

De la fiabilité des données d'identification et de quantification de protéines par MS

Myriam Ferro
BGE/EDyP Laboratory
CEA/Grenoble
29/11/2012

Atelier Prospectom 29-30 novembre 2012, Grenoble

Definition of proteomics

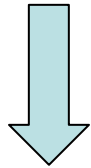
GENES



TRANSCRIPTS (RNA)



PROTEINS



Genomics



Transcriptomics



Proteomics

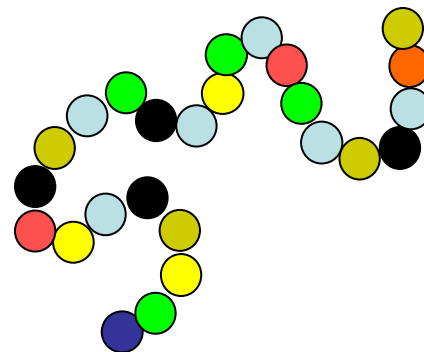


Dynamics

(differential expression: localization, time)



- Processing
- PTMs
- Localization
- Partners
- Etc.



Proteomics: what's for ?

Some questions

- mutation or environmental conditions
- protein-protein interactions
- searching for biomarkers
- etc.

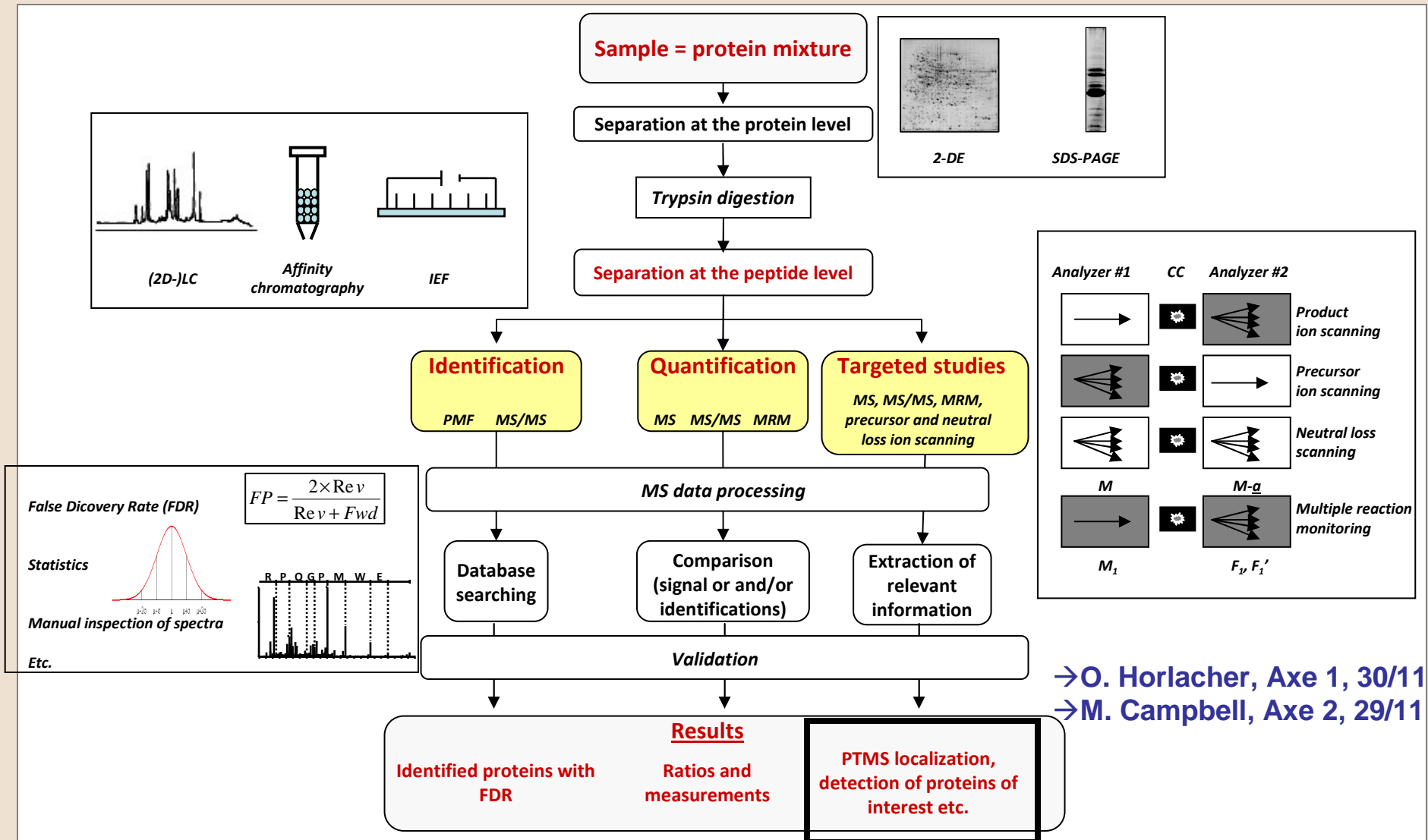


Data visualization

New knowledge and questions !

From Käll & Vitek, 2011

Proteomics analyses using MS



→ O. Horlacher, Axe 1, 30/11
 → M. Campbell, Axe 2, 29/11

Some reminders about mass spectrometry measurements



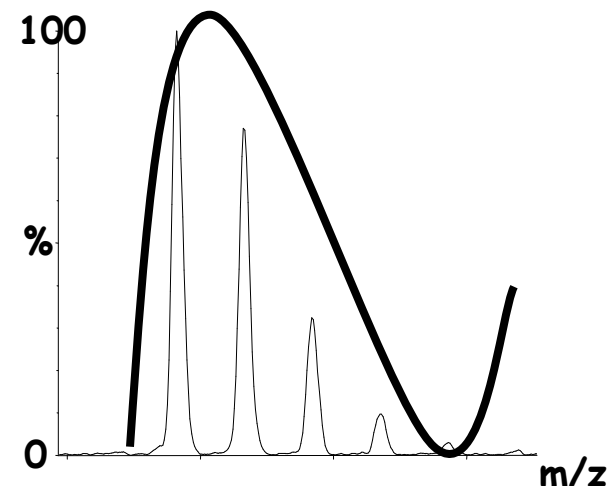
Reminder # 1: the molecular mass

- **Ex: EWMPGQPR = C₄₄ H₆₅ N₁₃ O₁₂ S**

- **Nominal mass (approximation)**

C=12; H=1; N= 14; O=16; S=32

→ M = 999



- **Monoisotopic mass (the more stable isotopes)**

C=12; H=1.007825; N= 14.003074; O=15.9949146; S=31.9720718

→ M = 999.4596

- **Average mass (barycentre of all masses)**

C=12.011; H=1.007994; N=14.006674; O=15.9994; S=32.066

→ M= 1000.1464

Reminder # 3: m/z ratio

- **M = molecular mass (ex: peptide)**
- **H⁺ = proton**
- **z = state of charge**

$$m/z = \frac{M + z H^+}{z}$$

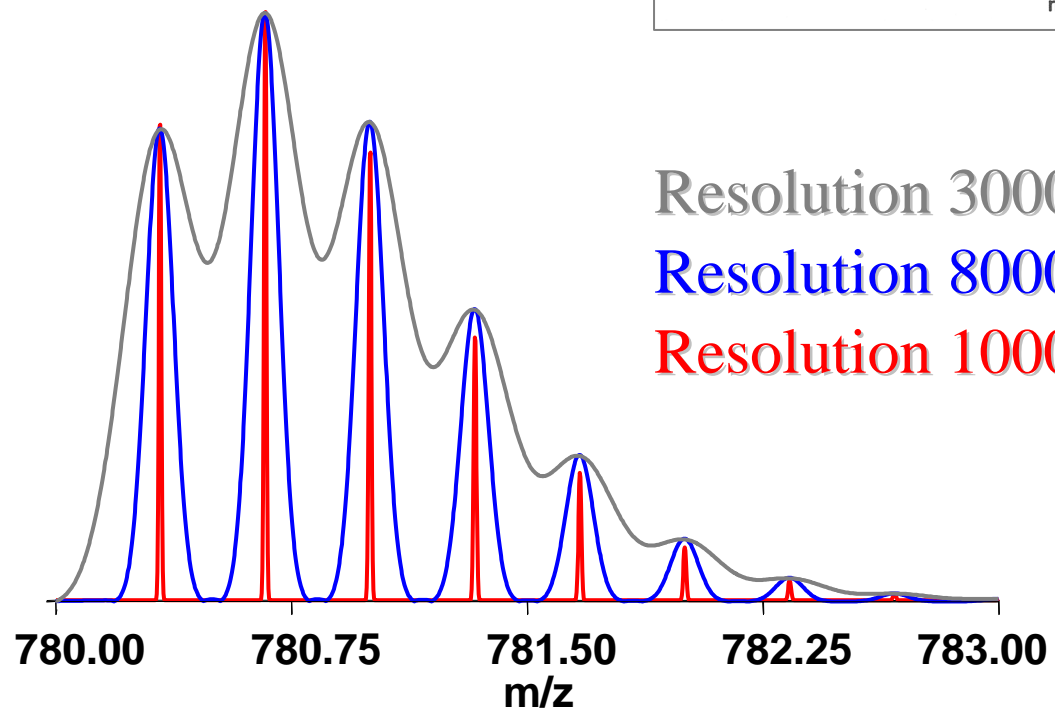
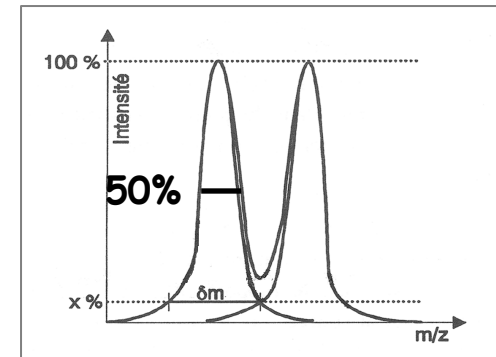
Reminder # 4: resolution

- **Characterize peaks separation**

(FWHM: Full Width at Half Maximum)

$$R = \frac{\Delta m}{m}$$

Insulin (Chain A) 3+



Resolution 3000

Resolution 8000

Resolution 100000

Accuracy and precision



High precision but low accuracy



High precision but high accuracy

Precision is related to reproducibility → careful control of instrument settings

Accuracy is related to calibration and resolution

Instrument with high resolving power (e.g OrbiTrap VELOS)

Resolution = 60 000

Mass accuracy (Orbitrap analyser) ~ 1-5 ppm

$$\frac{|M_{\text{exp}} - M_{\text{th}}|}{M_{\text{th}}}$$

ppm = part per million

Features of mass spectrometers

Table 1 Performance comparisons of the mass spectrometry instruments

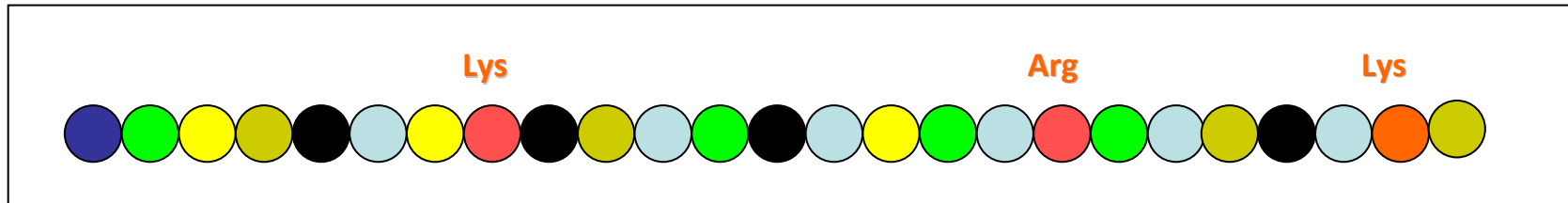
| Instrument | Applications | Resolution | Mass accuracy | Sensitivity | Dynamic range | Scan rate |
|--------------------|---|------------|---------------|-------------|---------------|----------------|
| LIT (LTQ) | Bottom-up protein identification in high-complexity, high-throughput analysis, LC-MS ⁿ capabilities | 2000 | 100 ppm | Femtomole | 1e4 | Fast |
| TQ (TSQ) | Bottom-up peptide and protein quantification; medium complexity samples, peptide and protein quantification (SRM, MRM, precursor, product, neutral fragment monitoring) | 2000 | 100 ppm | Attomole | 1e6 | Moderate |
| LTQ-Orbitrap | Protein identification, quantification, PTM identification | 100,000 | 2 ppm | Femtomole | 1e4 | Moderate |
| LTQ-FTICR, Q-FTICR | Protein identification, quantification, PTM identification, top-down protein identification | 500,000 | <2 ppm | Femtomole | 1e4 | Slow, slow |
| Q-TOF, IT-TOF | Bottom-up, top-down protein identification, PTM identification | 10,000 | 2–5 ppm | Attomole | 1e6 | Moderate, fast |
| Q-LIT | Bottom-up peptide and protein quantification; medium complexity samples, peptide and protein quantification (SRM, MRM, precursor, product, neutral fragment monitoring) | 2,000 | 100 ppm | Attomole | 1e6 | Moderate, fast |

The strategies for protein identification



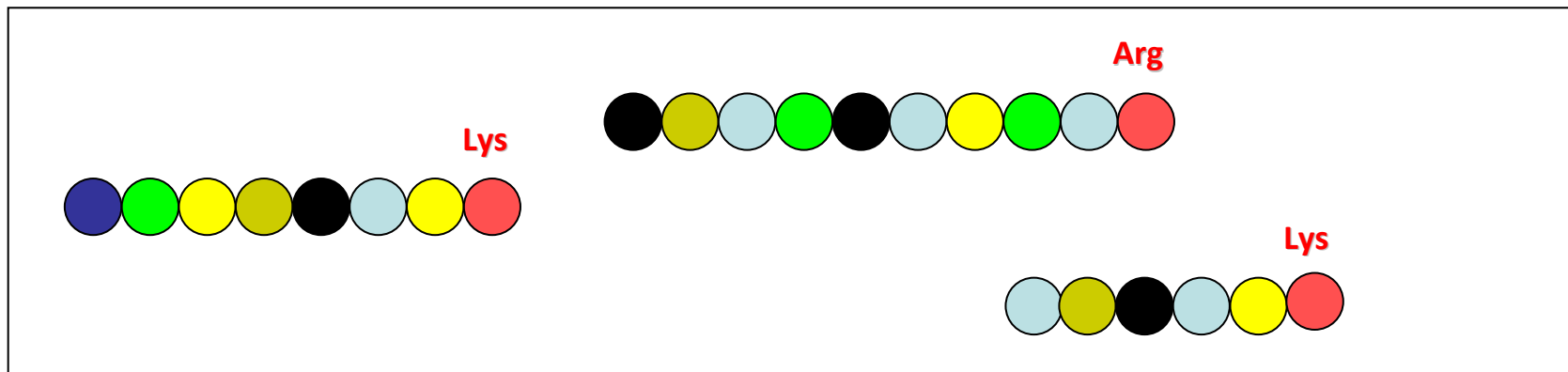
Trypsin digestion

Protein



Trypsin

Peptides



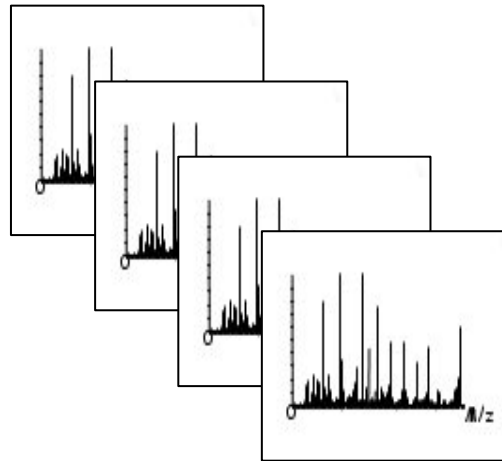
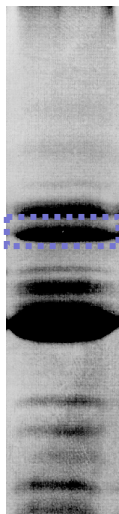
A proteomics workflow

Proteins

Tryptic peptides

MS(/MS) spectra
(peak picking)

Identification &
quantification results



```

ACTB_HUMAN          Mass: 43710   Score: 1390   Queries matched: 135
(P40709) Actin, cytoplasmic 1 (Beta-actin)
[ Check to include this hit in error tolerant search or archive report ]

Query   Observed   H(eqpt)   H(eqdb)   Delta   Miss   Score   Expect   Rank   Peptide
2552    795.37      794.36    794.47   -0.10   0       29      7.2    1   K.LIAPPER.K.1515
2372    850.19      849.18    847.42   -1.76   1       5       2.8e+02  0   K.CVDVIRK.D
2085    473.24      444.47    444.34   -0.08   0       51      0.0082  1   H.AVFPFIVGR.P.1972.3284.4001.4005
4475    449.49      916.97    915.44   -1.52   0       89      0.8e-06  1   K.AGPAAGDQGR.A.5132.4184.4183.4110.4113.4261
4704    459.88      997.75    997.48   -0.27   0       52      0.0053  1   H.DLIVDMK.I.4704.4704.4704.4707.4714.4
6631    567.16      1132.31   1131.52   -0.79   0       86      2.2e-06  1   R.GYSFTTAEK.R.6604.6613.6617.6618.6621.661
2164    386.27      1170.93   1170.36   -0.04   0       77      1.5e-05  1   H.HQVWVWQGR.D.2133.7155.7176.7181.7180.718
2189    549.47      1116.93   1116.61   -0.32   0       66      0.0023  1   H.DITALAPDHR.I.2040.2044.2046.2047.2049.201
2542    609.04      1198.87   1197.51   -0.55   0       78      1.3e-05  1   K.DSYVDEAGSR.R.2532.7549.7547.7556
13054   759.39      1516.77   1515.70   -0.07   0       46      0.022   1   K.QEYDESGPVIWR.H.13034.15044.13047
16511   883.06      1764.11   1763.77   -0.33   0       36      0.16    1   H.DDDIAALVWDSGCK.A.16947.16175.16516
16600   896.09      1790.17   1789.88   -0.28   0       103     3.1e-08  1   K.SYELDQGVYIYGRK.P.16619.16835.16847.168
17125   904.49      1806.77   1804.98   -1.79   2       8       1.1e+02  4   H.HQKITALAPDHR.K.17112
17846   624.47      1870.39   1868.80   -1.59   0       15      22     2   -.HDDIAALVWDSGCK.A.17923.18366
18649   977.67      1953.33   1953.06   -0.27   0       89      6.9e-07  1   H.VAFEDVPLIYAPLIRK.A.18635.18616.18645.1
18642   686.37      2092.89   2093.05   -0.97   1       7       1.1e+02  6   K.SYELDQGVYIYGRK.C
20625   1188.09     2314.17   2314.86   0.10   0       112     4e-09   1   H.DLVANTVLSGDTDFQIADR.H.20619.20624.2061
21381   782.07      2343.19   2342.16   -1.03   1       21      4.8    1   H.DLVANTVLSGDTDFQIADR.H
21469   789.34      2365.88   2366.25   -1.25   2       10      55     6   H.AVFPFIVGR.P.1972.3284.4001.4005
22621   948.26      2841.76   2842.32   -0.46   1       8       74     10  -.HDDIAALVWDSGCK.A.17923.18366
23220   1151.72     3452.14   3451.57   -0.57   1       2       2.4e+02  5   R.FKCPALQPSFLQESGDIHTFTFRSIRK.D.K
23280   1321.78     3962.32   3961.82   -0.50   2       6       72     2   H.CFALQPSFLQESGDIHTFTFRSIRK.D.K
    
```

