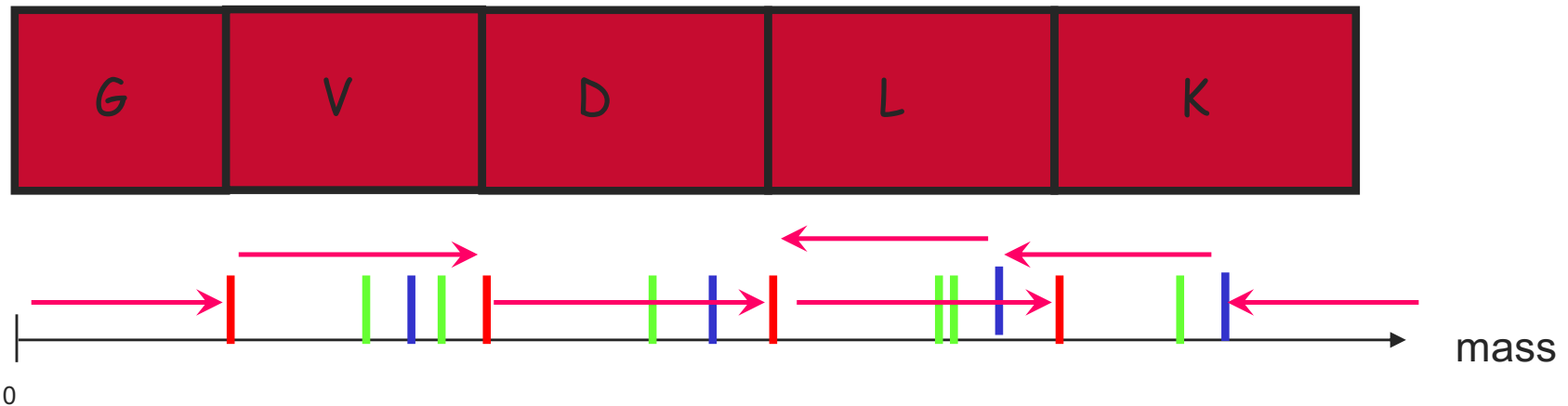
Paths do not cross the same vertex twice.
There are no loops in a path

Amino acids have specific masses

1-letter code	3-letter code	Chemical formula	Monoisotopic	Average
A	Ala	C ₃ H ₅ ON	71.03711	71.0788
R	Arg	C ₆ H ₁₂ ON ₄	156.10111	156.1875
N	Asn	C ₄ H ₆ O ₂ N ₂	114.04293	114.1038
D	Asp	C ₄ H ₅ O ₃ N	115.02694	115.0886
C	Cys	C ₃ H ₅ ONS	103.00919	103.1388
E	Glu	C ₅ H ₇ O ₃ N	129.04259	129.1155
Q	Gln	C ₅ H ₈ O ₂ N ₂	128.05858	128.1307
G	Gly	C ₂ H ₃ ON	57.02146	57.0519
H	His	C ₆ H ₇ ON ₃	137.05891	137.1411
I	Ile	C ₆ H ₁₁ ON	113.08406	113.1594
L	Leu	C ₆ H ₁₁ ON	113.08406	113.1594
K	Lys	C ₆ H ₁₂ ON ₂	128.09496	128.1741
M	Met	C ₅ H ₉ ONS	131.04049	131.1926
F	Phe	C ₉ H ₉ ON	147.06841	147.1766
P	Pro	C ₅ H ₇ ON	97.05276	97.1167
S	Ser	C ₃ H ₅ O ₂ N	87.03203	87.0782
T	Thr	C ₄ H ₇ O ₂ N	101.04768	101.1051
W	Trp	C ₁₁ H ₁₀ ON ₂	186.07931	186.2132
Y	Tyr	C ₉ H ₉ O ₂ N	163.06333	163.1760
V	Val	C ₅ H ₉ ON	99.06841	99.1326

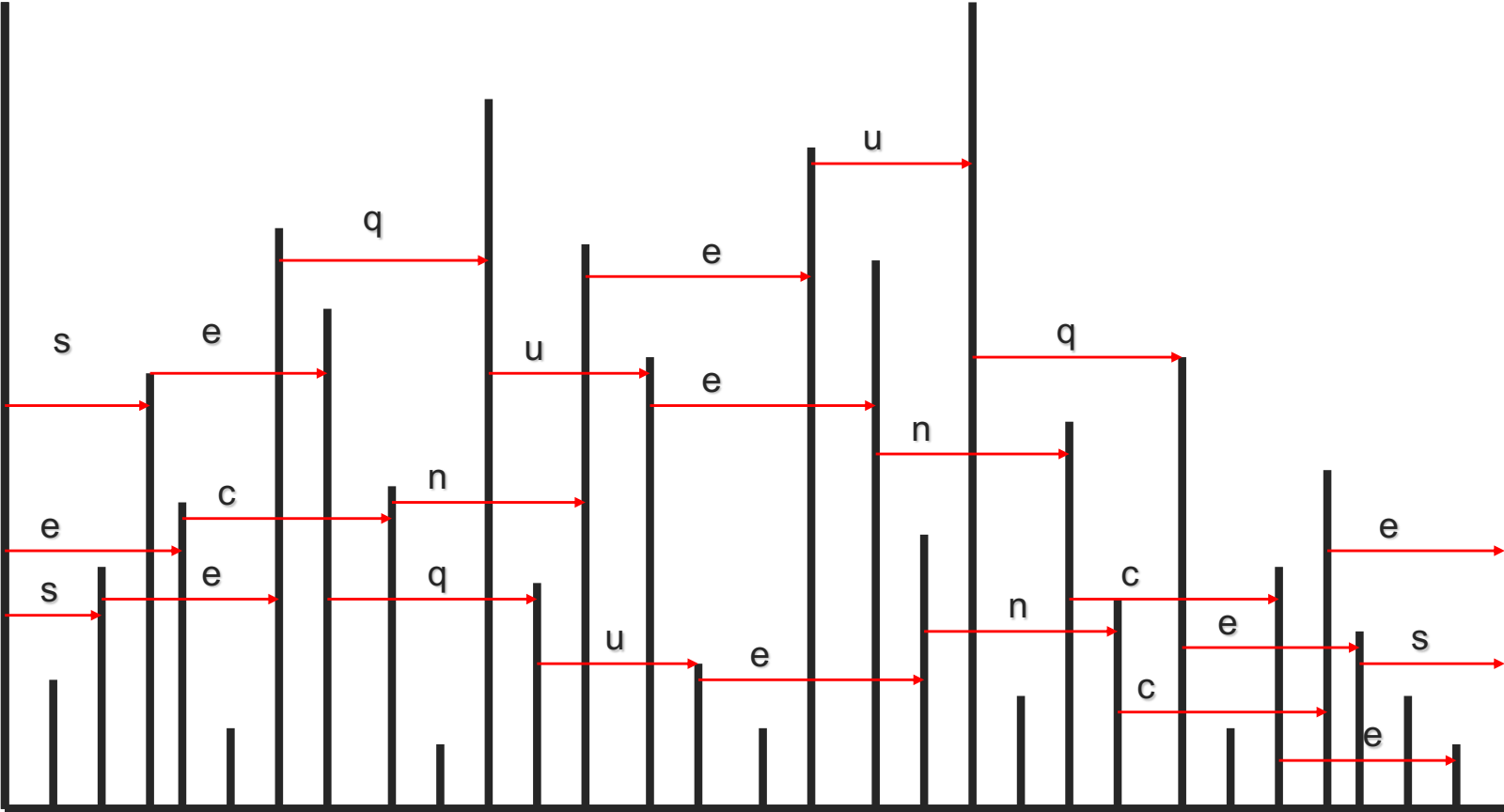
We can use this mass info to reconstruct the peptide sequence! Which amino acid cannot we not distinguish?

Spectrum graph for a peptide

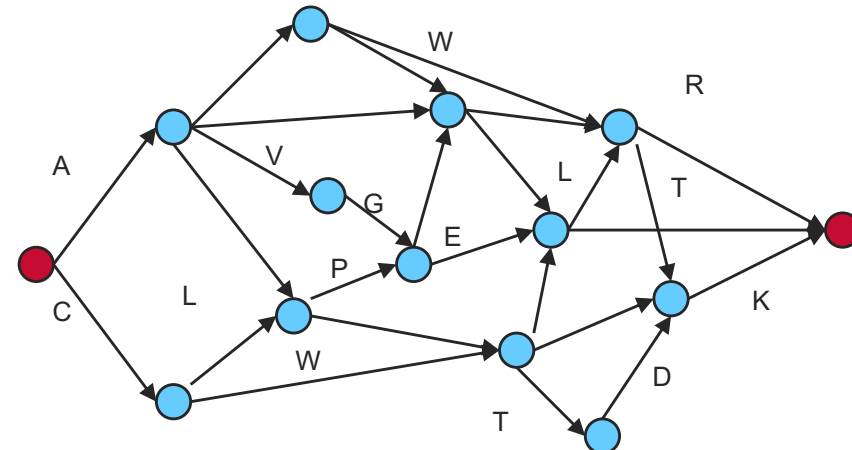
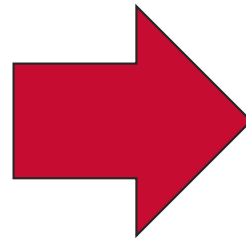
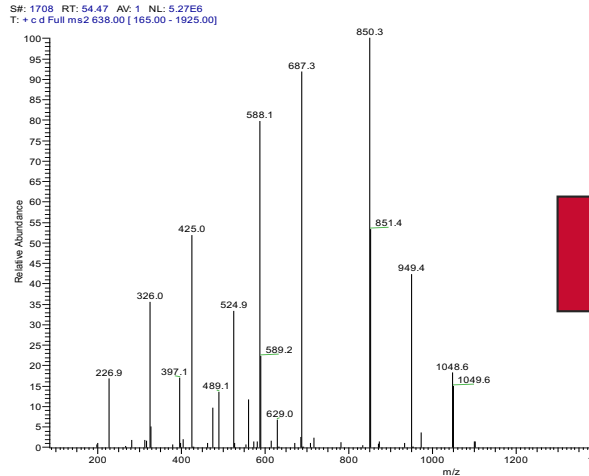


- Connect peaks together
 - If their mass difference = mass of an amino acid
- Theoretical spectrum is dependent on
 - Set of ion-types considered
 - Larger if multi-charge ions are considered

Building a graph from a spectrum



De novo sequencing algorithms

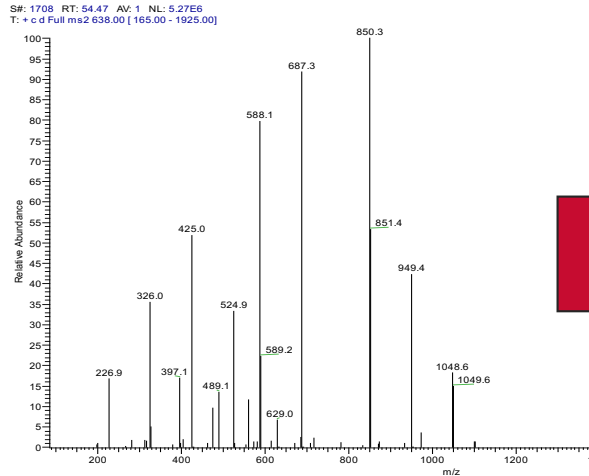


Given a series of peaks...

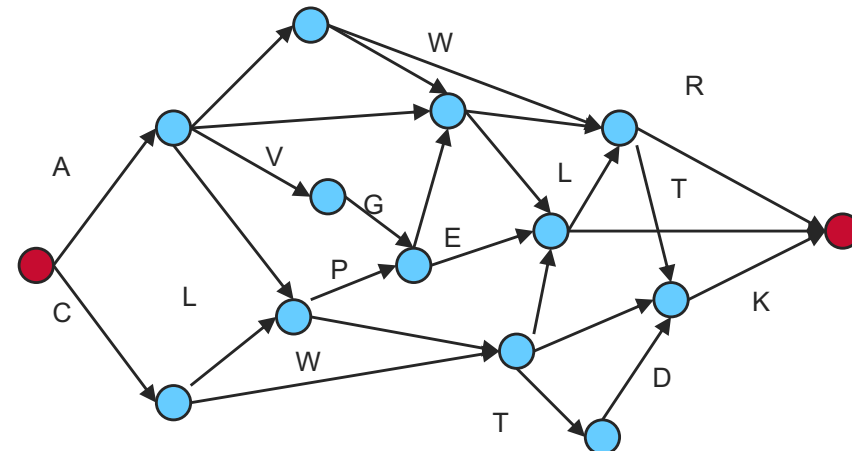
We generate a DAG (Directed Acyclic Graph) from the spectra

Assuming we know that A is the start, and the K is the terminal
Try working out the possible sequences (paths) from start to end.

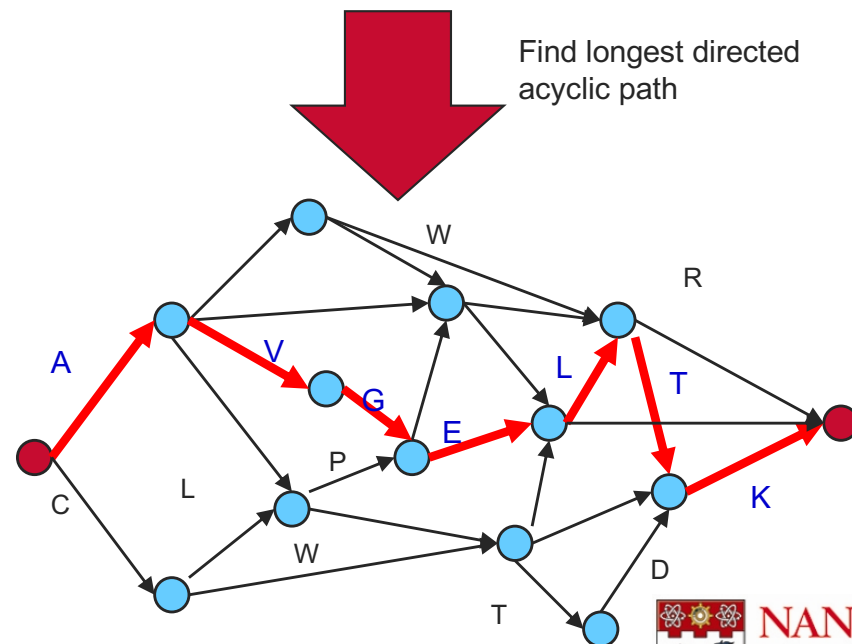
De novo sequencing algorithms



Given a series of peaks...



Find longest directed acyclic path



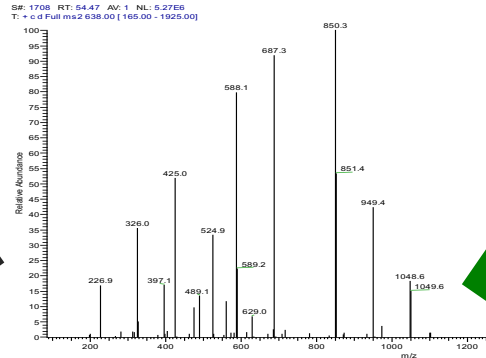
AVGELTK

What is the rationale for the longest possible path being correct? (The best explanation is the one that can explain the presence of all the peaks)

Is the longest path necessary the correct sequence?
What could go wrong?

Rank	Responses
1	
2	
3	
4	
5	
6	Other

2 paths to the same answer ---De novo vs. database search

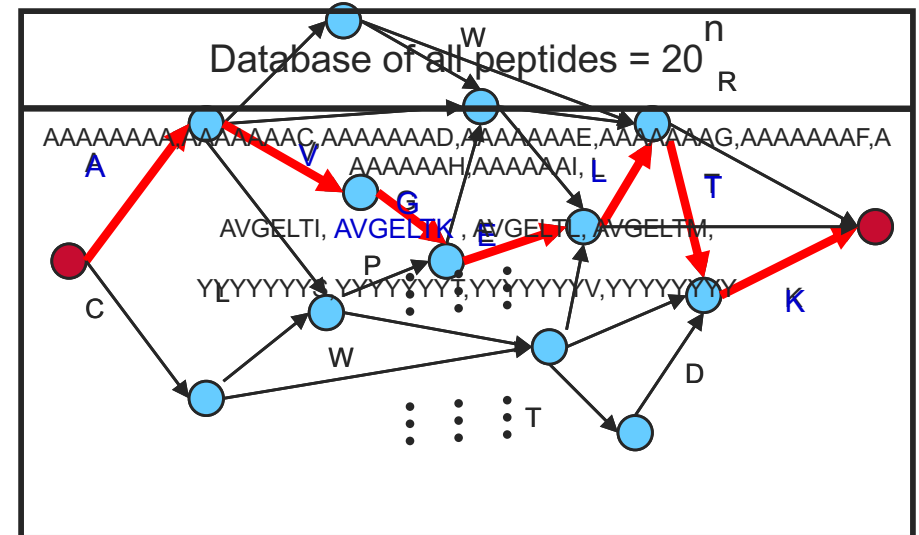


Database Search

De Novo

Database of known peptides

MDERHILNM, KLQWVCS DL, PTYWASDL, ENQIKRSACVM, TLACHGGEM, NGALPQWRT, HLLERTKMN VV, GGPASSDA, GGLITGMQSD, MQPLMNWE, ALKIIMNVRT, **AVGELTK**, HEWAILF, GHNLWAMNAC, GVFGSVLRA, EKLNKAATYIN..