Trypsin
Digestion

LC-MS(/MS)
analysis

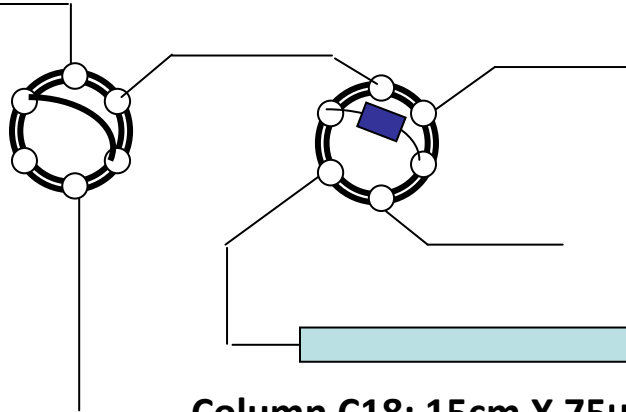
Database searching &
data analysis

Sample preparation
Biochemistry

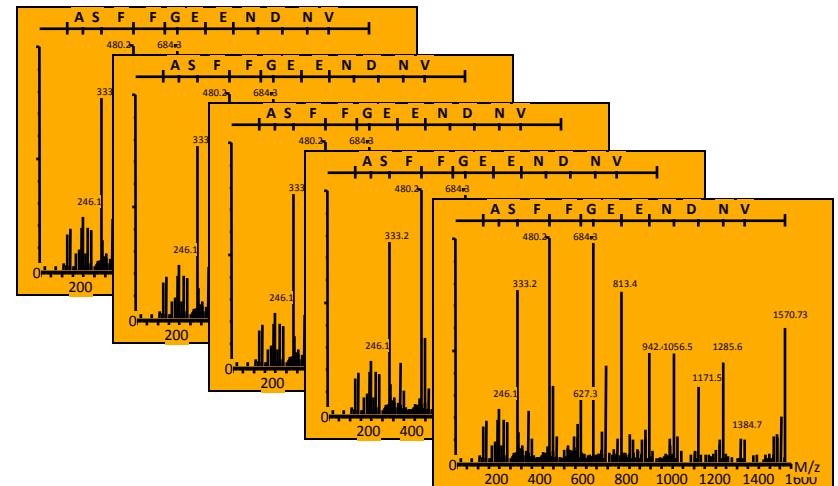
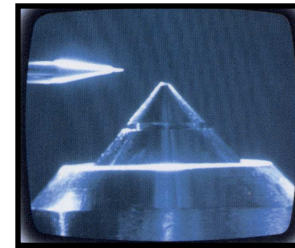
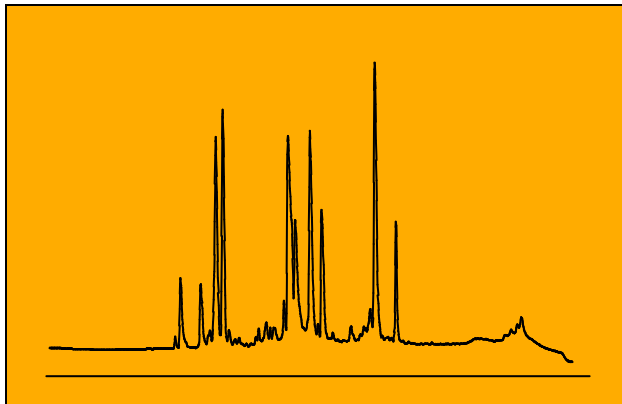
Mass spectrometry
Liquid chromatography

Informatics
Bioinformatics

Shotgun analyses



Column C18; 15cm X 75 μ mID
200 nL/min

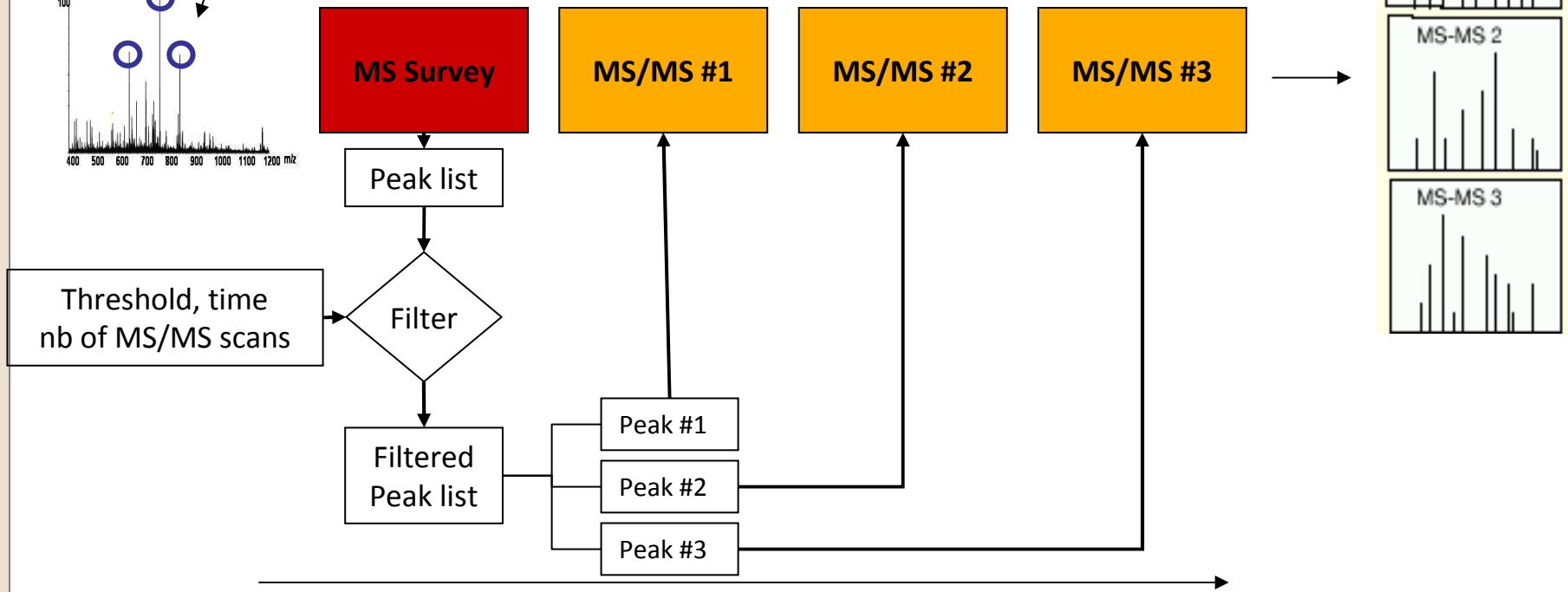
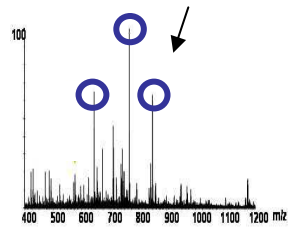
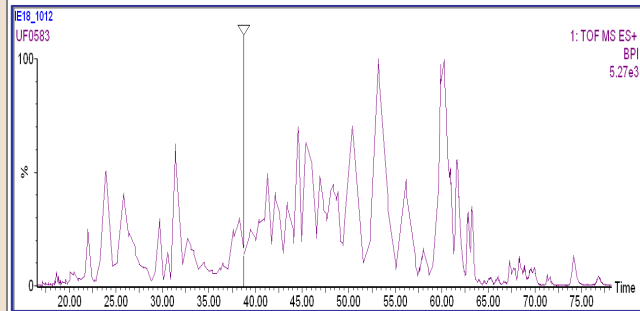


Ex: Orbitrap Velos ~ 5000 spc/hr

Separation and concentration

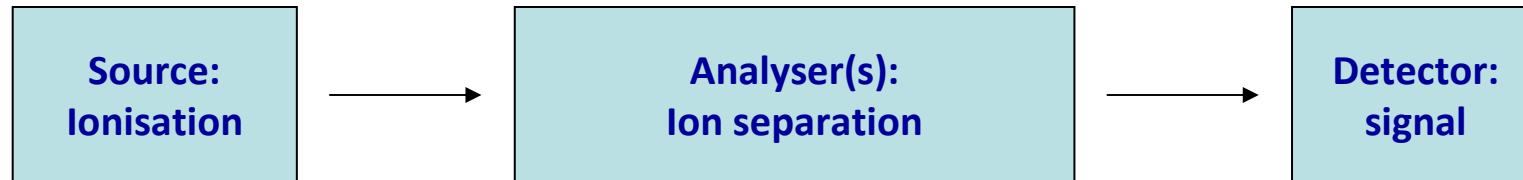
Automated MS/MS

Data dependent acquisition

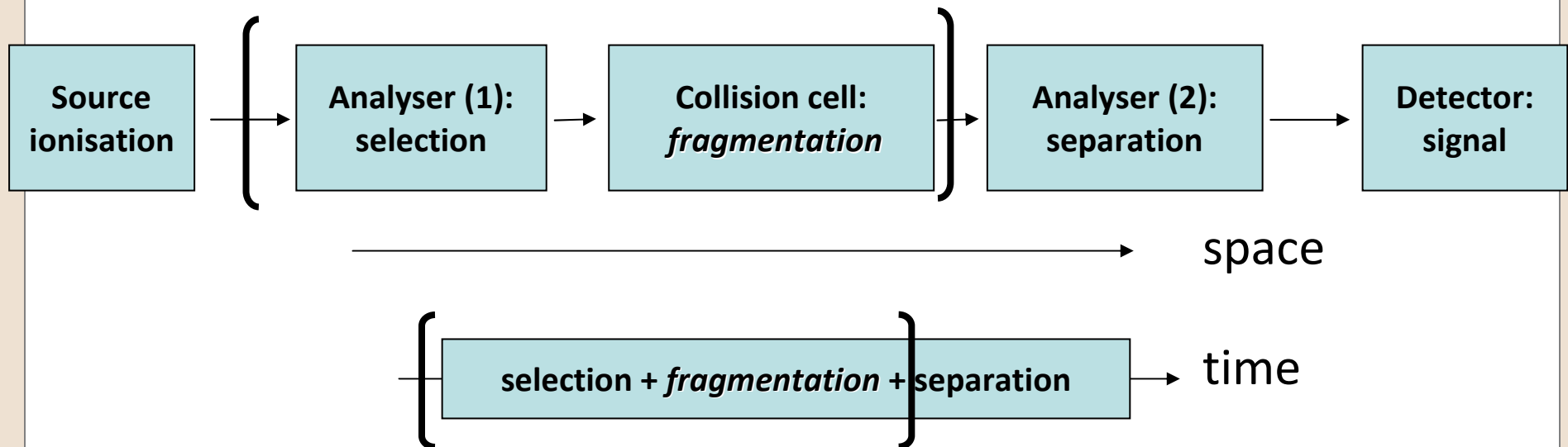


1-2 seconds (for the Orbitrap)

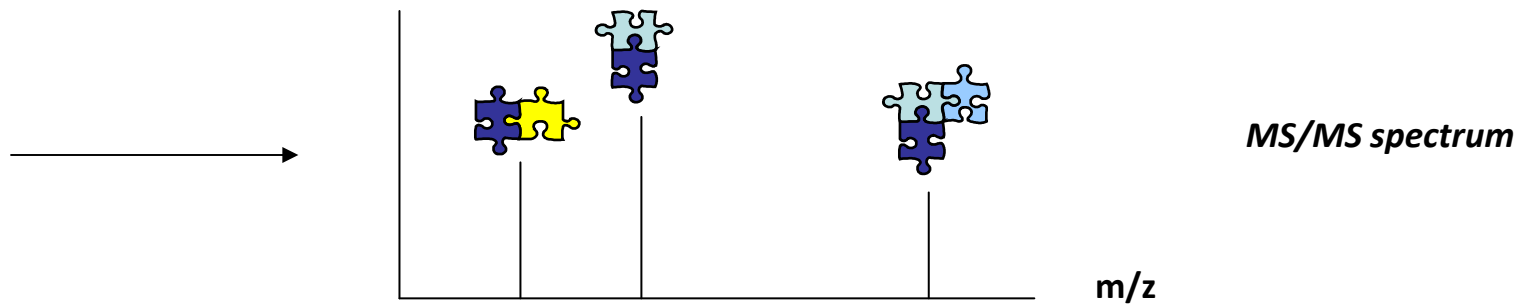
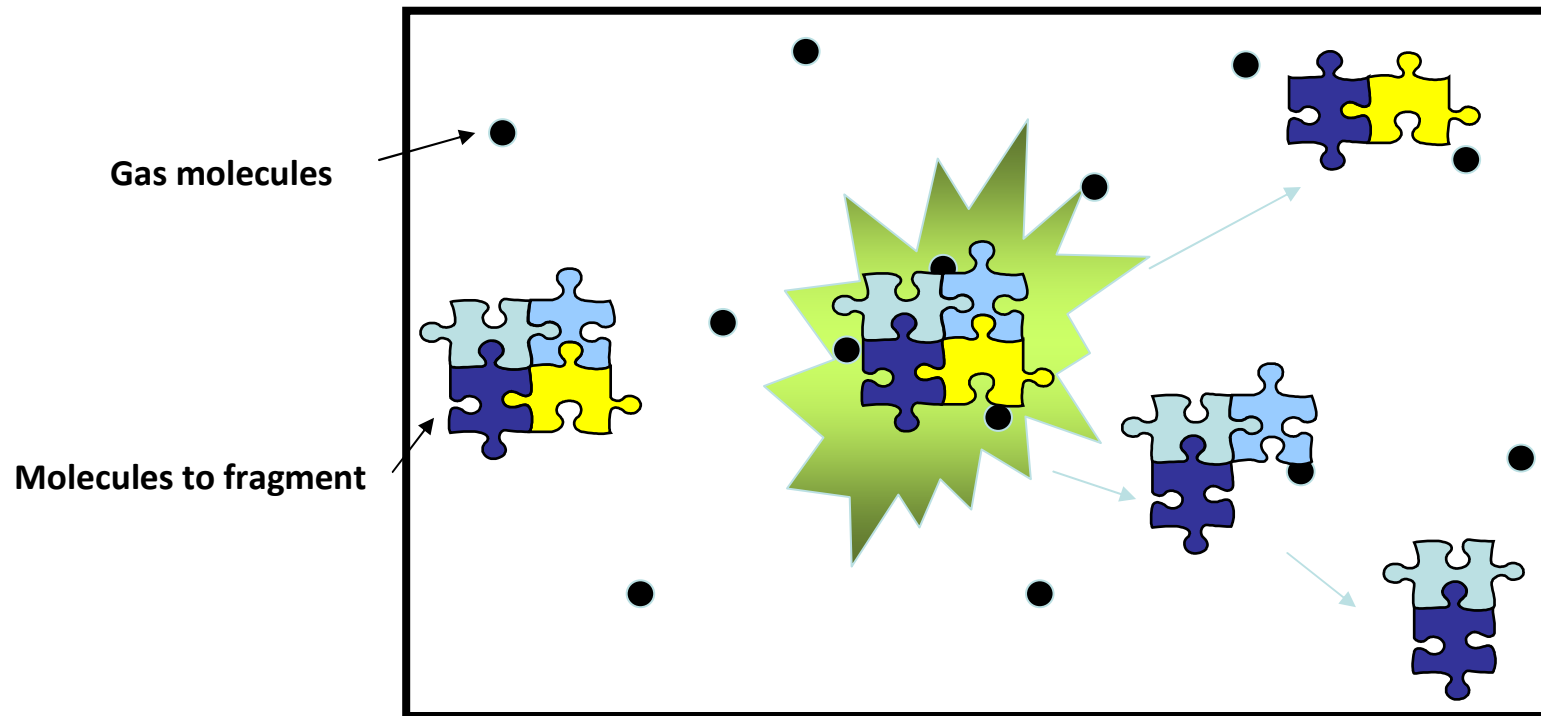
Tandem mass spectrometry



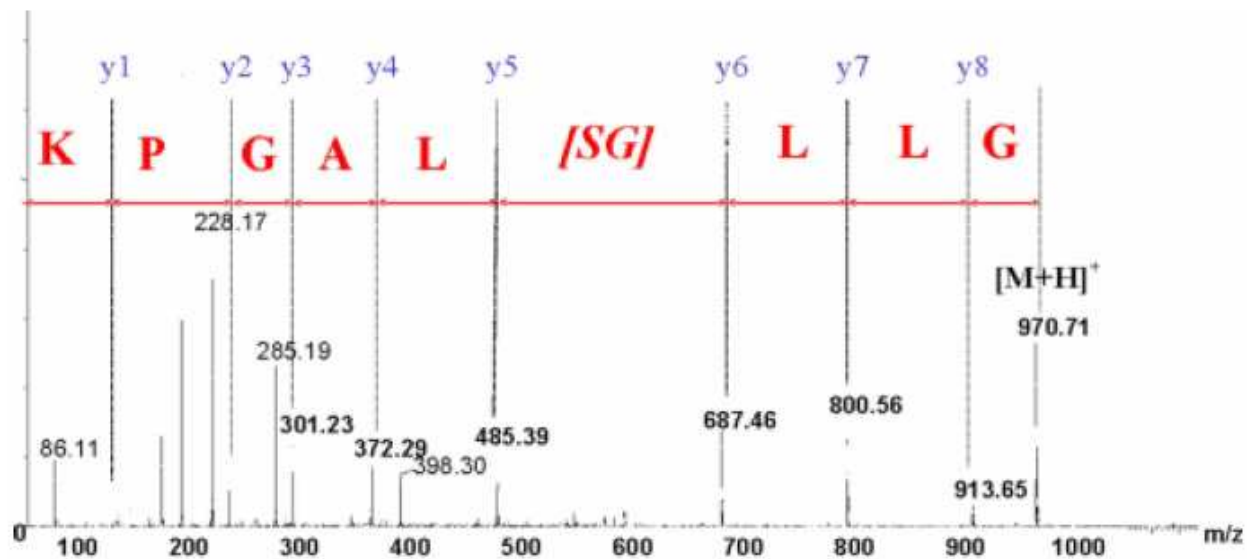
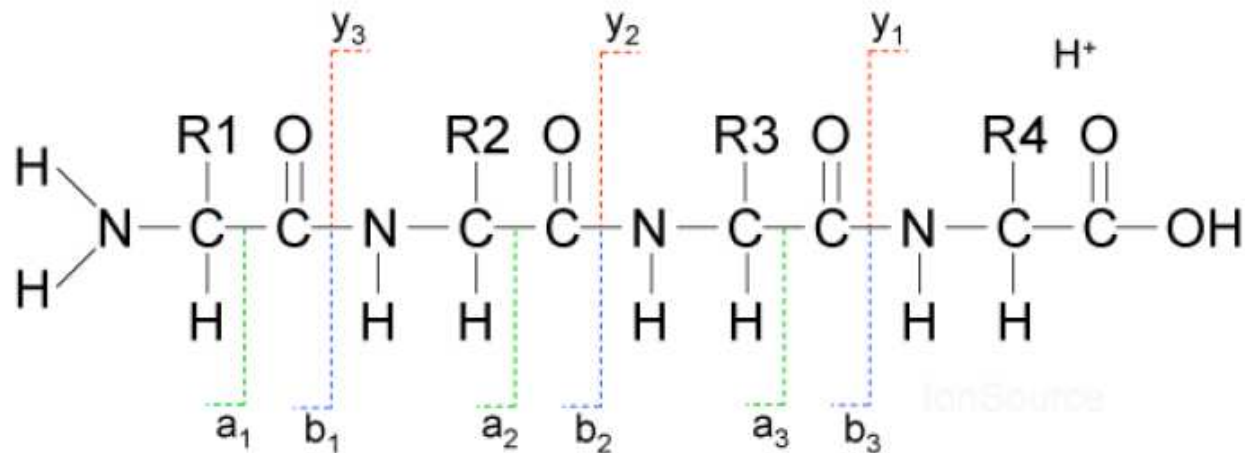
Tandem mass spectrometry (MS/MS) → peptide fragmentation