AVGELTK

De novo vs. database search: A paradox

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

Break

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

Notice we haven't said anything about quantitation yet

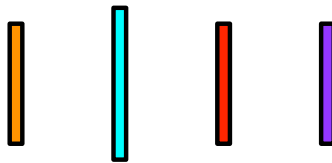
General rules for confident protein identification

- At least 2 unique peptides *
- Coverage of at least 50-70% total protein sequence
- High abundance

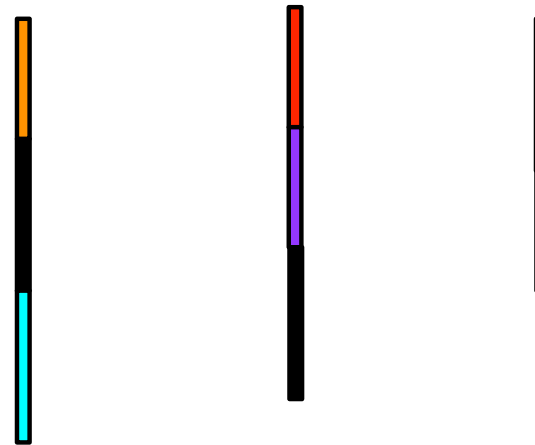
* Commonly used rule.

Peptides to proteins

Identified peptides



known protein sequences



Protein A Protein B Protein C

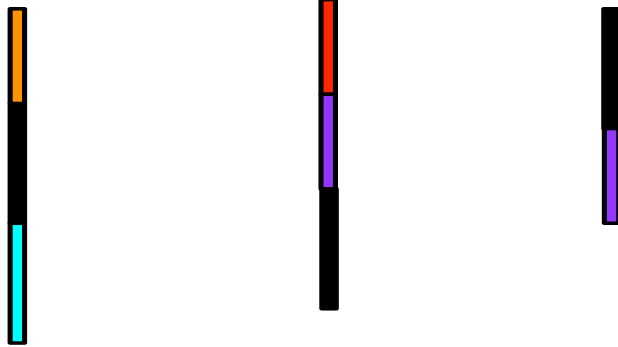


Sequence that is not supported by any spectra

Order the proteins from highest to lowest confidence ?

Order the proteins from highest to lowest confidence ?

known protein sequences



Protein A

Protein B

Protein C

- A. B -> C -> A
- B. A -> B -> C
- C. A -> C -> B
- D. There is no clear answer

Food for thought --- meaningful ambiguity?

Identified peptides

known protein sequences



Protein A

Protein B

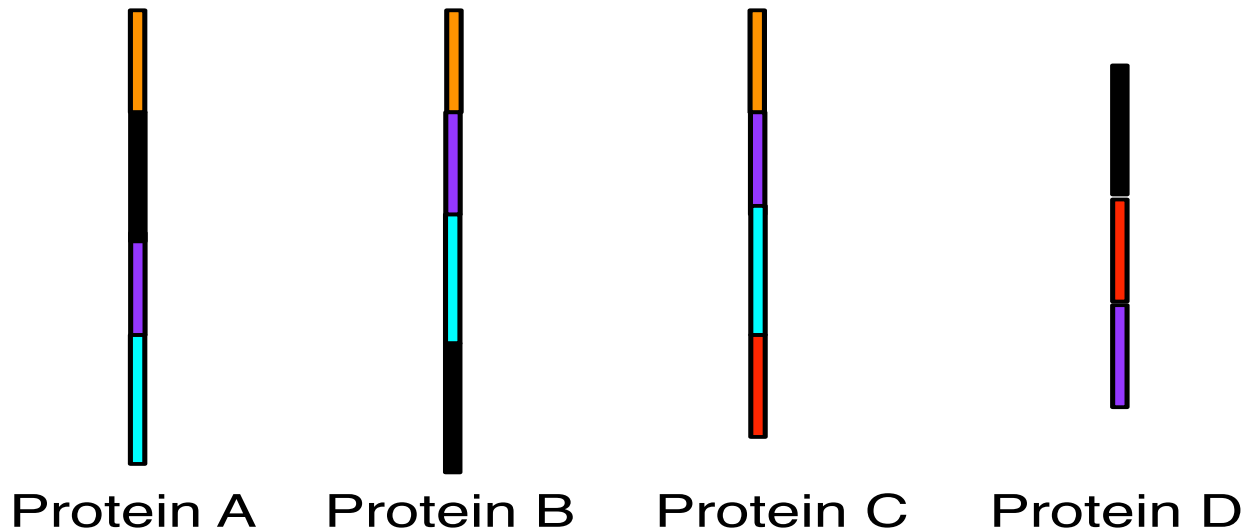
Protein C

Protein D

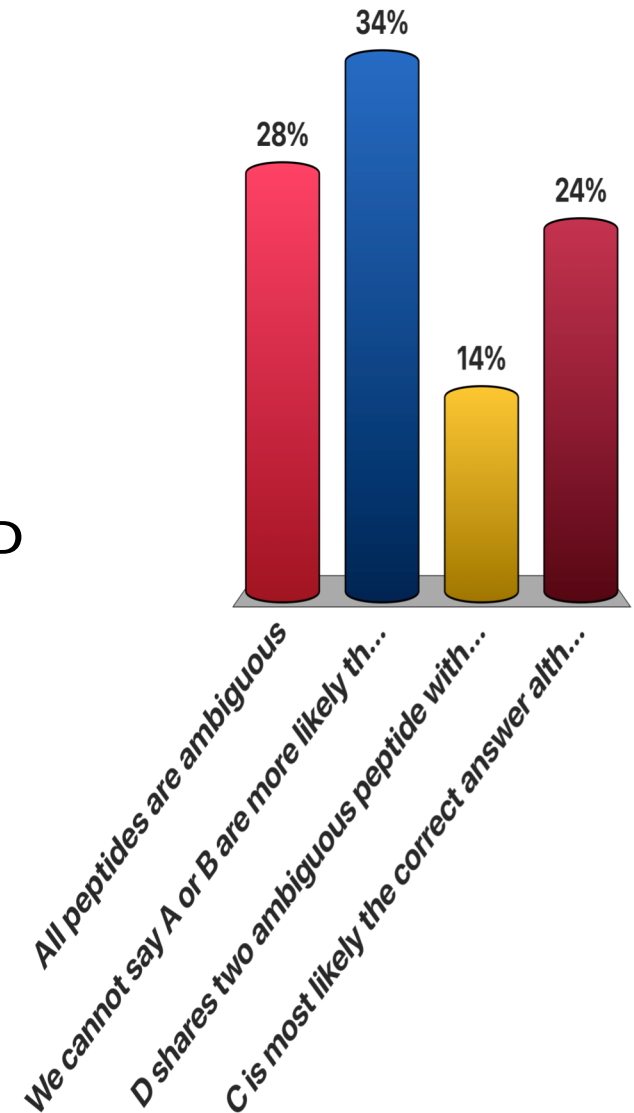
Try to order the proteins from highest to lowest confidence ?