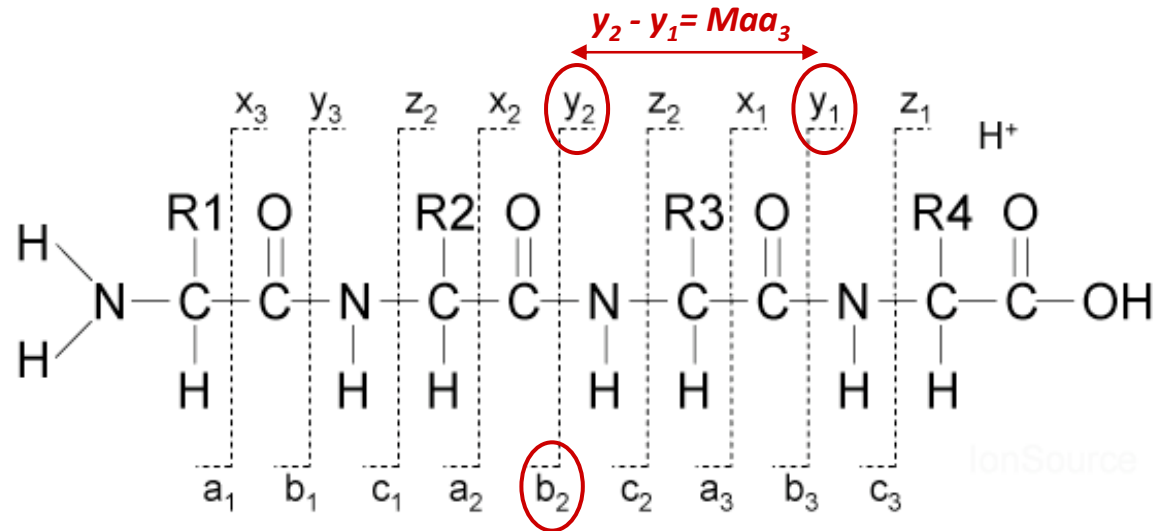
Peptide fragmentation in the mass spectrometer



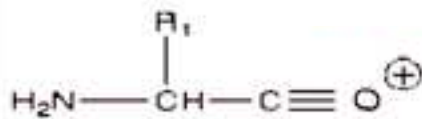
A typical MS/MS spectrum



Rules of fragmentation



Peptide bond cleaved: complementary ions y'' et b



b ion → peptide N-ter



y ion → peptide C-ter

Other ions

Immonium ions : specific of each amino acid ex: m/z 120 --> Phe

Internal fragments

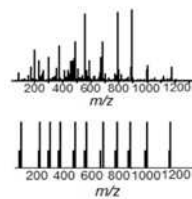
Specific fragmentations of PTMs (cf M. Jaquinod presentation)

Neutral loss ex: - 18 Da --> perte d'H₂O

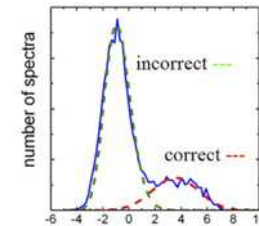
Schematic overview of a typical workflow of the proteomics informatics processing of a data set



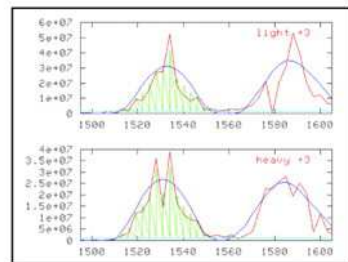
1. Conversion to and use of open data formats



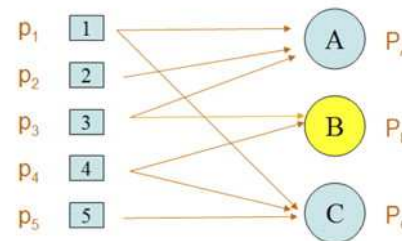
2. Spectrum identification with a search engine



3. Validation of identifications



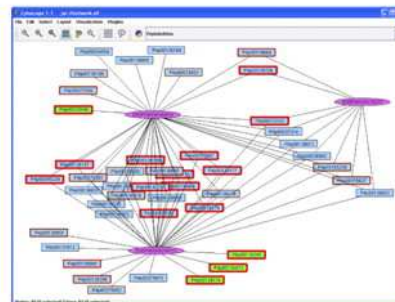
5. Quantification



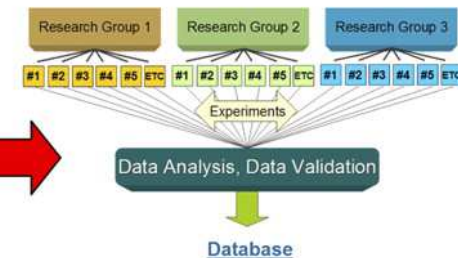
4. Protein inference



6. Organization in local data management systems

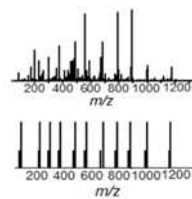
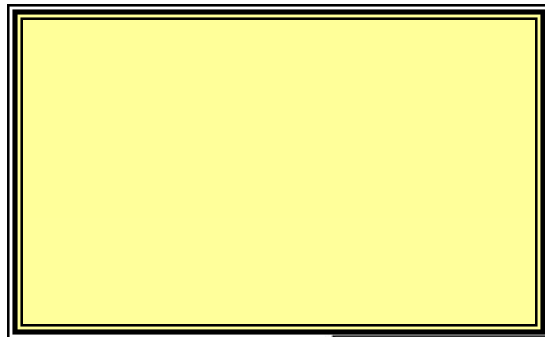


7. Interpretation of the protein lists

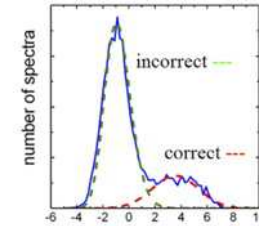


8. Transfer to public data repositories

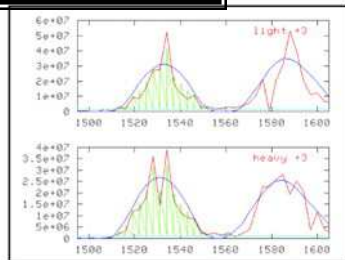
Schematic overview of a typical workflow of the proteomics informatics processing of a data set



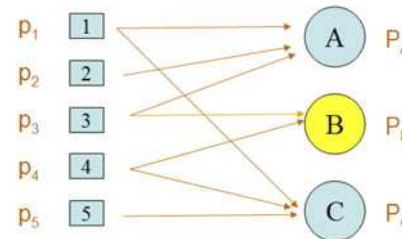
2. Spectrum identification with a search engine



3. Validation of identifications



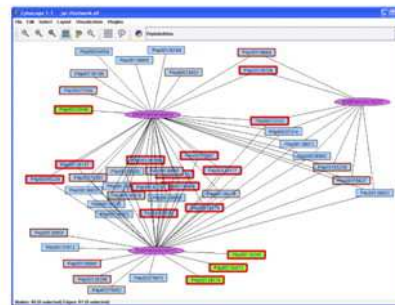
5. Quantification



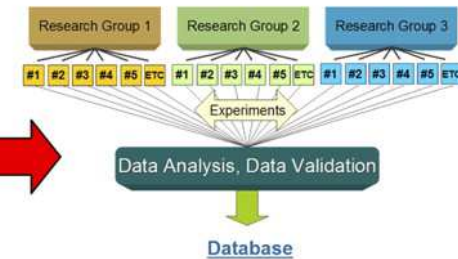
4. Protein inference



6. Organization in local data management systems

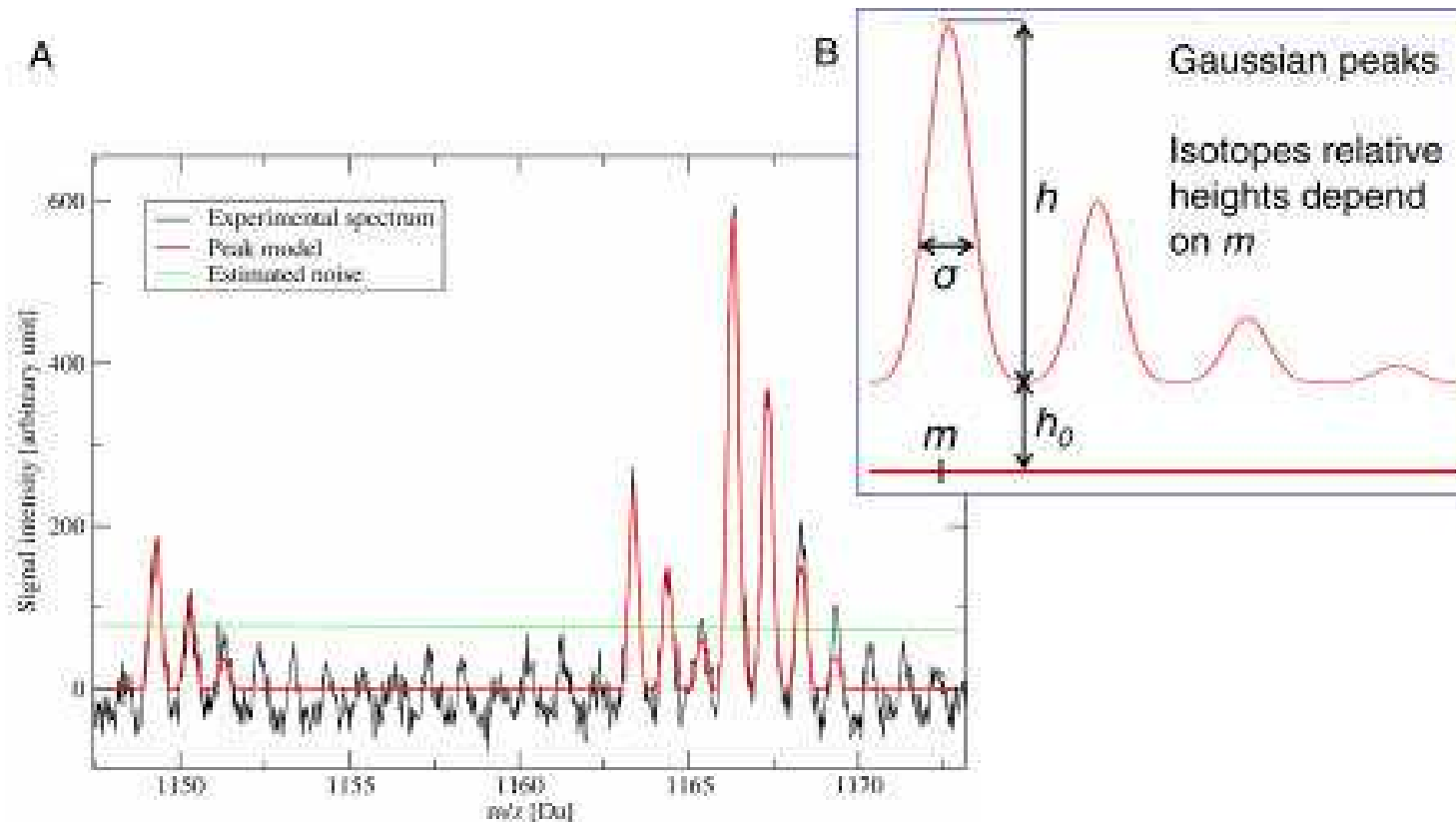


7. Interpretation of the protein lists



8. Transfer to public data repositories

Peak detection

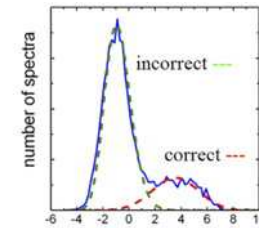
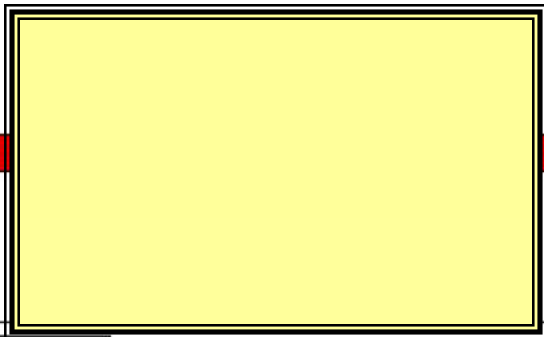


From Colinge & Bennett, 2007

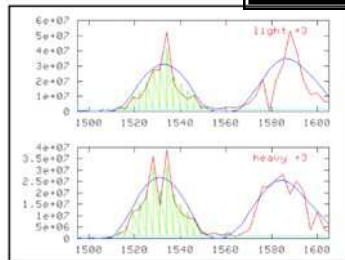
Schematic overview of a typical workflow of the proteomics informatics processing of a data set



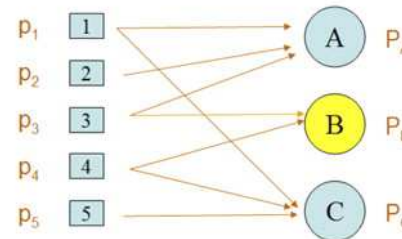
1. Conversion to and use of open data formats



3. Validation of identifications



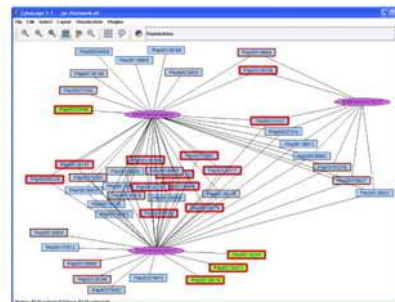
5. Quantification



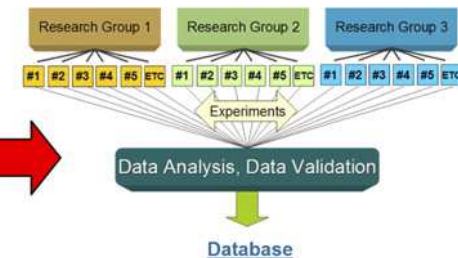
4. Protein inference



6. Organization in local data management systems



7. Interpretation of the protein lists



8. Transfer to public data repositories

The paradigms for protein identification using MS/MS data



- **Database searching (protein sequences)**
- **Interpretation of MS/MS spectra**
 - De novo sequencing
 - Peptide Sequence Tags (PSTs)
- **Using a reference database (analytical data)**
 - AMT database
 - Spectra libraries

The paradigms for protein identification using MS/MS data

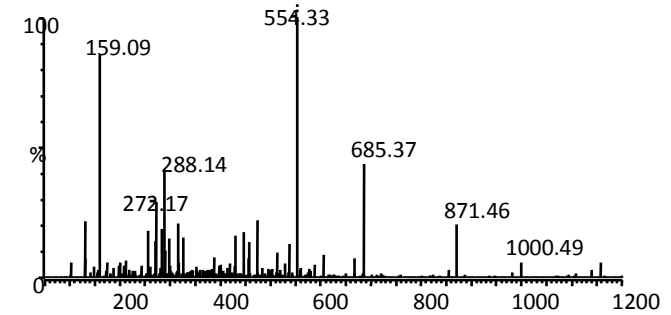


- **Database searching (protein sequences)**
- **Interpretation of MS/MS spectra**
 - De novo sequencing
 - Peptide Sequence Tags (PSTs)
- **Using a reference database (analytical data)**
 - AMT database
 - Spectra libraries

Searching uninterpreted MS/MS spectra

Experimental data :

- 1- M = mass of the peptide
- 2- MS/MS spectrum

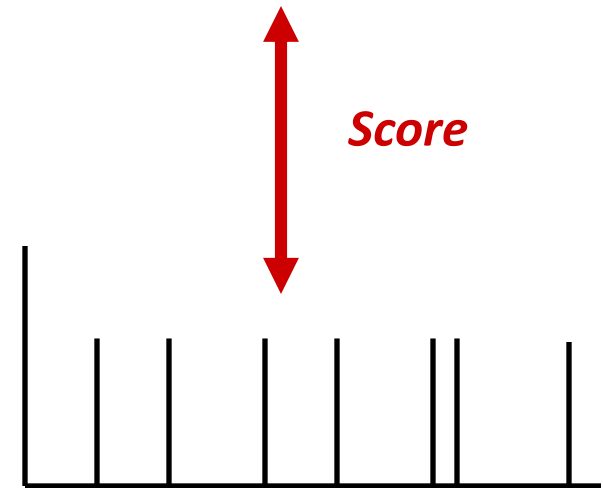


Mascot, Sequest, X! Tandem, OMSSA...