Which of the following statements regarding ambiguous peptides could potentially be true

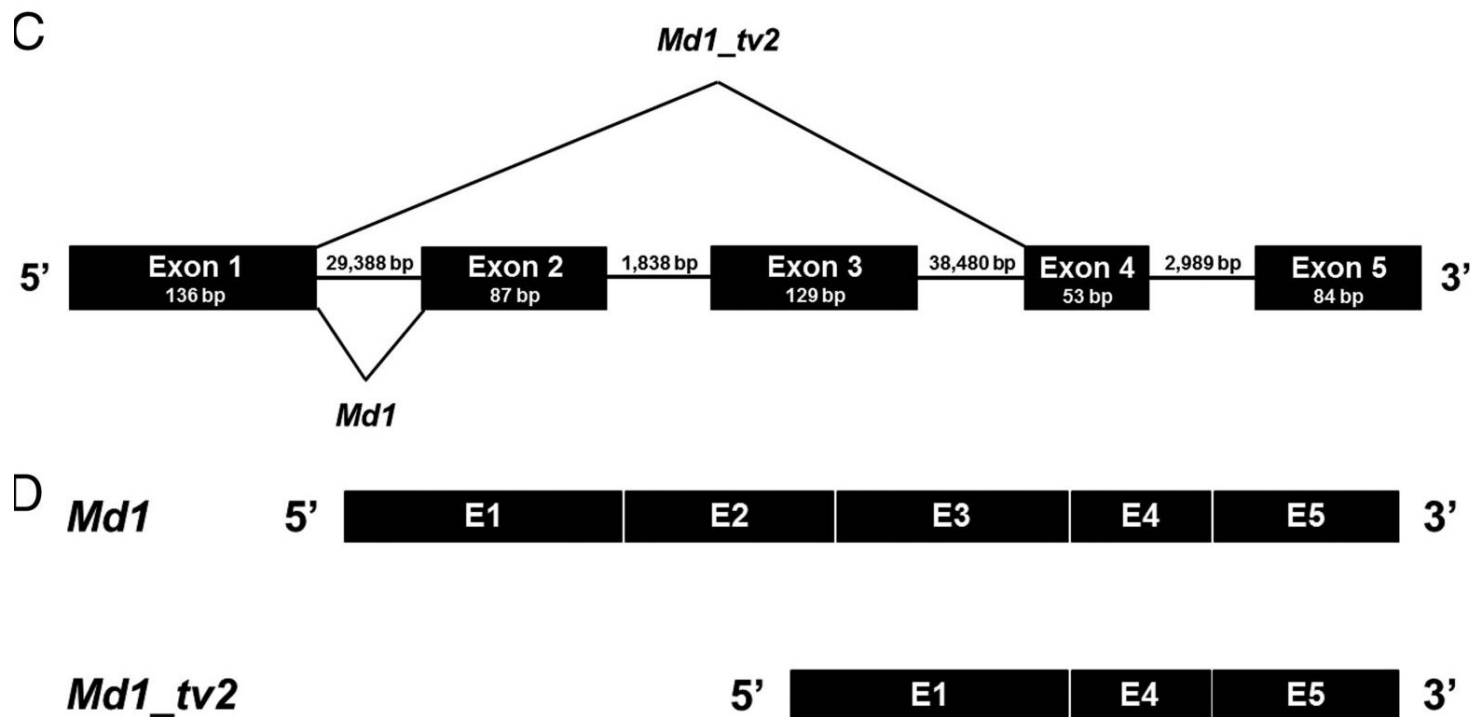
known protein sequences



- A. All peptides are ambiguous
- B. We cannot say A or B are more likely than the other because they share all same peptides
- C. D shares two ambiguous peptide with C, and one with A/B.
- D. C is most likely the correct answer although it comprises only ambiguous peptides



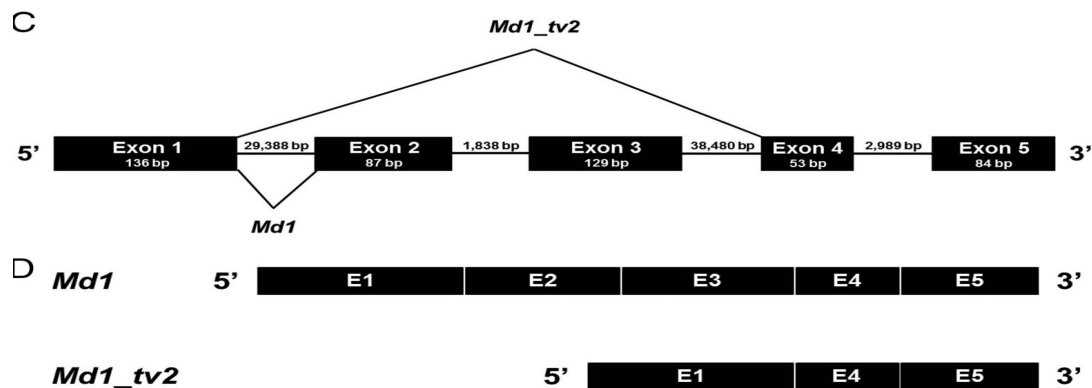
Representative sequences and splice variants



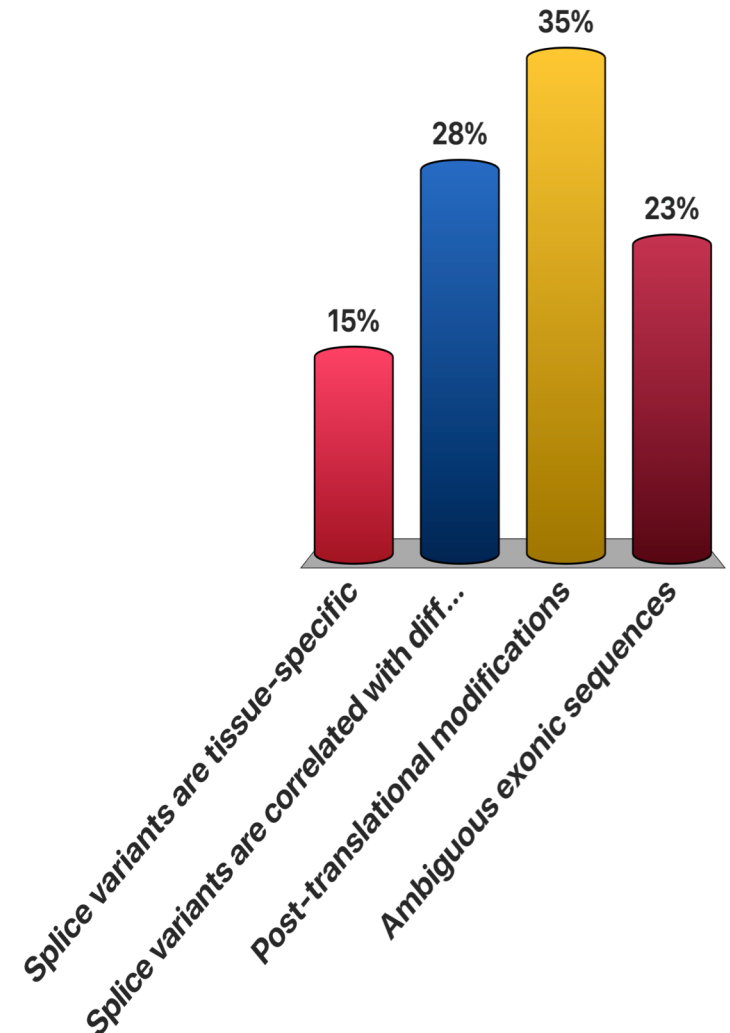
Splice variants are subsets of the full sequence

In proteomics, we usually take the full sequence as the “search sequence”

The following are potential problems when using full sequence for library search

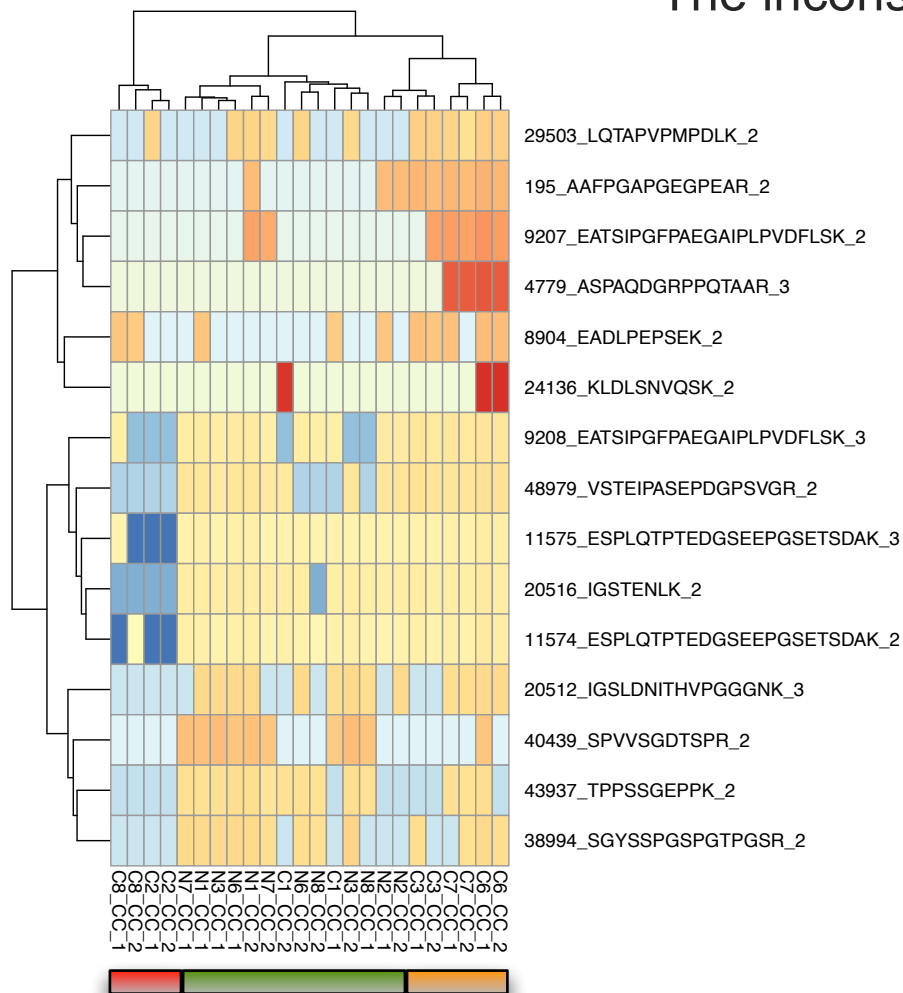


- A. Splice variants are tissue-specific
- B. Splice variants are correlated with different phenotypes
- C. Post-translational modifications
- D. Ambiguous exonic sequences



Food for thought --- splice variants

The inconsistencies corresponded to splice variants



MAPT, P10636

Exon4 - AEEAGIGDTPSLEDEAAGHVTQ^

Exon5 -
 EPESGKVVQEGFLREPGLSHQLMSGMPGAPLLPEGPREATRQPSGTGPEDTEGGR
 HAPPELLKHQLLDLHQEGPPLKGAGGKERPGSKEEVEDDRDVESSPQDSPPSK**ASPA**
QDGRPPQTAAREATSIPGFPAEGAIPLVDFLSK**VSTEIPASEPDGPSVGR**AKGDAPL
 EFTFHVEITPNVQKEQAHSEEHLGR**AAFFGAPGEGPEAR**GPSLGEDT**KEADLPEPSEK**
 QPAAAPRGKPVSRVPQLK^

Exon6 - ARMVSKSKDGTGSDDKKAK^

Exon7 -
 TSTRSSAKTLKNRCLSPKHPTPGSSDPLIQSSPAVCPEPPSSPKYVSSVTSRTGSSGA
 KEMK**LK**^

Exon8 -
 GADGKTKIATPRGAAPPQKGQANATRIPAKTP**PAPKTPSS**^

Exon9 -
GEPKSGDRSGYSSPGSPGTPGSRSRTPSLPTPPTREP**KKVAVVRTPPKSPSSAKSRLQ**
TAPVMPDLKNV**KSKIGSTENLKHQ**GGGK^

Exon10 -
 VQ**IINKKLDLSNVQSK**CGSKDNIK**HVP**GGGS^

Exon11 -
 VQ**I**VYK**PV**DLSK**VTSKCG**SLGNI**HHK**P^

Exon12 -
 GGGQ**VEVKSEK**LD**FKDRVQSKIGSLDNITHV**PGGG**NK**K^

Exon13 -
 IETH**KLTFRENAKAKTDHGA**EIVY**KSPVSGDTS**PR**HLSNV**SSTGSID**MVDSPQLATLA**
 DE**VCASLAK**CI

In cancer, a gene marker MAP was reported to be down-regulated in severe cancer but this gene has 11 splice variants

Which of the following is most similar to piecing protein information from spectra

- A. Putting together one jigsaw puzzle
- B. Putting together many jigsaw puzzles with the pieces mixed up
- C. Putting together many jigsaw puzzles with the pieces mixed up, and some pieces missing
- D. Putting together many jigsaw puzzles with the pieces mixed up, some pieces missing, and some reference picture boxes are wrong

Which of the following is incorrect regarding the Protein Identification through library Search?

- A. MS characterization of proteins is highly dependent on bioinformatic analysis
- B. Bioinformatics programs can be used to search for the identity of a protein in a database of theoretically digested proteins
- C. Even in reality, the protease digestion is always perfect in MS
- D. The purpose of the database search is to find exact or nearly exact matches

Representation of proteomics data

- Many format types
- Very complex
- But can be categorized

Format	Vendor/Creator	Comments
Wiff	Applied Biosystems	Proprietary
RAW	Thermo Fisher	Proprietary
raw	Waters	Proprietary
d	Agilent	Proprietary
dta	-	Text-based format
MGF	Matrix Science	Text-based format
mzData	PSI	Markup language; Superseded by mzML
mzML	PSI	Markup language; Current
mzXML	ISB	Markup language; Superseded by mzML

Based on what you know... how many different data formats exist for genomics data?

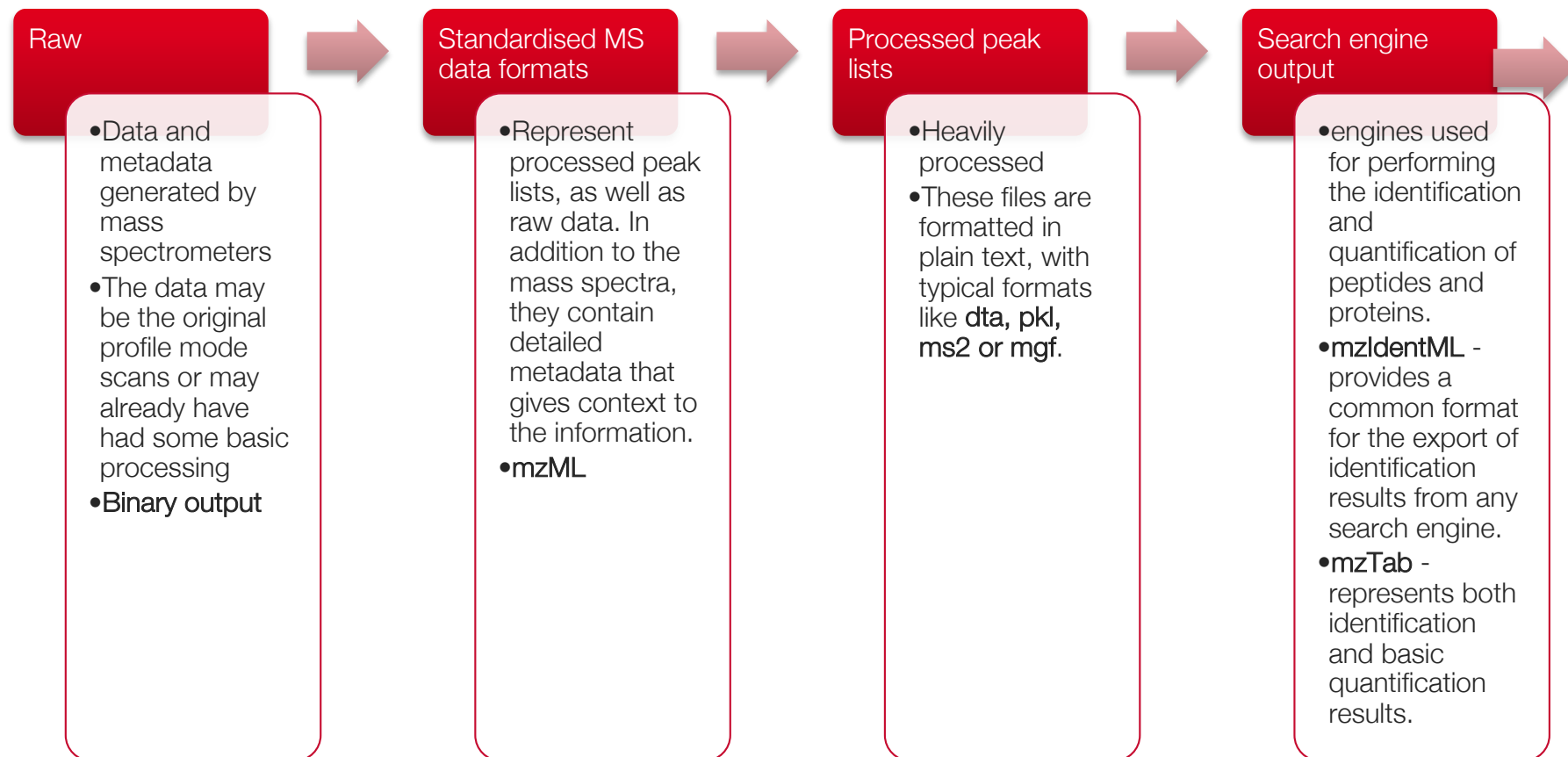
Let's see what happens when we are messy