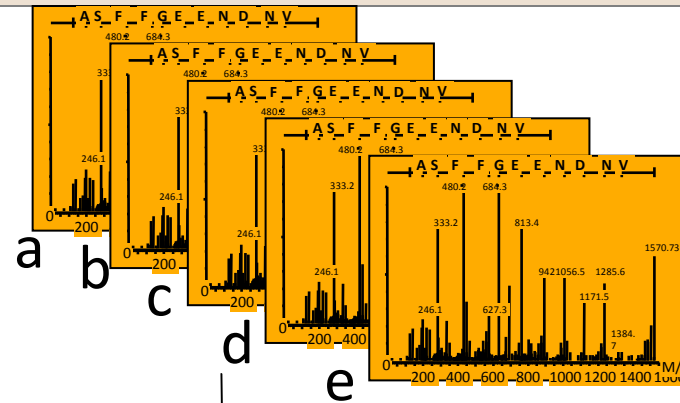
Virtual data:

- 1- trypsin digestion n M →
- 2- n MS/MS spectra

Protein databanks

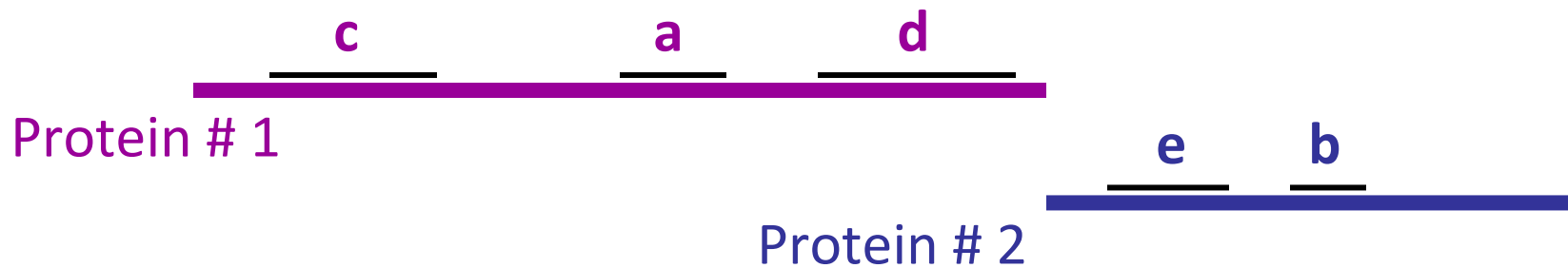


Database searching for shotgun analyses



*Database searching
(uninterpreted MS/MS spectra)*

Clustering



- **Open source tools**

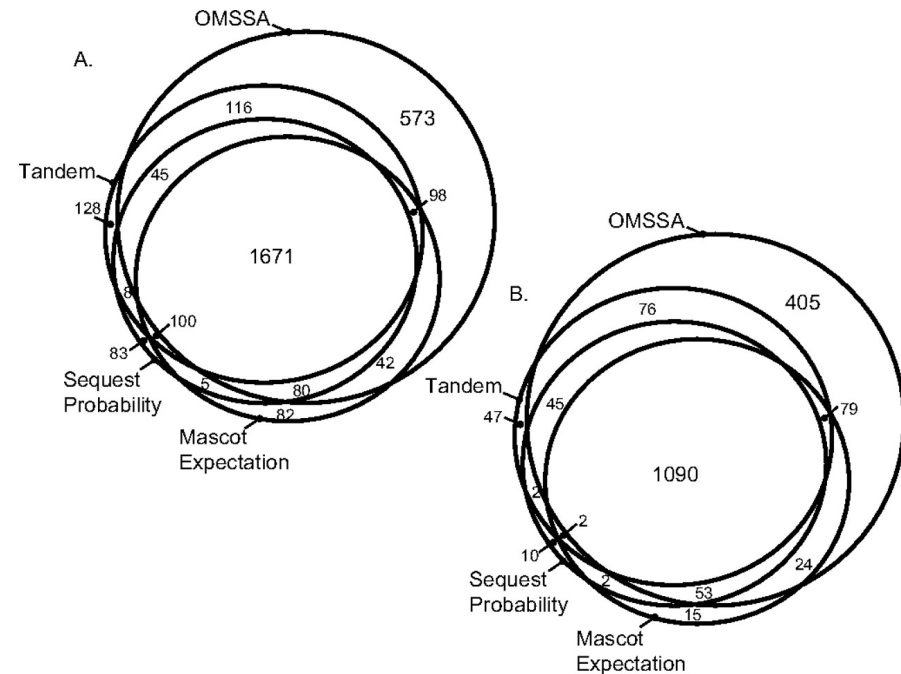
- X! Tandem
- OMSSA
- Comet ...

- **Commercial softwares**

- Mascot
- Phenyx ...

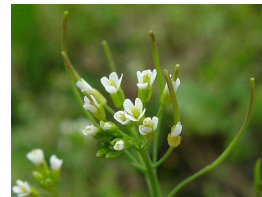
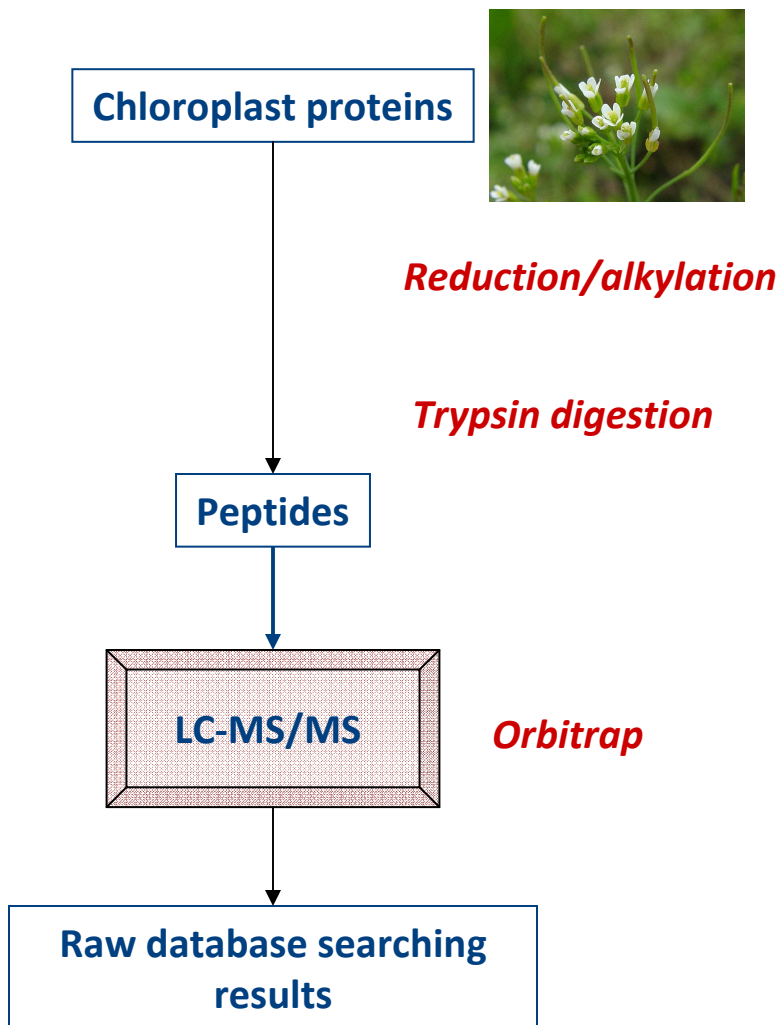
- **Available with Mass spectrometers**

- Sequest
- ProteinLynx Global Server
- ProteinPilot ...



→ A. Cornuéjols, Axe 1, 30/11

Search constraints



A. thaliana

Choices to be done

Protein databank (species)
(contaminants)

Chemical modifications
(e.g. alkylated methionine)

Enzyme specificity
Missed cleavage ?

Mass accuracy
For Orbitrap instruments
Peptide: 1-10 ppm
Fragments: 0.6 Da

Constraints in database searching

MASCOT MS/MS Ions Search

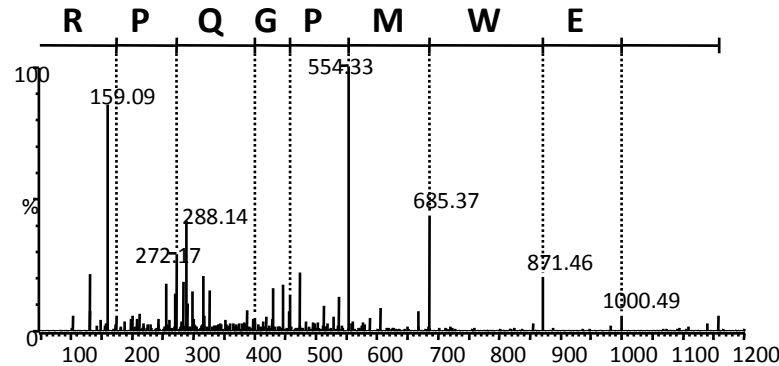
| | | | |
|---|--|---|---|
| Your name | <input type="text" value="Myriam"/> | Email | <input type="text" value="myriam.ferro@cea.fr"/> |
| Search title | <input type="text"/> | | |
| Database | ATH_Cplet_AMT <input type="button" value="v"/> | | |
| Taxonomy | All entries <input type="button" value="v"/> | | |
| Enzyme | Trypsin <input type="button" value="v"/> | Allow up to | 1 <input type="button" value="v"/> missed cleavages |
| Fixed modifications | <input type="text" value="Acetyl (K)"/> Acetyl (N-term) Acetyl (Protein N-term) Amidated (C-term) Amidated (Protein C-term) <input type="button" value="v"/> | Variable modifications | <input type="text" value="Acetyl (K)"/> Acetyl (N-term) <input type="button" value="v"/> Acetyl (Protein N-term) Amidated (C-term) Amidated (Protein C-term) <input type="button" value="v"/> |
| Quantitation | None <input type="button" value="v"/> | | |
| Peptide tol. ± | 20 <input type="button" value="ppm v"/> # ¹³ C 0 <input type="button" value="v"/> | MS/MS tol. ± | 1.2 <input type="button" value="Da v"/> |
| Peptide charge | 2+ <input type="button" value="v"/> | Monoisotopic | <input checked="" type="radio"/> Average <input type="radio"/> |
| Data file | <input type="text" value="D:\pkl_files\EUJU780b.mgf"/> <input type="button" value="Browse..."/> | | |
| Data format | Mascot generic <input type="button" value="v"/> | Precursor | <input type="text"/> m/z |
| Instrument | ESI-QUAD-TOF <input type="button" value="v"/> | Error tolerant | <input type="checkbox"/> |
| Decoy | <input type="checkbox"/> | Report top | AUTO <input type="button" value="v"/> hits |
| <input type="button" value="Start Search ..."/> | | <input type="button" value="Reset Form"/> | |

The paradigms for protein identification using MS/MS data



- Database searching (protein sequences)
- **Interpretation of MS/MS spectra**
 - De novo sequencing
 - Peptide Sequence Tags (PSTs)
- Using a reference database (analytical data)
 - AMT database
 - Spectra libraries

De novo sequencing



PepNovo; PEAKS etc.

MS/MS interpretation

EWMPGQPR → whole or partial sequence

MSBlast, BLASTP, TBLASTN ...

Alignment with protein, EST and genomic databases

EWMPGQPR

... .VGRRHGSRVSKSAEWMPGQPRPPHLDGSAPGD... .

Searching for similar proteins

EWMPGQPR

... .VGATNSMSRFSMSADW**MPGQPRPSYLDGSAPGD... .**

Dancik V, Addona TA, Clauser KR, Vath JE, Pevzner PA. *De novo peptide sequencing via tandem mass spectrometry*. J Comput Biol. 1999 Fall-Winter;6(3-4):327-42.

Shevchenko A, Chernushevich I, Shevchenko A, Wilm M, Mann M. "De novo" sequencing of peptides recovered from in-gel digested proteins by nanoelectrospray tandem mass spectrometry. Mol Biotechnol. 2002 Jan;20(1):107-18.

The paradigms for protein identification



- **Interpretation of MS/MS spectra**
 - De novo sequencing
 - Peptide Sequence Tags (PSTs)
- **Database searching (protein sequences)**
- **Spectral libraries**
 - Spectra libraries
 - Prerequisite: MS/MS spectra already “annotated”

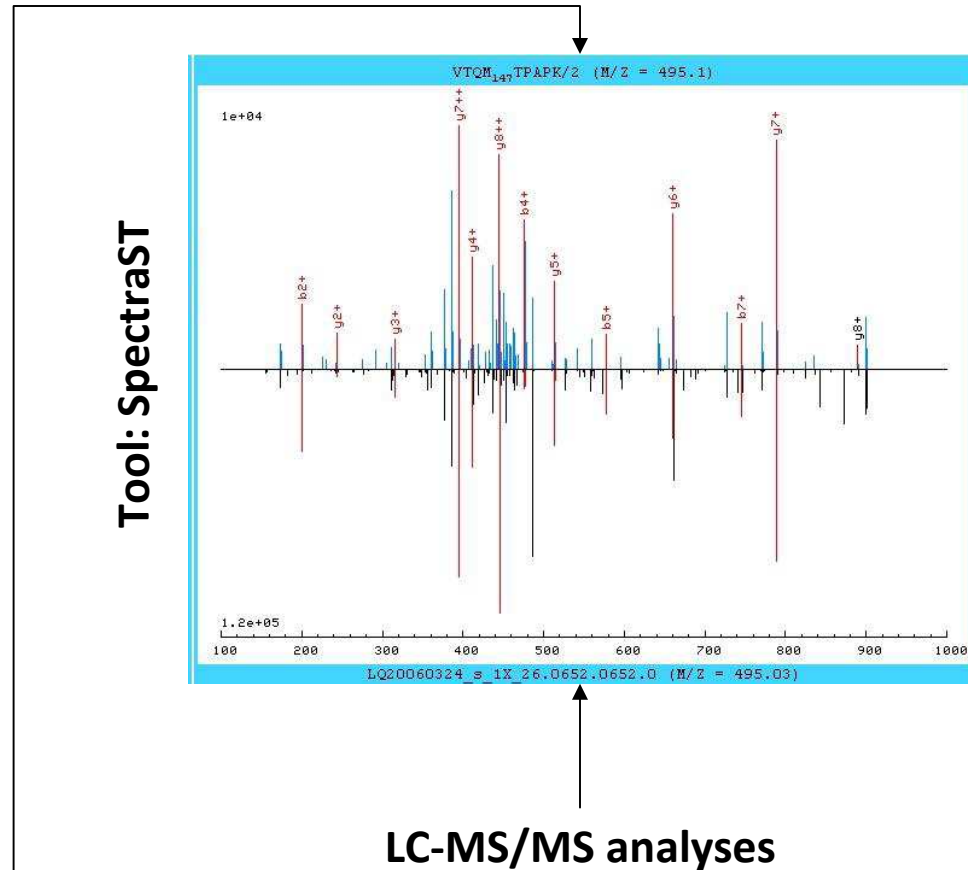
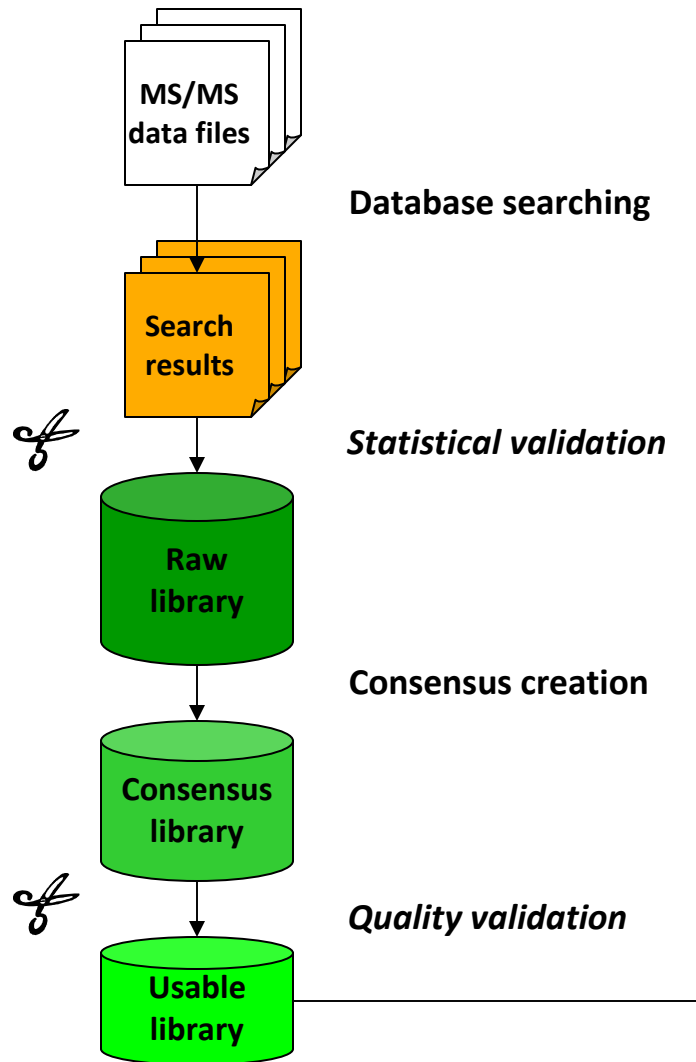
- **Concept**

- Capitalize on validated experimental data
 - Use peptide fragmentation « pattern »

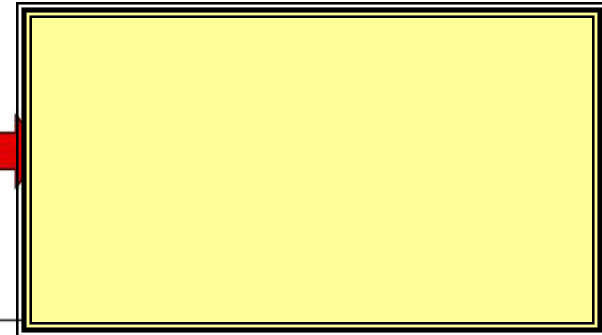
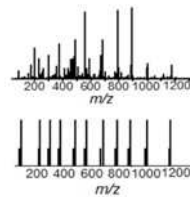
- **Principle**

- Create a « consensus » spectrum for all peptides identified from experimental spectra
- Compare this « consensus » spectrum with new experimental spectrum

Pipeline of a spectral library search

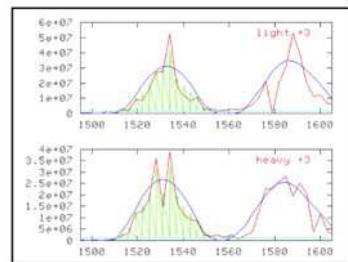


Schematic overview of a typical workflow of the proteomics informatics processing of a data set

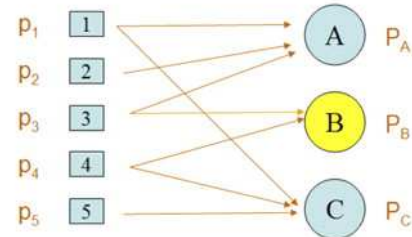


1. Conversion to and use of open data formats

2. Spectrum identification with a search engine



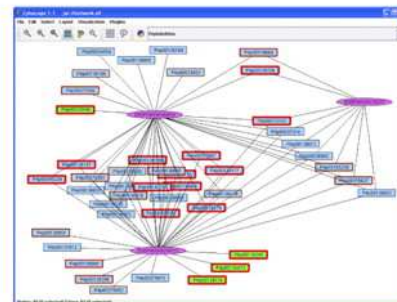
5. Quantification



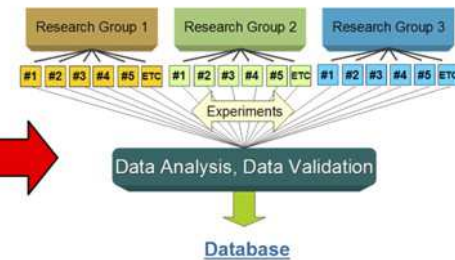
4. Protein inference



6. Organization in local data management systems



7. Interpretation of the protein lists



8. Transfer to public data repositories

Validation of database searching results



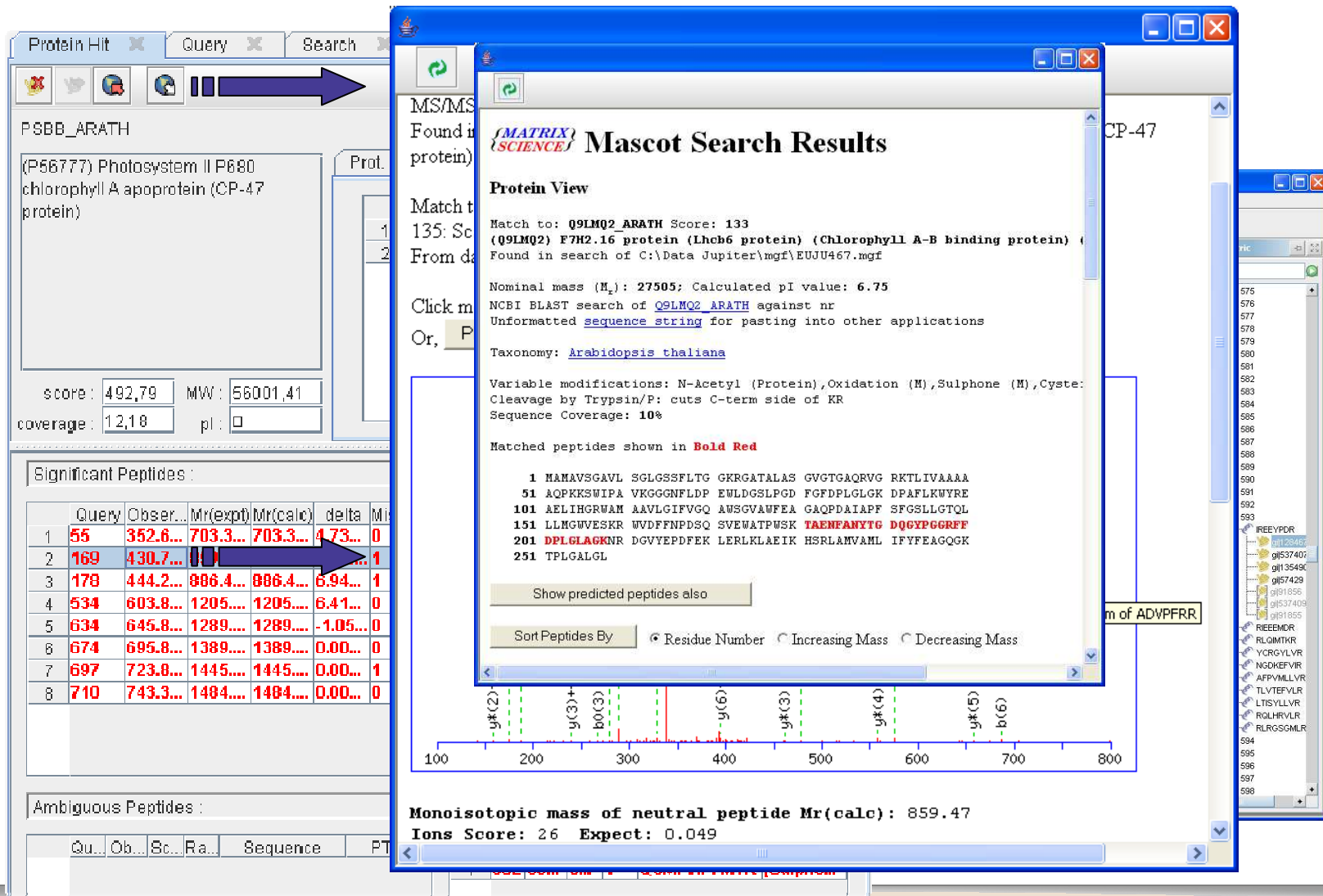
- **Tools**

- Many!

- **Needs**

- Results filtering
 - Using predefined rules
- Manual filtering
 - Validate spectrum-peptide match
 - View « matched » spectrum on peptide sequence
 - Access protein description (NCBI, SwissProt, etc)
- Refine a result by re-submitting spectra
- Save on-going validation
- Generate reports
 - Filtered ... but consistent

Validation of the results: an example using IRMa



The screenshot displays the IRMa software interface. On the left, a 'Protein Hit' window shows the search results for 'PSBB_ARATH', identifying it as '(P56777) Photosystem II P680 chlorophyll A apoprotein (CP-47 protein)'. Below this, a table of 'Significant Peptides' is shown, with a blue arrow pointing to the second row (peptide 169).

| | Query | Observed | Mr(expt) | Mr(calc) | delta | Mod |
|---|-------|----------|----------|----------|----------|-----|
| 1 | 55 | 352.6... | 703.3... | 703.3... | 4.73... | 0 |
| 2 | 169 | 430.7... | 859.4... | 859.4... | 0.00... | 1 |
| 3 | 178 | 444.2... | 886.4... | 886.4... | 6.94... | 1 |
| 4 | 534 | 603.8... | 1205.... | 1205.... | 6.41... | 0 |
| 5 | 634 | 645.8... | 1289.... | 1289.... | -1.05... | 0 |
| 6 | 674 | 695.8... | 1389.... | 1389.... | 0.00... | 0 |
| 7 | 697 | 723.8... | 1445.... | 1445.... | 0.00... | 1 |
| 8 | 710 | 743.3... | 1484.... | 1484.... | 0.00... | 0 |

The main window, 'Mascot Search Results', displays the 'Protein View' for the match: '09LMQ2_ARATH Score: 133'. It includes the protein name '(09LMQ2) F7H2.16 protein (Lhcb6 protein) (Chlorophyll A-B binding protein)', nominal mass (27505), calculated pI (6.75), and taxonomy ('Arabidopsis thaliana'). A list of 'Matched peptides shown in Bold Red' is provided, with peptide 201 'DPLGLAGKNR' highlighted. Below the list is a mass spectrum plot showing peaks for y*(2) through b(6). At the bottom, the 'Monoisotopic mass of neutral peptide Mr(calc): 859.47' and 'Ions Score: 26 Expect: 0.049' are displayed.

False discovery rate (FDR)

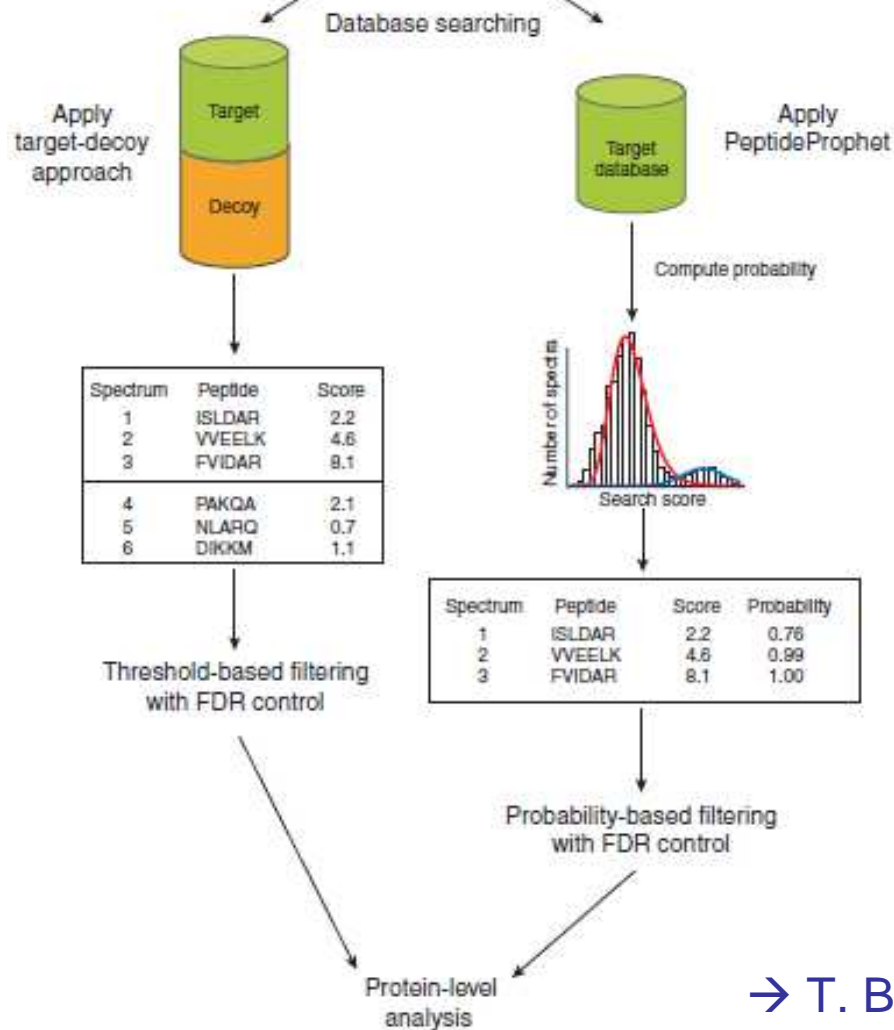
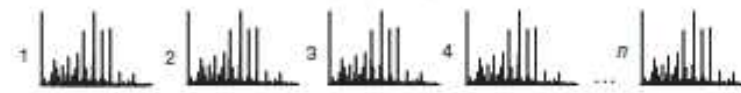
- **False discovery rate**

- Percentage of incorrect peptide-spectra matches in identification result

- **Problem**

- How to distinguish « true » and « false » PSMs

FDR



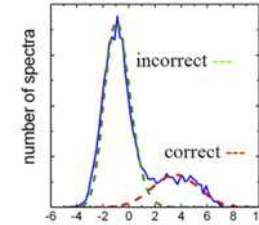
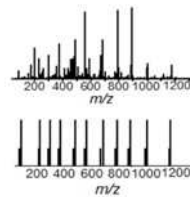
$$FP = \frac{2 \times Rev}{Rev + Fwd}$$

→ T. Burger, Axe1, 29/11

Schematic overview of a typical workflow of the proteomics informatics processing of a data set



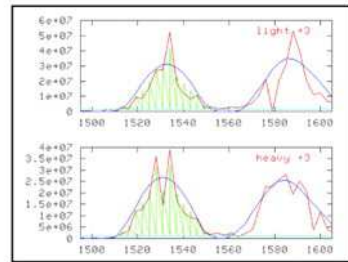
Common Data Format



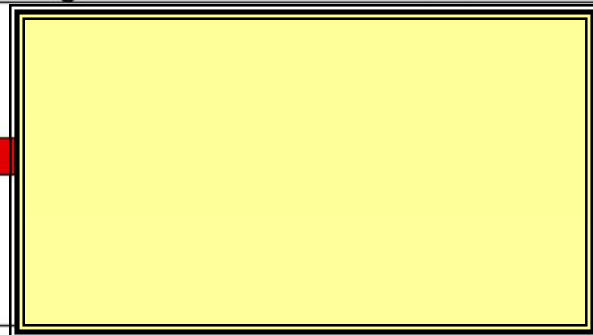
1. Conversion to and use of open data formats

2. Spectrum identification with a search engine

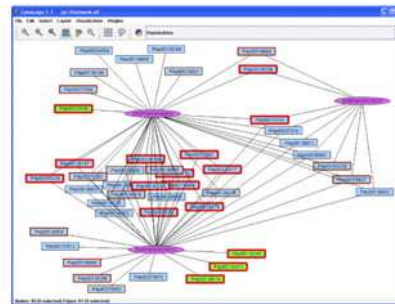
3. Validation of identifications



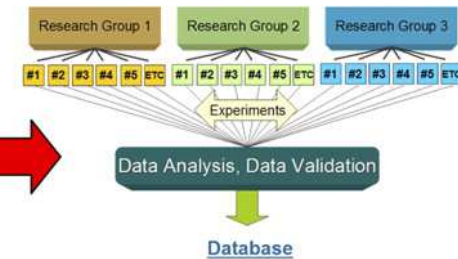
5. Quantification



6. Organization in local data management systems



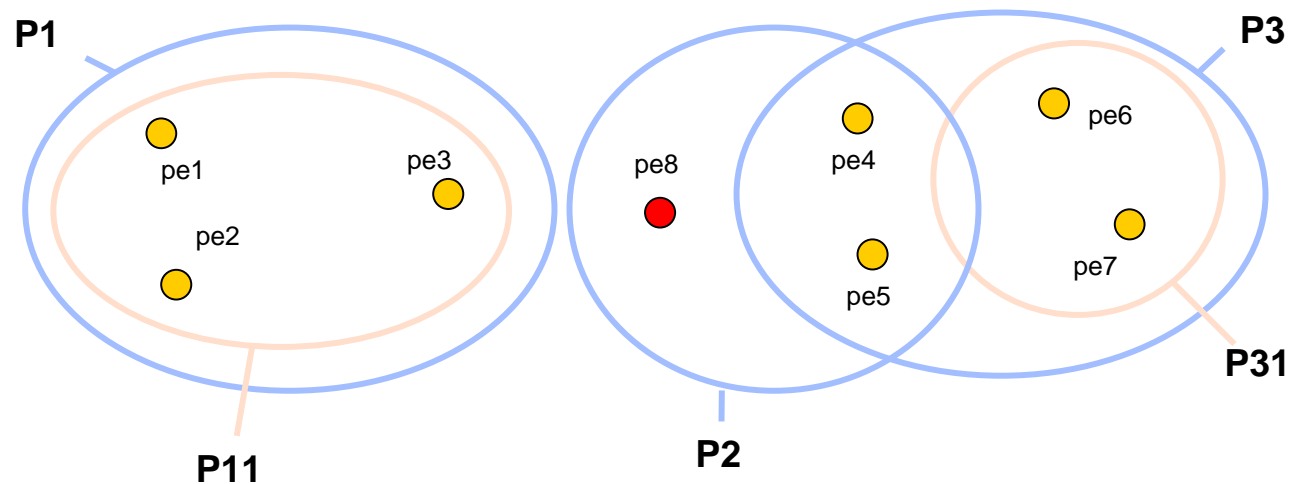
7. Interpretation of the protein lists



8. Transfer to public data repositories

The protein inference problem

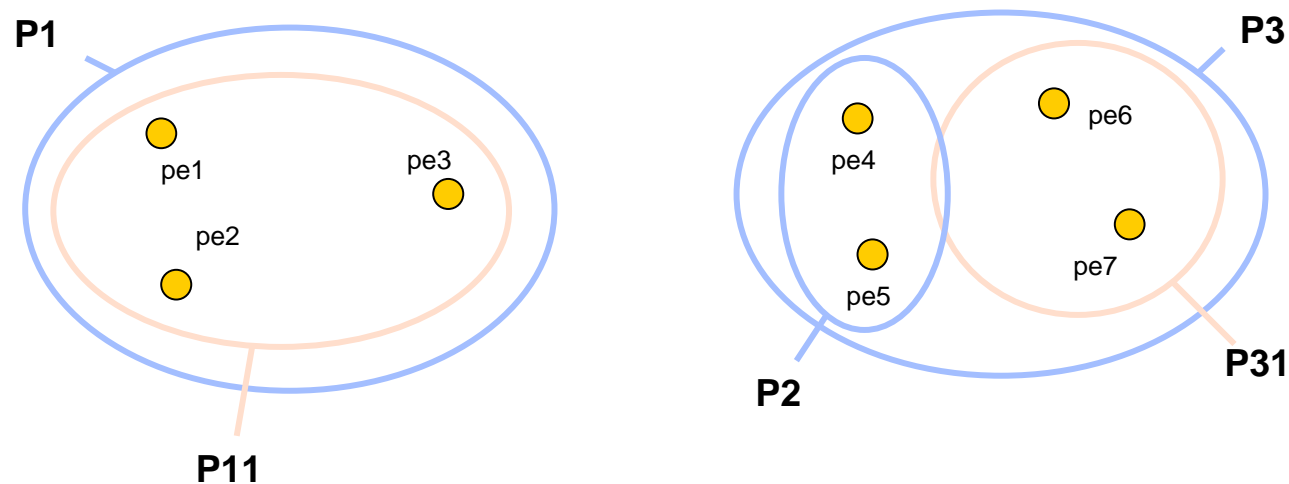
- **The result is a list of protein groups and not proteins**
 - Proteins that share a same set of peptides (P1 and P11)
 - Proteins that share a subset of peptides (P3 and P31)



The protein inference problem

- **List of protein groups and not proteins**

- Proteins that share a same set of peptides (P1 and P11)
- Proteins that share a subset of peptides (P3 and P31)
- After validation the list of protein groups can be modified (consistency to be assured)

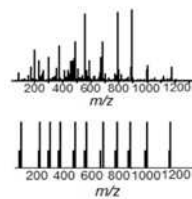


—————> B. Valot, S. Bouveret, axe 1, 29/11

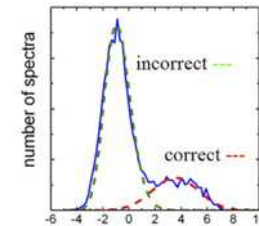
Schematic overview of a typical workflow of the proteomics informatics processing of a data set



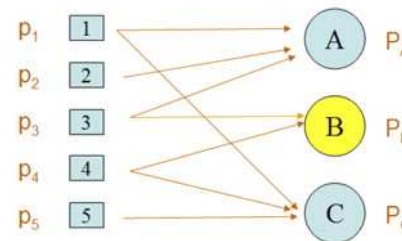
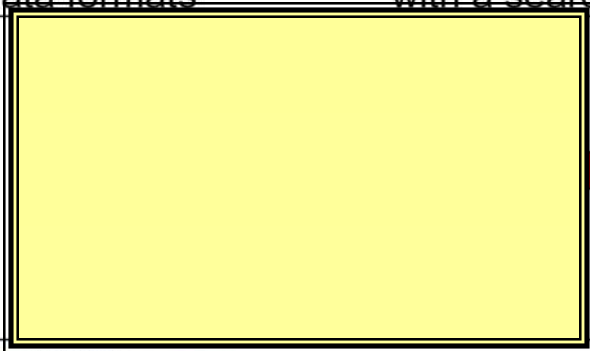
1. Conversion to and use of open data formats



2. Spectrum identification with a search engine



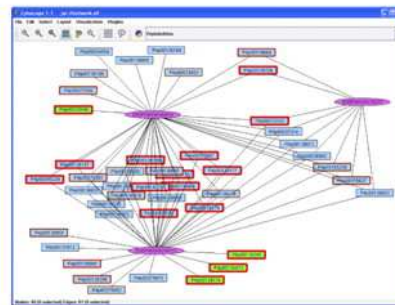
3. Validation of identifications



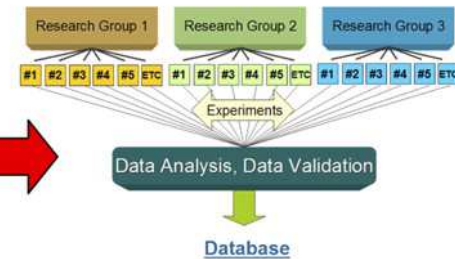
4. Protein inference



6. Organization in local data management systems



7. Interpretation of the protein lists



8. Transfer to public data repositories

Quantitative proteomics



-Comparison of peptide/protein levels in 2, 3 ... n samples

→ **Relative quantification**

(up- or down-regulation of a protein in a sample relative to an other, results expressed as a 'fold' increase or decrease)

Determination of an exact amount (concentration) of peptide/protein in a sample

→ **Absolute quantification**

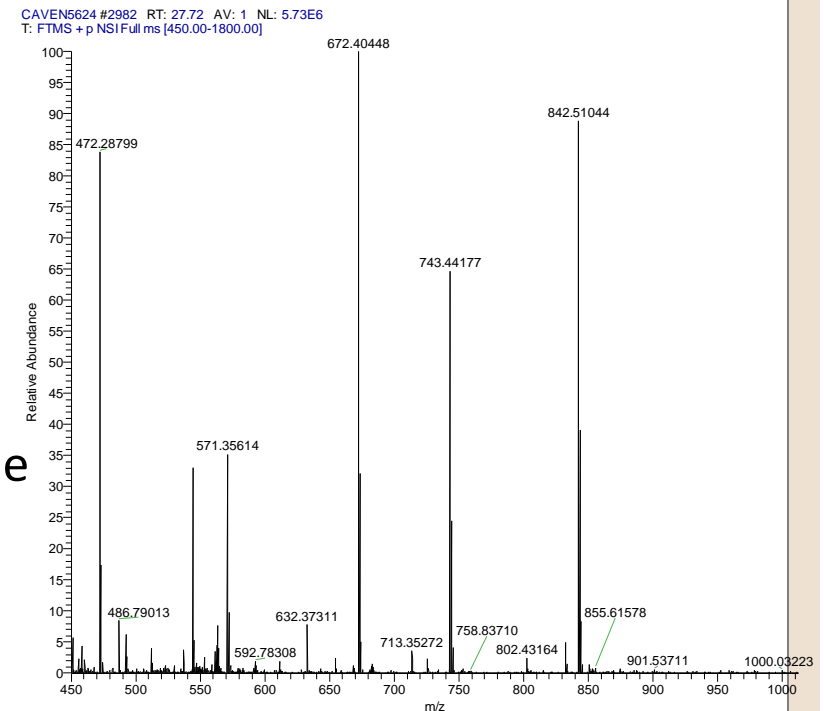
(e.g. ng or nmoles per gram of tissue, or ng or nmoles/ml of plasma ; use of an internal standard)