<https://www.youtube.com/watch?v=1LfgbB0Mcqo>

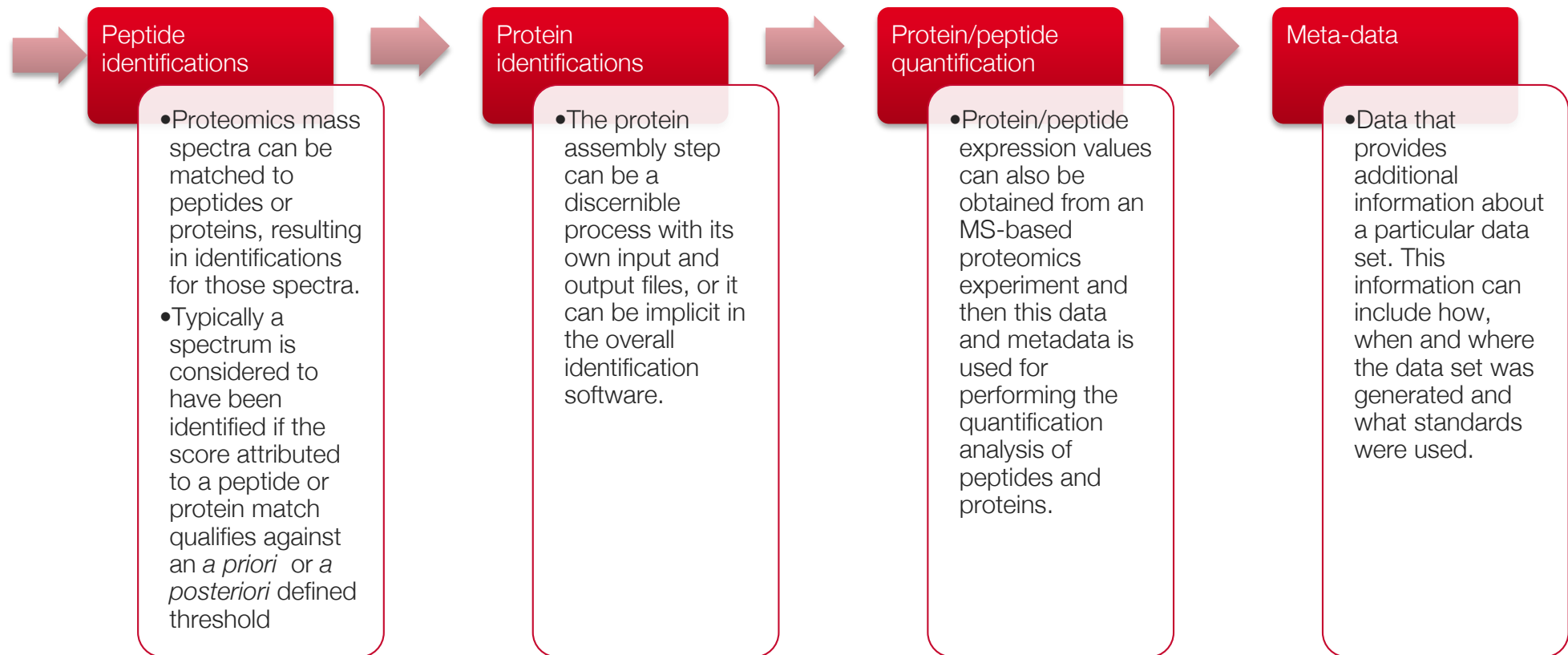
Organizing and formatting proteomics data

- Proteomics data is complex
- We need to organize data so that others can easily read and get the information they need

Proteomics data has many layers



Proteomics data has many layers



List-based formats

Some examples:

mgf (mascot generic file)

```
BEGIN IONS
PEPMASS=406.283
CHARGE=2+,3+
TITLE=Experiment_1
145.119100 8
217.142900 75
409.221455 11
438.314735 46
567.400183 24
714.447552 31
116.113400 72
91.2165000 32
405.288933 94
39.3021000 12
549.379462 21
715.466300 81
15.1098000 62
45.1358430 28
```

Parameters

pkl (peak list)

```
814.27 22673800 1
221.06 2529.3
223.84 220.9
226.91 1026.9
227.97 1037.9
231.06 110.6
239.05 7193.1
239.74 2513.3
240.27 363.4
240.79 1314.7
241.45 629.9
254.85 332.5
259.71 200.5
260.93 2437.7
```

mz intensity

Which list-based format is more informative? What information can you tell from these lists?

Non-rich descriptions
Lacking meta-data

mz intensity

List-based formats are easy to read?

A. True

B. False

List-based formats are an efficient storage of information?

A. True

B. False

List-based formats are information-rich (i.e., allow me to store a lot of data on running parameters, machines used, etc etc?)

A. True

B. False

Open-source formats

- Proprietary formats locks proteomics analysis to packaged software with instrument platform.
- Not all software are available
- Solution: Standardize formatting to open-source format
- All open source formats are XML-based.

What is XML?

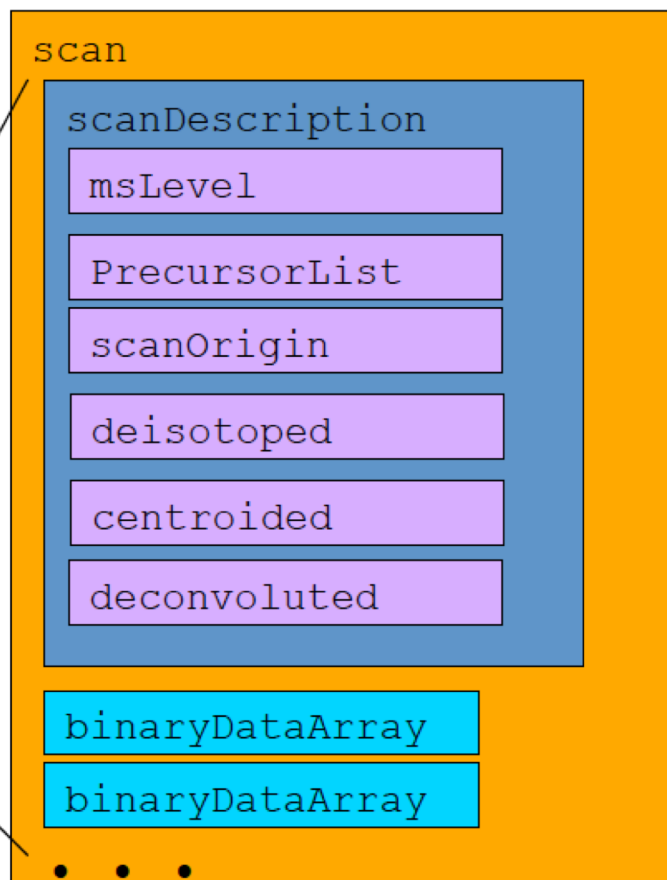
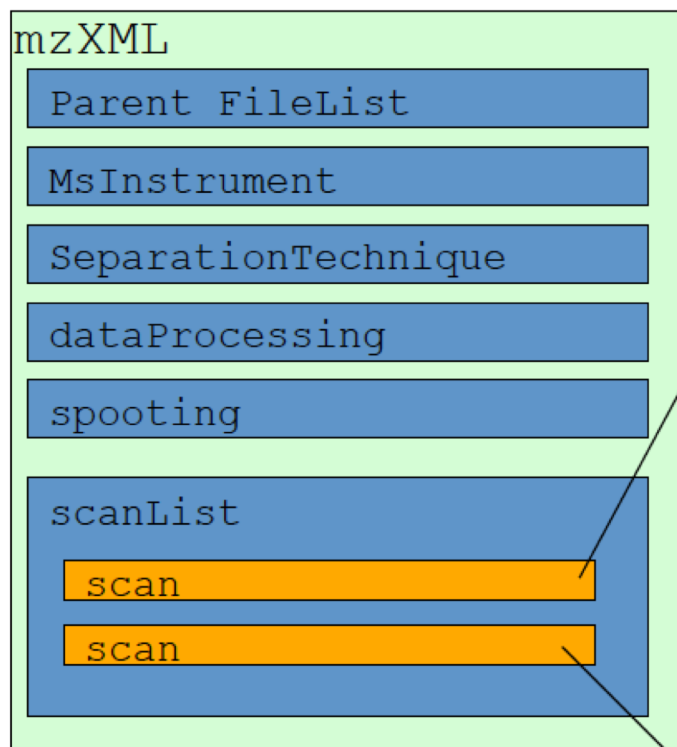
- Markup language
- Controlled vocabulary
- Have you heard of HTML?
- Features you need to know:
 - Controlled structure
 - Nesting (multi-layered information)
 - Meta-data
 - Machine and human-readable syntax

Do not memorize the formats!

You only need to appreciate why it is better than list-based formats and basic features

mzXML

Rich meta-data



mzXML was the first xml based file format developed for proteomics experiments. It was developed by the System Biology Group, USA.

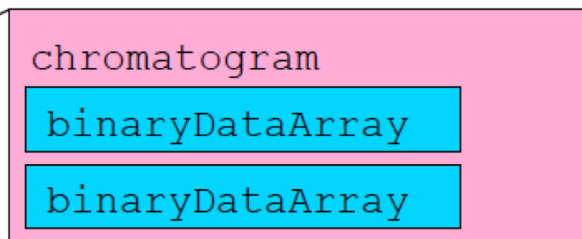
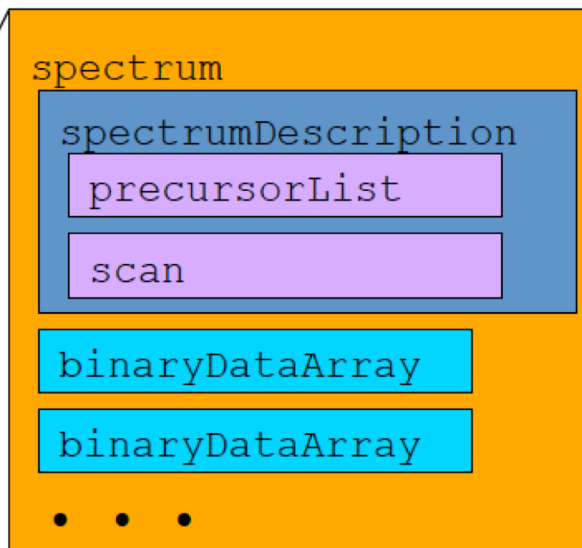
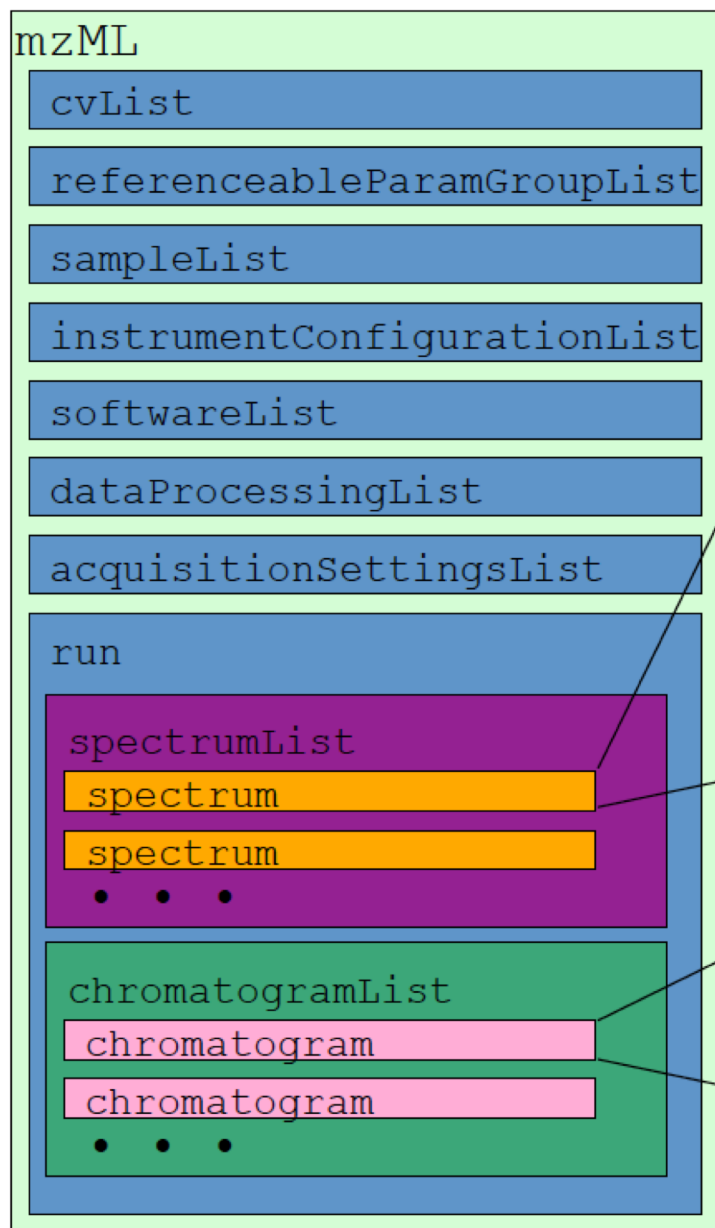
The annotations in the file are string based. It means, they are in this way: (**Name Attribute, Value**).

Do not support chromatograms information.

Is very difficult to extend. The structure of the file don't allow to define new parameter or features for each elements. For example, msInstrument are defined only by the name of the instrument. Also, if the spectrum is preprocessing with any program, is difficult to incorporate the information.

Actually exist more than 4 versions of the schema. The schema is supported by the System Biology Group, USA-Zurich.

mzML



Meta data about the spectra plus all the spectra themselves.

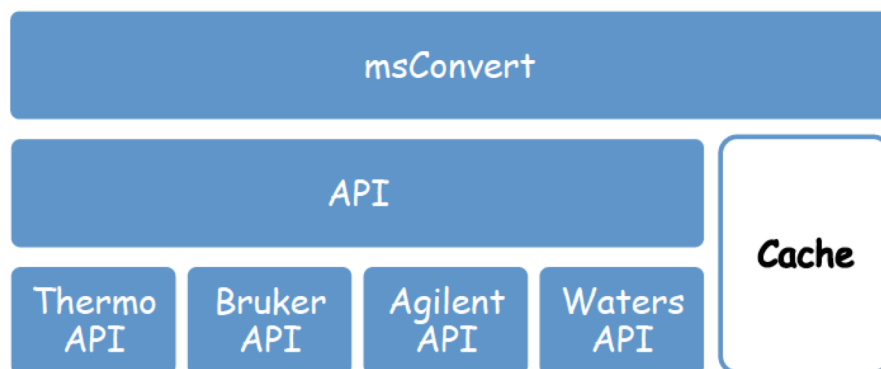
The header at the top of the file encodes information about: **the source** of the data as well as information about **the sample**, instrument and **software** that processed the data.

Cvterms are used to define the metadata and the properties of each element (**software**, **instrument**, **sample**, **scansetting**, etc.

Chromatograms may be encoded in mzML in a special element that contains one or more cvParams to describe the type of chromatogram, followed by two base64-encoded binary data arrays.

Format conversion

If you want to use open source or analyze data from different platforms,
Then need to standardize format (to MZXML or MZML). Use **ProteoWizard**



File Input Supported:

- Thermo
- Bruker
- Agilent
- Waters
- Pkl
- mgf,
- dta
- ms2

Includes both
proprietary and
open formats

File Output Supported:

- mzML
- mzXML
- mzData
- Pkl
- mgf

Cross-platform !!!!

BIOINFORMATICS APPLICATIONS NOTE Vol. 24 no. 21 2008, pages 2534–2536
doi:10.1093/bioinformatics/btn323

Genome analysis

ProteoWizard: open source software for rapid proteomics tools development

Darren Kessner^{1,*}, Matt Chambers², Robert Burke¹, David Agus¹ and Parag Mallick^{1,3,*}

¹Spielberg Family Center for Applied Proteomics, Cedars-Sinai Medical Center, ²Department of Biochemistry, Vanderbilt University, Nashville, TN and ³Department of Chemistry & Biochemistry, University of California, Los Angeles, CA, USA

Received on April 18, 2008; revised on May 21, 2008; accepted on June 18, 2008

Advance Access publication July 7, 2008

Associate Editor: John Quackenbush

<http://proteowizard.sourceforge.net/>

Open analysis platforms

- OpenMS (<https://www.openms.de/>)
- Galaxy
(<https://galaxyproject.org/proteomics/>)
- ExPASy (<https://www.expasy.org/>)
- TransProteomics Pipeline
(<http://tools.proteomecenter.org/>)
- Maxquant
([http://www.biochem.mpg.de/5111795/
maxquant](http://www.biochem.mpg.de/5111795/maxquant))

Not quite there yet... what are we missing?

Samples

Proteins

The image shows a Microsoft Excel spreadsheet titled 'nm.3807-S4.xls [Read-Only] [Compatibility Mode] - Microsoft Excel'. The spreadsheet contains a large table of data with columns labeled A through AC. Row 1 is the header row, with columns A-C containing protein information (GeneSy, mbol, kidneyTisue1) and columns D-AC containing quantitative measurements (ue2-ue27). The data is organized into a grid with various cell styles and colors. The status bar at the bottom indicates 'Ready', 'EN', and the date '6/28/2016'.

We haven't talked about quantitation!