Even **Absolute quantification is Relative** – relative to an internal standard

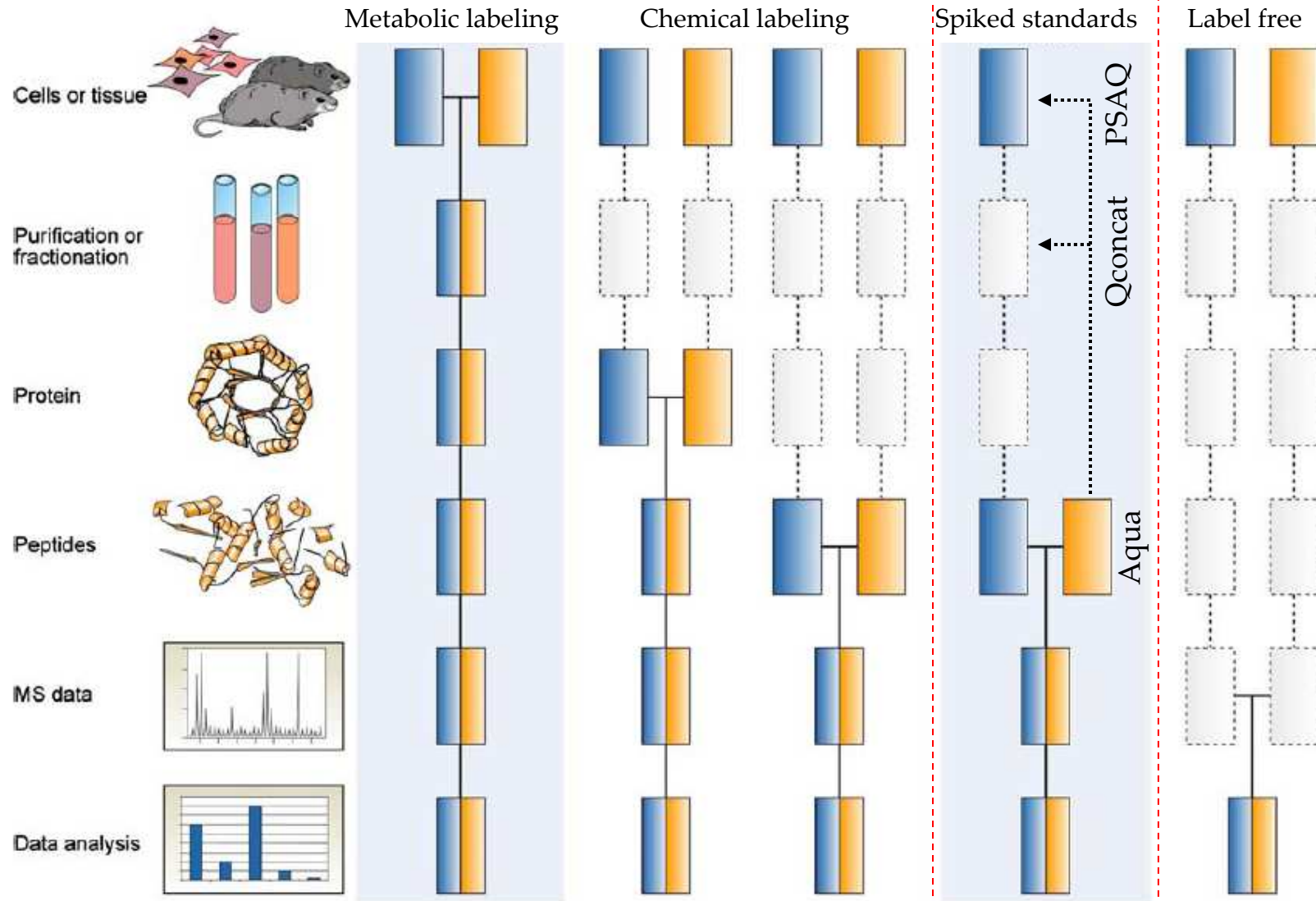
MS-based quantitative proteomics

!!! For several reasons, MS is not a quantitative tool !!!

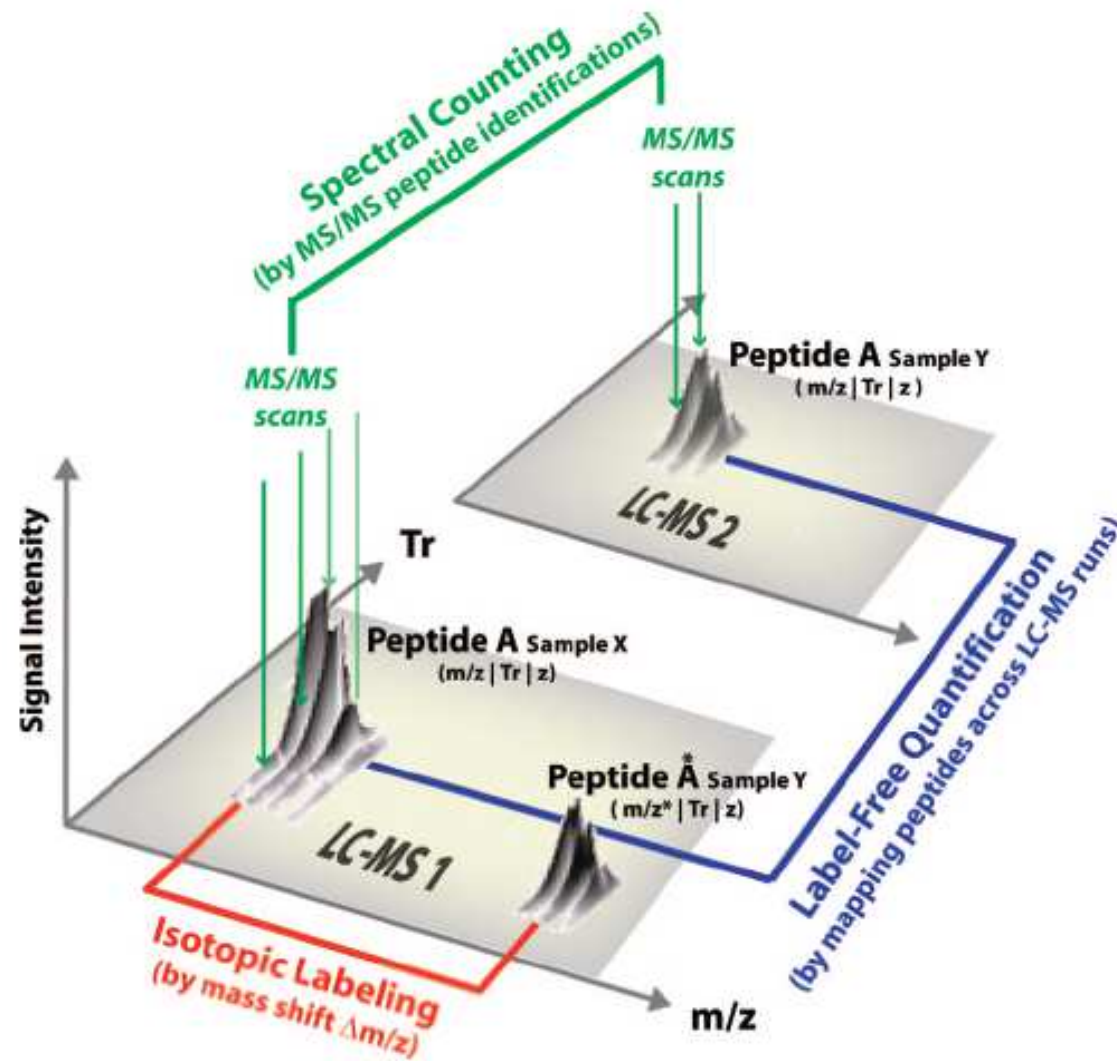
- Ionization is function of the structure
 - Physico-chemical properties of amino acids present in peptide
- Ionization suppression effect
 - Competition between peptides
 - Presence of one peptide may suppress the ionization of another one
- Ionization depends of the sample matrix (salt, detergents, contaminants, ...)
 - Variation in samples will result in variation in signal intensities



Comparative proteomics workflows



Relative quantification



Spectral counts

- **Implementations:**

- Number of matched MS/MS spectra for a given protein

- **Normalization factors**

- Protein length
- Total number of spectra
- etc.

- **Derived metrics used to determine sample relative composition**

$$PAI = \frac{N_{Obsd}}{N_{Obsvbl}} \quad emPAI = 10^{PAI} - 1$$

$$\text{Protein } i \text{ molar content} = \frac{emPAI_i}{\sum emPAI}$$

$$APEX_i = \frac{n_i \times p_i}{O_i \times \sum_{k=1}^N \frac{n_k \times p_k}{O_k}} \times C$$

Spectral counts : a semi-quantitative approach



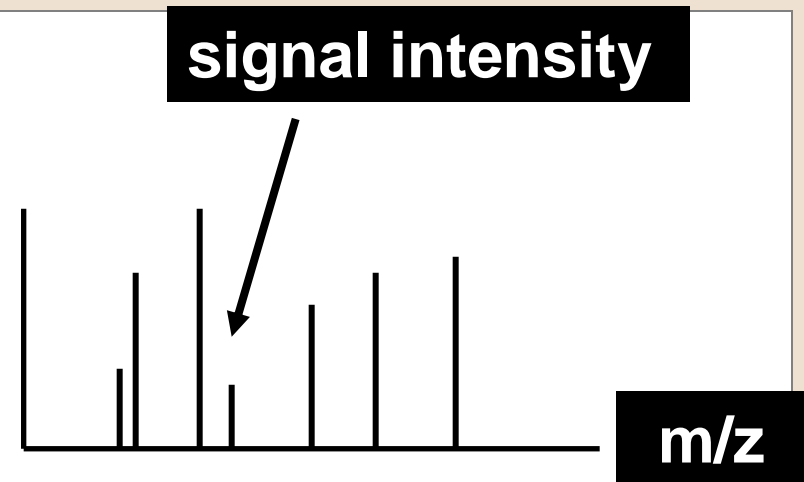
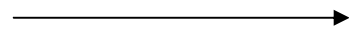
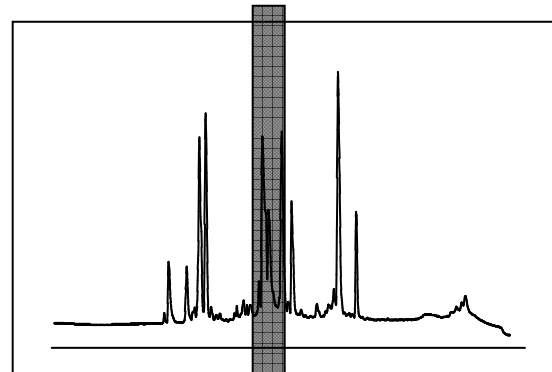
- **Advantages**

- Easy to implement.
- Do not require observation of the same peptides between experiments.

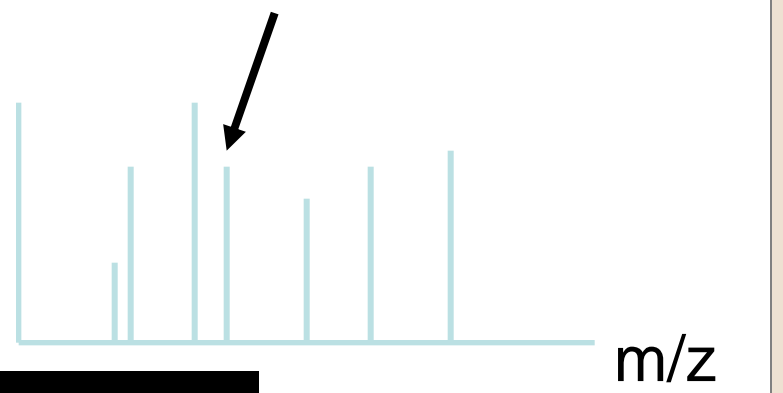
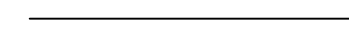
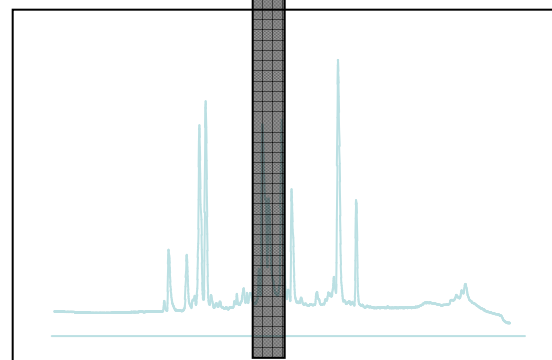
- **Drawbacks**

- Requires many spectra per protein (MudPIT based).
- Saturation at high SC
- Complicated in case of shared peptides (isoforms)

Label-free quantification



Mass accuracy, calibration, resolution

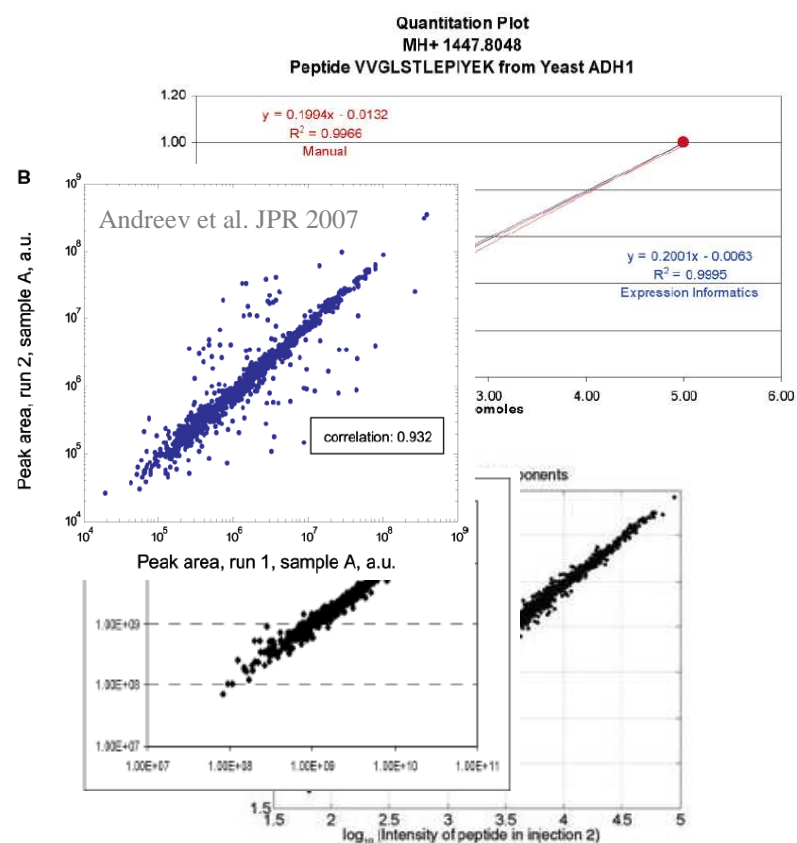


Retention time

LC reproducibility

EXtracted Ion Chromatogram (XIC) alias « MS trace »

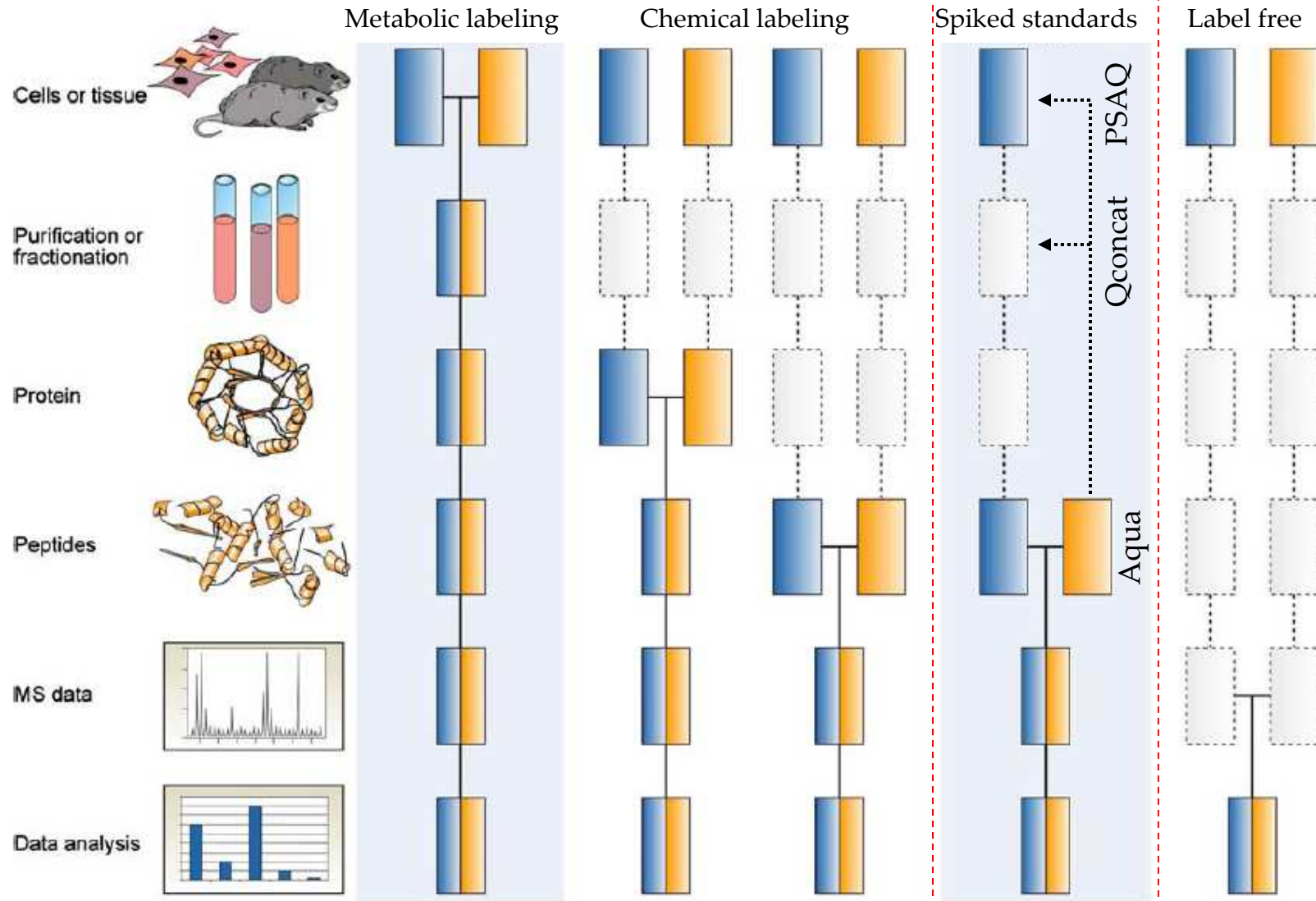
- Peptides abundances scale linearly with concentration
- Controlled sample handling and analytical procedures allow repeatable peptide abundance measurements



XIC: challenges

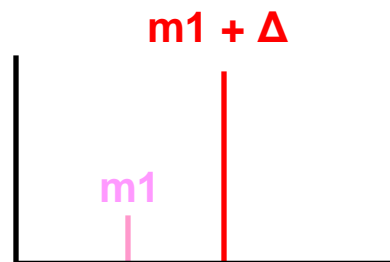
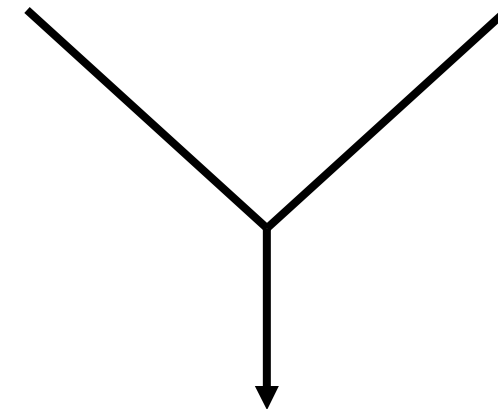
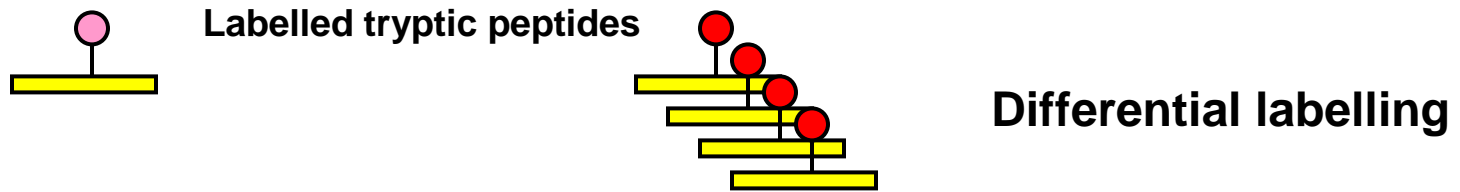
- **Quantitative readings must be extracted from MS or MS/MS spectra → intensities to be extracted**
- Peptide and protein identification must be performed
- The two types of information must be merged and quality controlled
- **Applicable statistical methods have to be identified**
- Individual steps have to be combined in an automated workflow bridging the gaps between commercially available software and custom-built tools.

Labeling approaches



Absolute

Labeling approaches : general principle



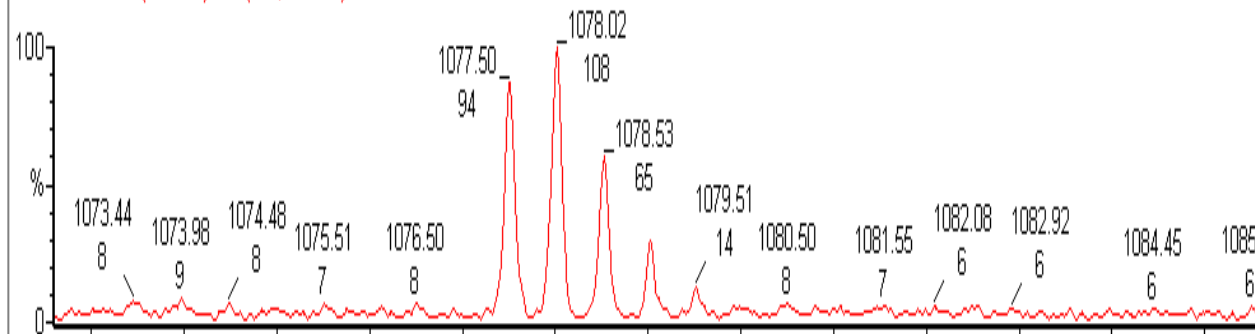
LC-MS/MS



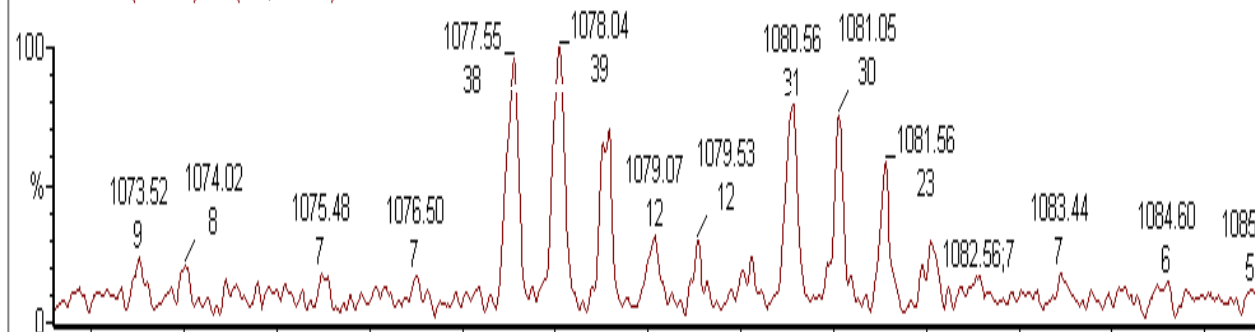
Quantification
(pairs detection)

SILAC: example

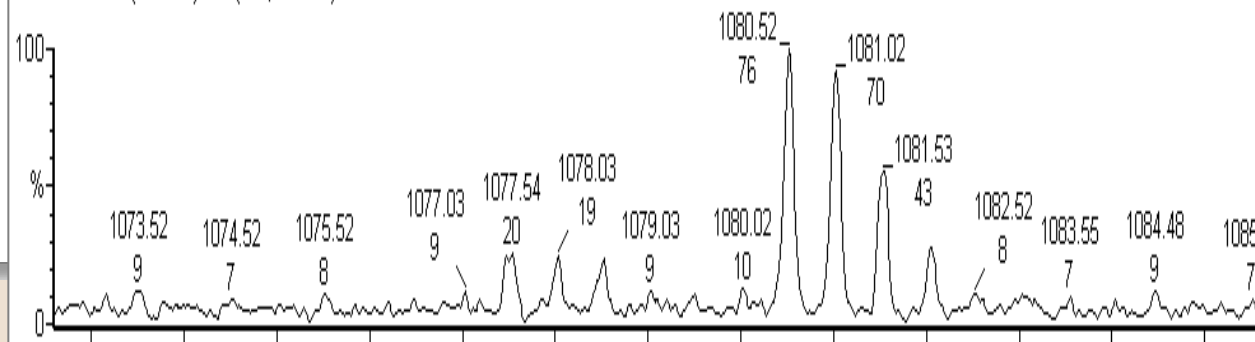
LCU5042 1069 (25.440) Sm (SG, 2x3.00)



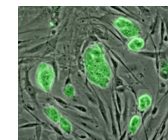
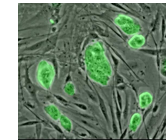
LCU5038 980 (23.748) Sm (SG, 2x3.00)



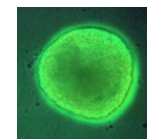
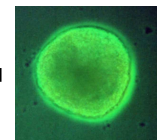
LCU5022 881 (21.865) Sm (SG, 2x3.00)



ES (^{12}C -Arg)

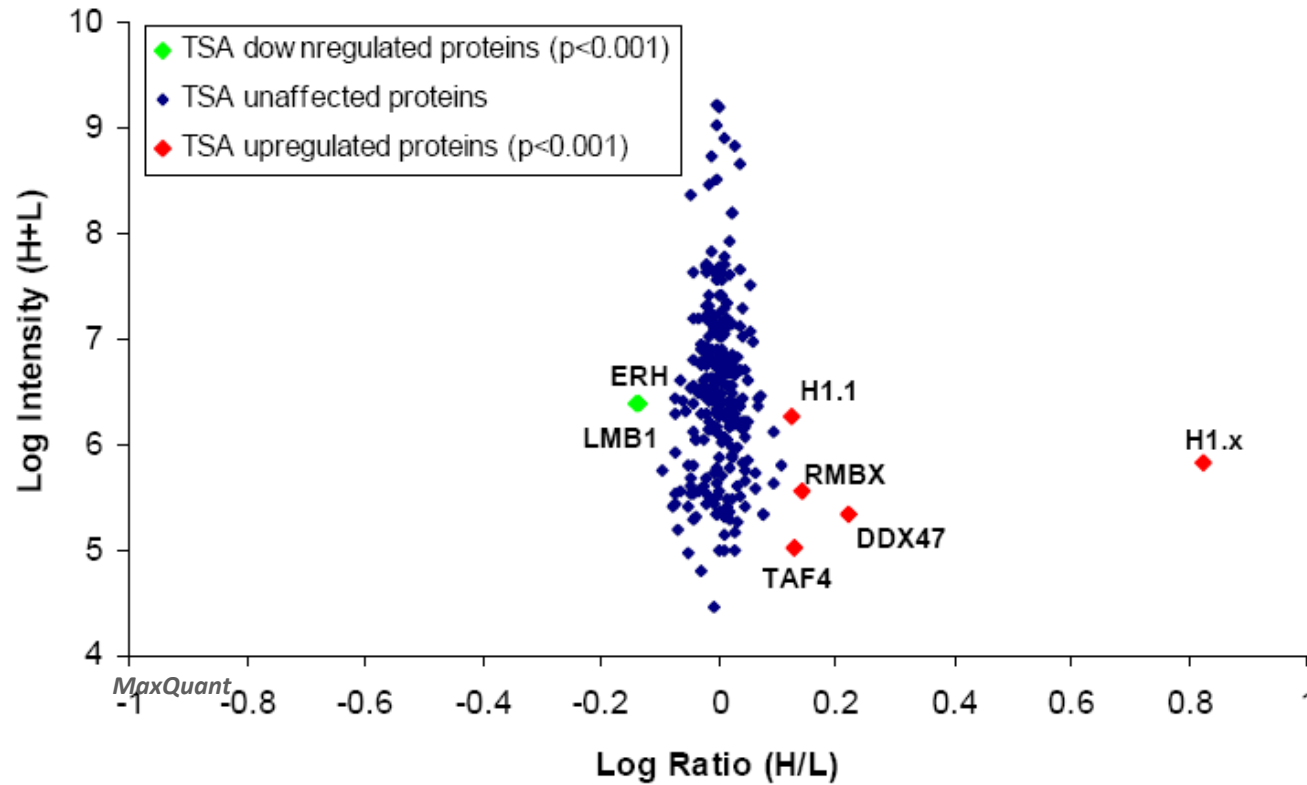


+



ESd (^{13}C -Arg)

SILAC: application exemple



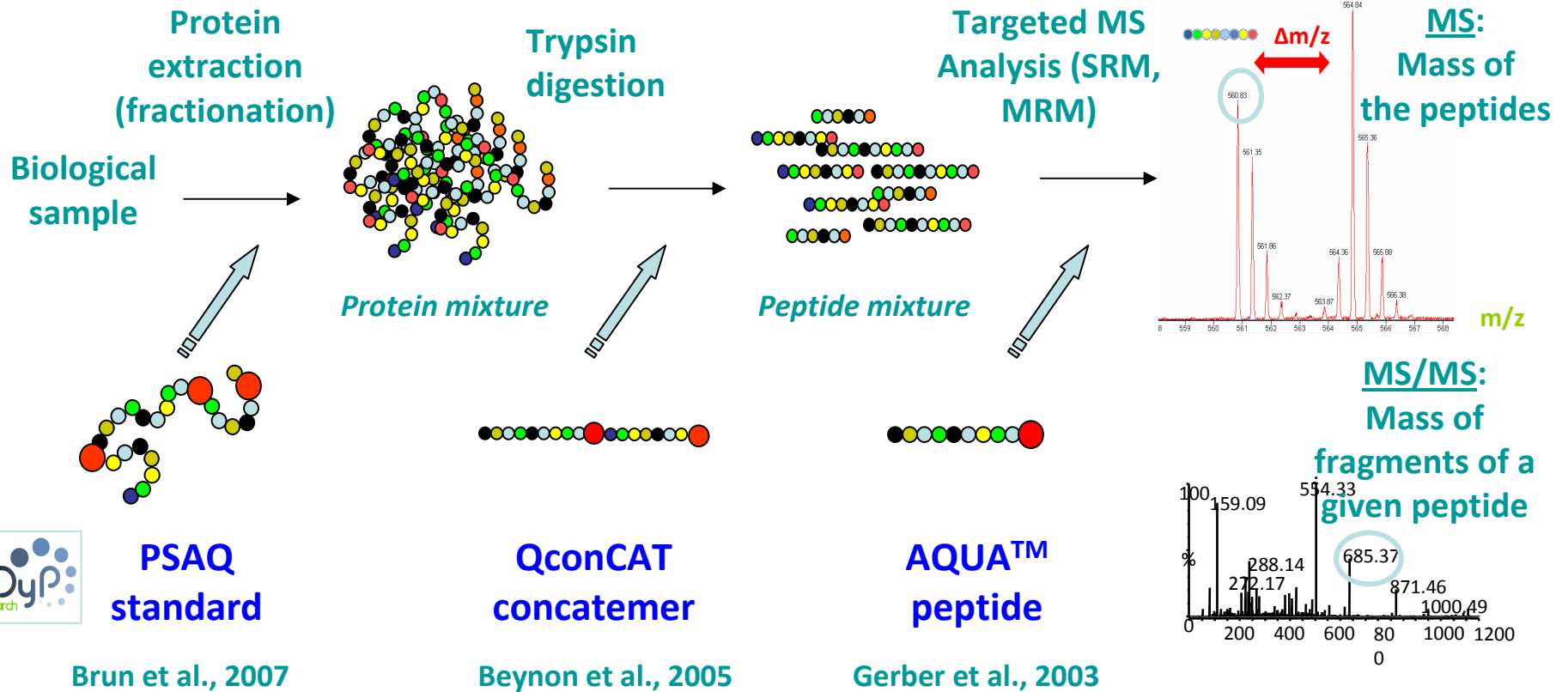
Biological duplicates + Cross-labelling

Identification : 1% False positive – **Quantification** : At least 2 peptides / protein

Methods with standards

- Towards « absolute » (accurate) quantification :
measuring the amount of a given protein in a
complex mixture
- The proteins to be quantified are already
identified (ex: validation de biomarkers)

Absolute quantification of proteins: use of labelled standards (isotope dilution)

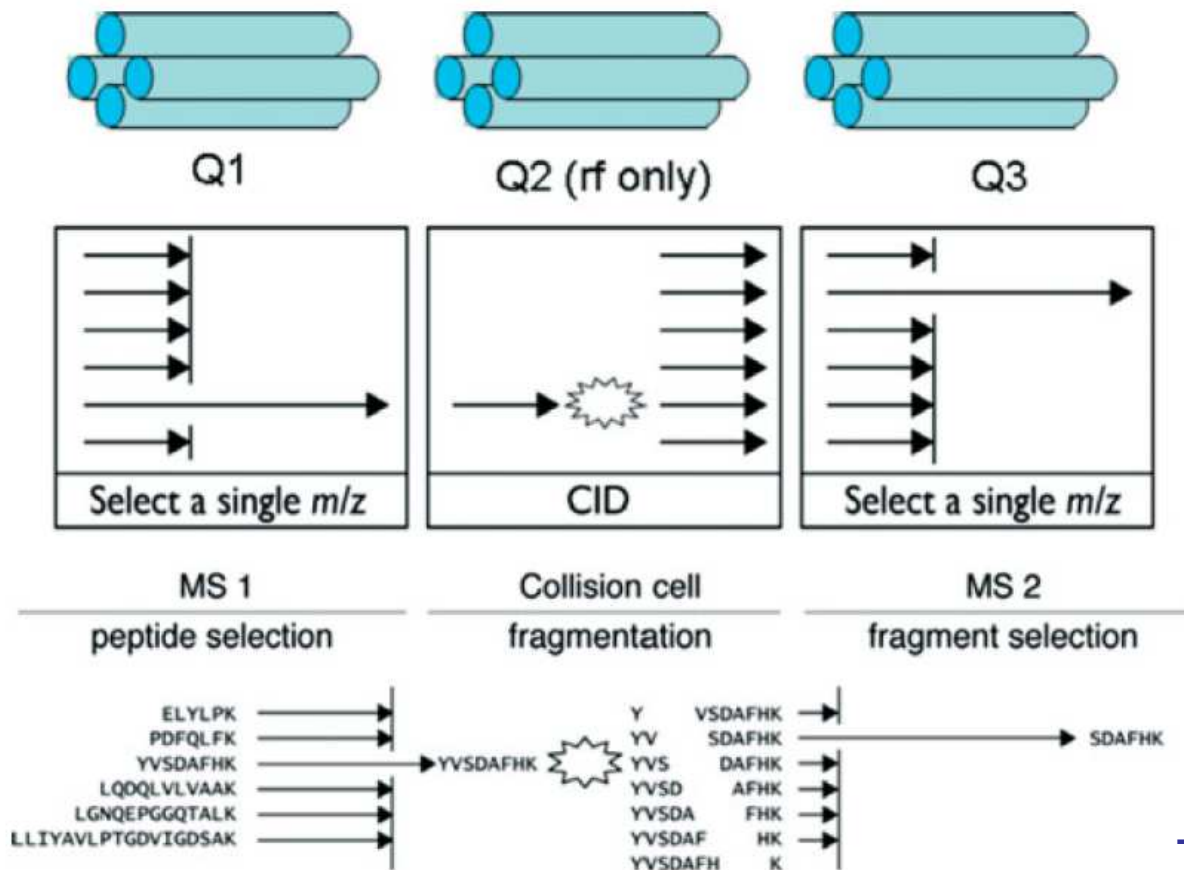


→ Specificity
→ Sensitivity

Accuracy

SRM

Monitoring of selected parent and specific fragment ions (enhanced accuracy and specificity).



Selection of proteotypic peptides and of good transitions.

Less automated process than discovery proteomics.

→ JF Giovanelli, Axe 1, 30/11

Quantitative proteomics: a question of choice



Label-free, labeling
Absolute, relative
Limited replicates

Orbitrap
Qtrap
Etc.