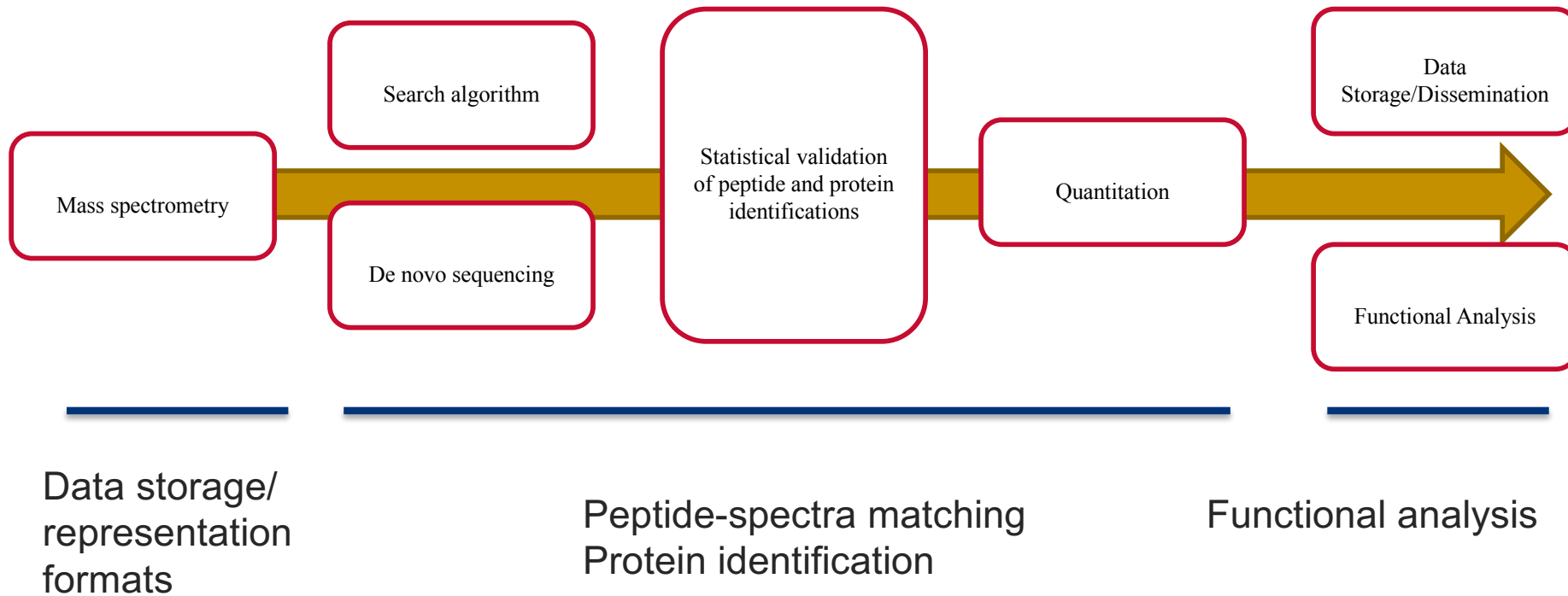
What have we learnt?

- Library search algorithms match theoretical spectra with observed spectra
- De novo sequencing algorithms use graph-theory methods to join observed spectra -> find longest path
- In proteomics, we observe peptides first. Proteins are inferred based on which peptides are detected
- A lot of useful information is lost during peptide to protein transition
- Proteomics data formats are diverse

You should now be able to

- Describe what an algorithm is, and how it differs from a heuristic or a computer program
- Describe the various levels of spectra data and their derivations in MS-based proteomics
- Describe the steps of a library search algorithm
- Describe the steps of a *de novo* sequencing algorithm
- Describe how peptides are assembled to proteins and associated problems
- Describe and evaluate the various levels and data representation formats in proteomics

Putting things into perspective



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.”

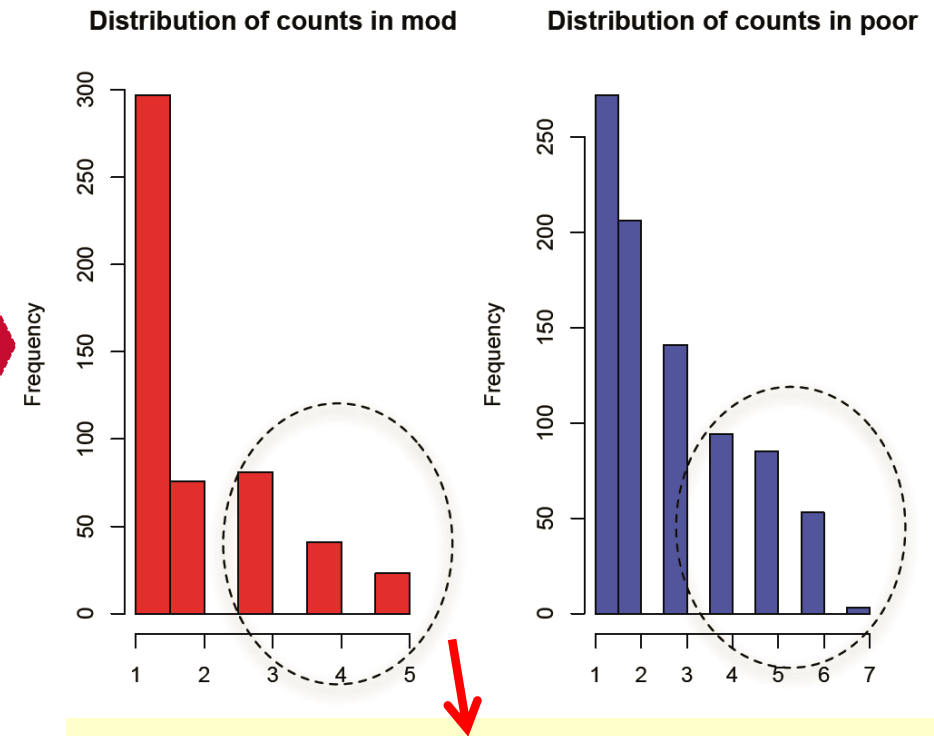
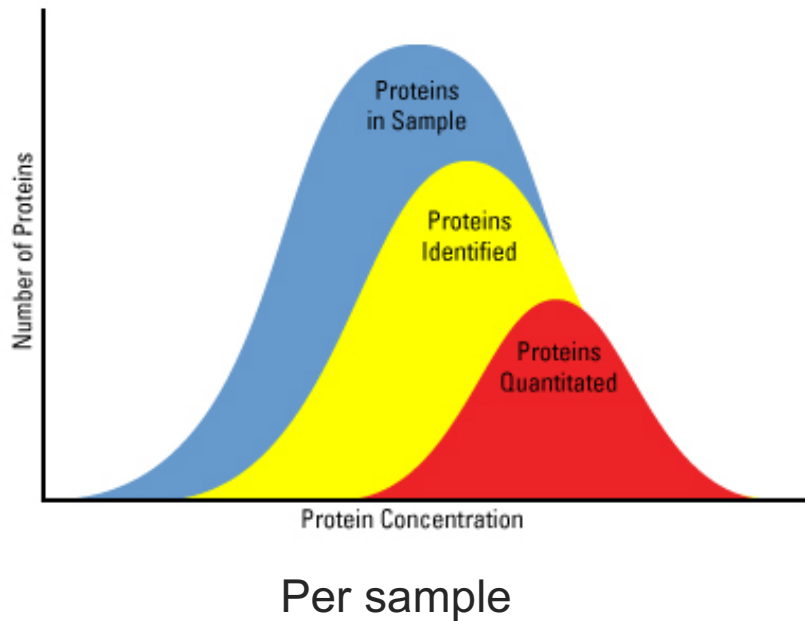
Steen & Mann. **The ABC's and XYZ's of peptide sequencing**. *Nature Reviews Molecular Cell Biology*, 5:699-711, 2004



Issues in proteomics: Coverage and consistency

Technical incompleteness

How it affects real data



Only 25 out of 800+ proteins are common to all 5 mod-stage HCC patients!

Going further

- There is a kind of proteomics (Data-Independent Acquisition) where there is no one-to-one correspondence between MS1 and MS2. What kind of problems do you think can happen. And why do you think they avoided collecting MS1?

Hint: In traditional proteomics, MS1 peaks are semi-randomly selected for MS2 to ensure one-to-one correspondence. But what is the unintended consequence of this procedure?

Readings (Encouraged)

- 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

Readings (Additional)

- Goh and Wong Spectra-first feature analysis in clinical proteomics—A case study in renal cancer, *JBCB*, 14 (05), 1644004, 2016
- 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