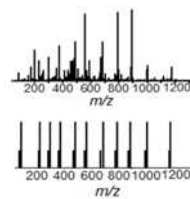
T-tests
Logistic regression
ML
Etc. → A. Klich, Axe1, 29/11
→ L. Gerfault, Axe1, 30/11

→ Not necessarily the optimal choice but the most pragmatic one

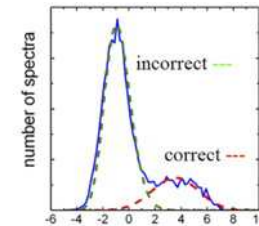
Schematic overview of a typical workflow of the proteomics informatics processing of a data set



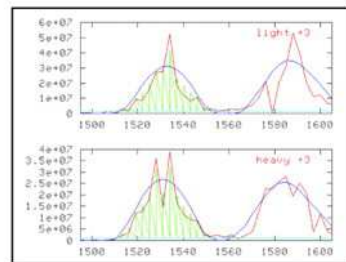
1. Conversion to and use of open data formats



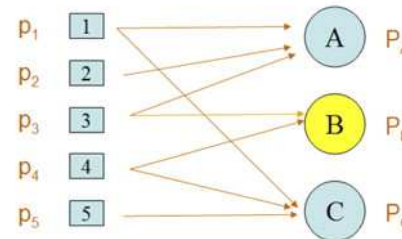
2. Spectrum identification with a search engine



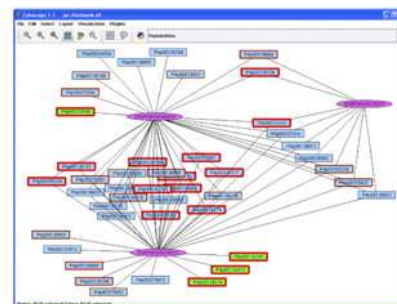
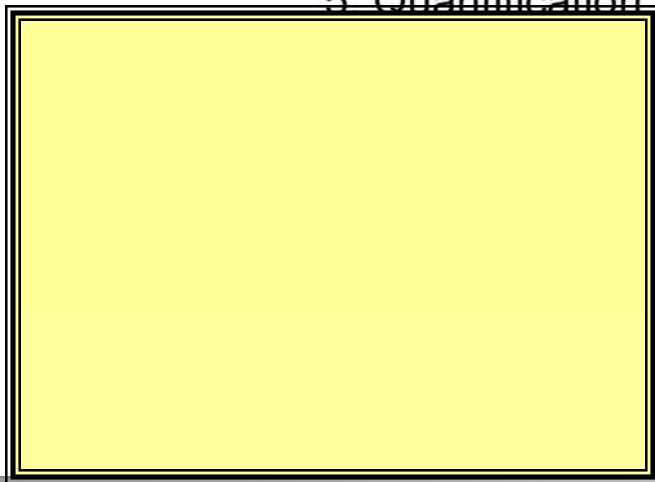
3. Validation of identifications



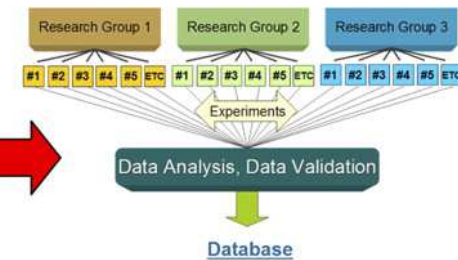
5. Quantification



4. Protein inference



7. Interpretation of the protein lists

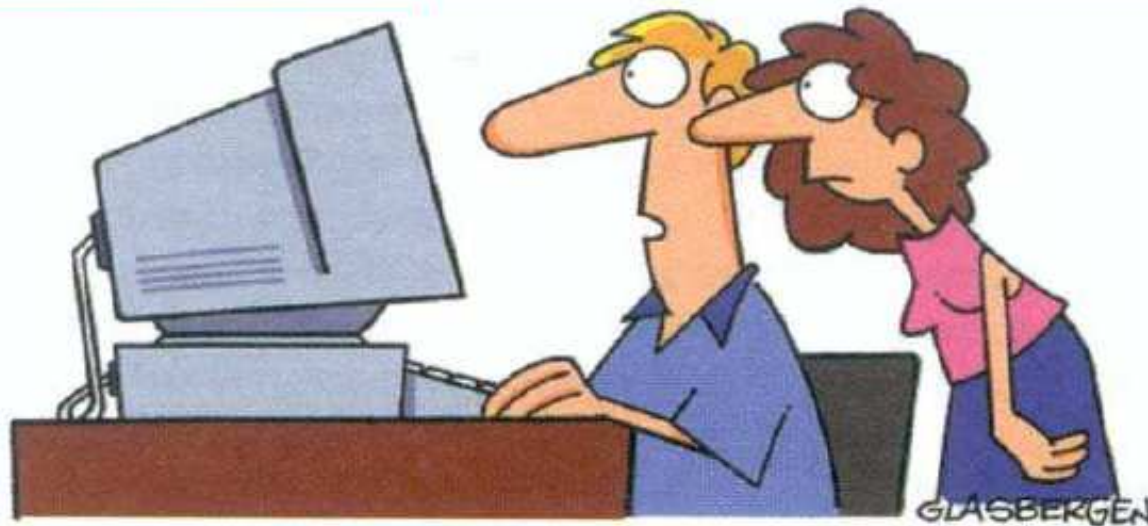


8. Transfer to public data repositories

Need for computing assistance !

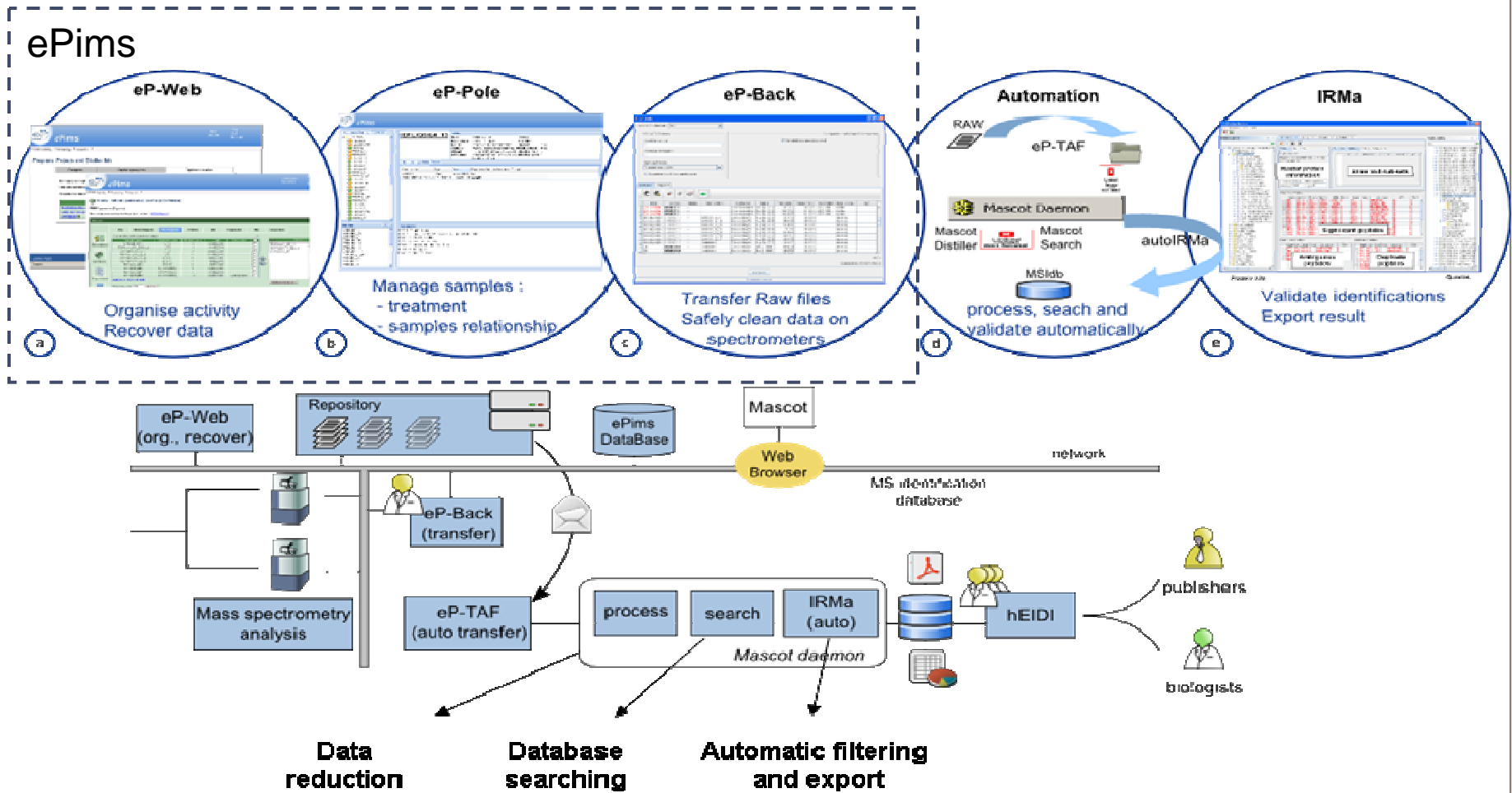


600 proteins/h/ instrument
Up to 20 000 MS-MS spectra/day/instrument
Up to 50 Go de données/ day



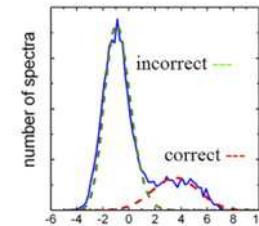
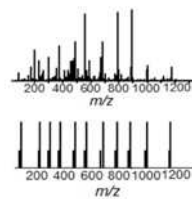
« The computer says I need to upgrade my brain
to be compatible with proteomic data analysis »

An IT workflow



ePIMS™: a dedicated LIMS for proteomics, open source (EDyP)
IRMa: automatic validation of Mascot results. Dupierris et al. 2009, Bioinformatics
MSldb: relational database for the storage of proteomics data (identification, quantification)

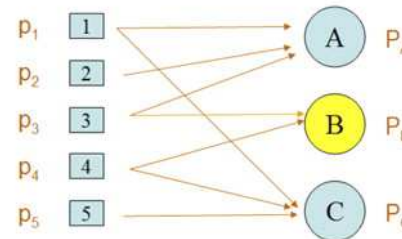
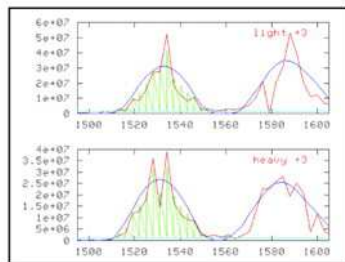
Schematic overview of a typical workflow of the proteomics informatics processing of a data set



1. Conversion to and use of open data formats

2. Spectrum identification with a search engine

3. Validation of identifications

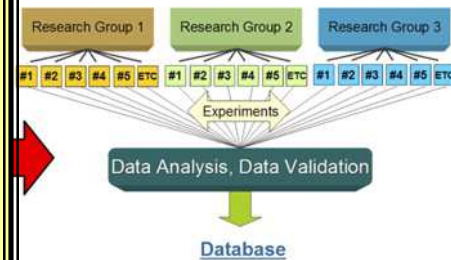
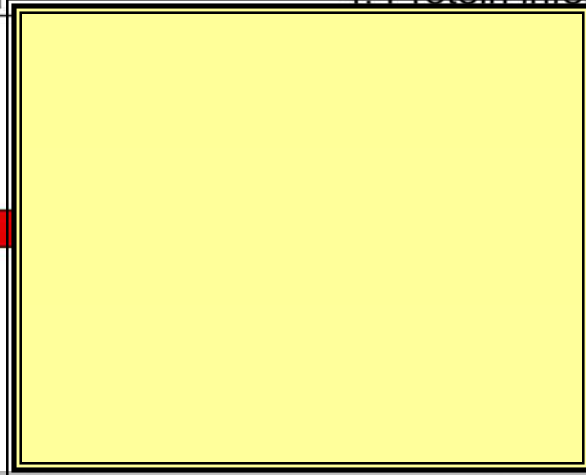


5. Quantification

4. Protein inference



6. Organization in local data management systems



8. Transfer to public data repositories

Looking for additional information

Technology-dependent

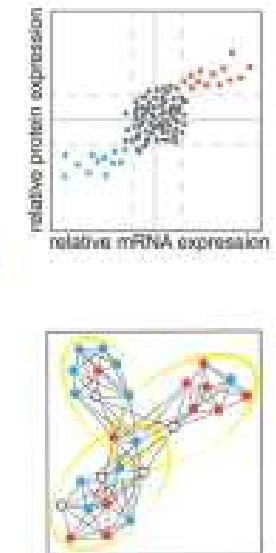
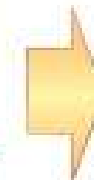
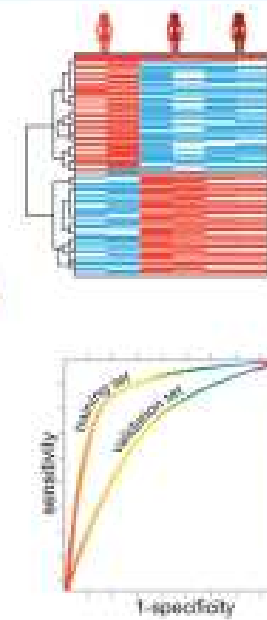
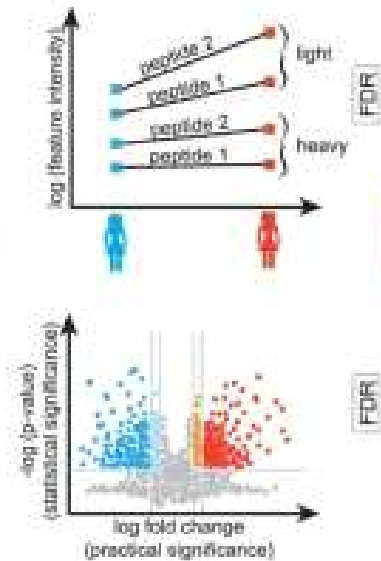
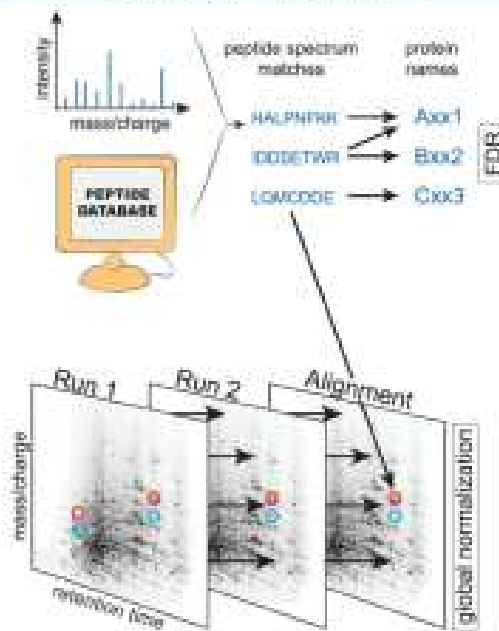
Technology-independent

a) peptide and protein identification from PSMs

c) peptide significance analysis

e) class discovery

g) data integration



b) feature detection, quantification, annotation, and alignment

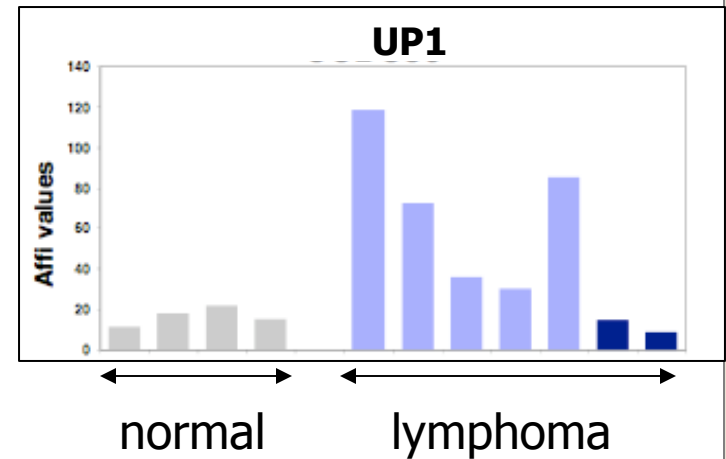
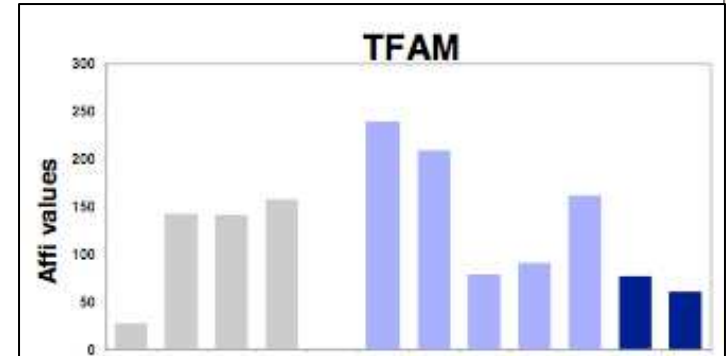
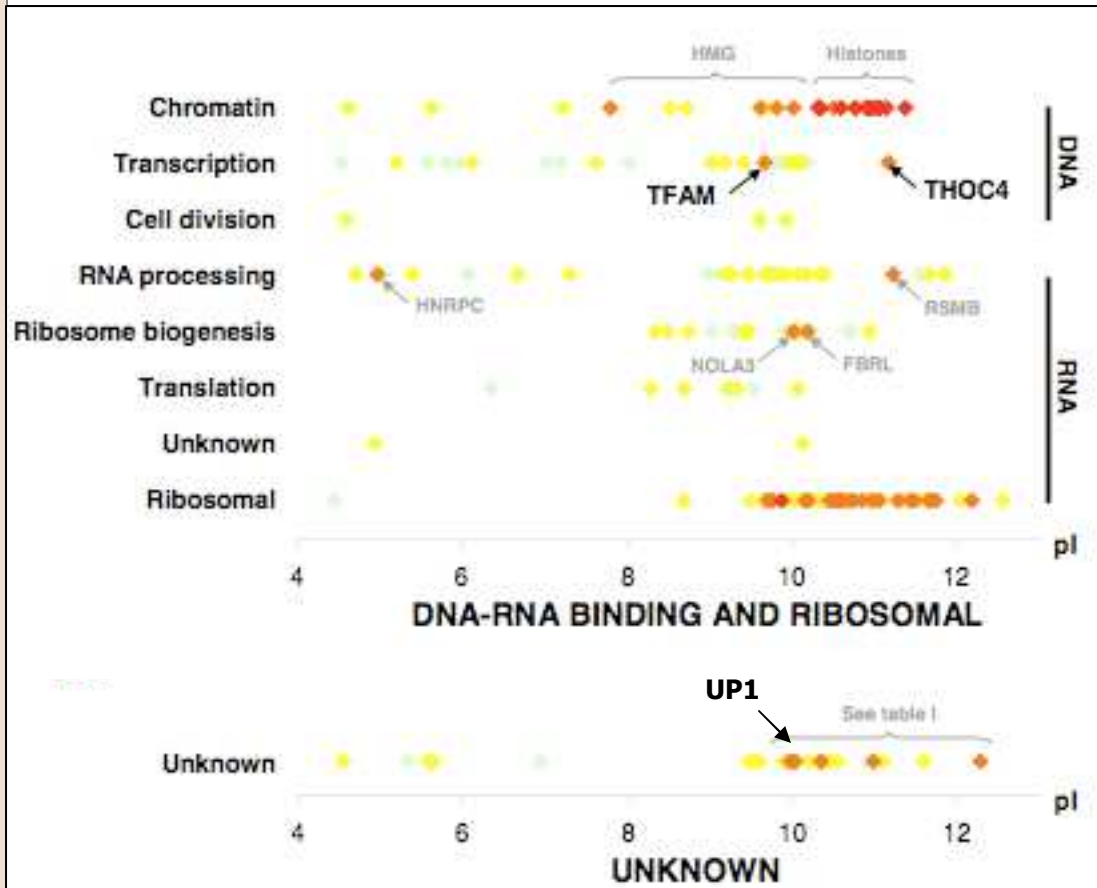
d) protein significance analysis

f) class prediction

h) pathway analysis

→ L. Gatto, Axe 2, 29/11

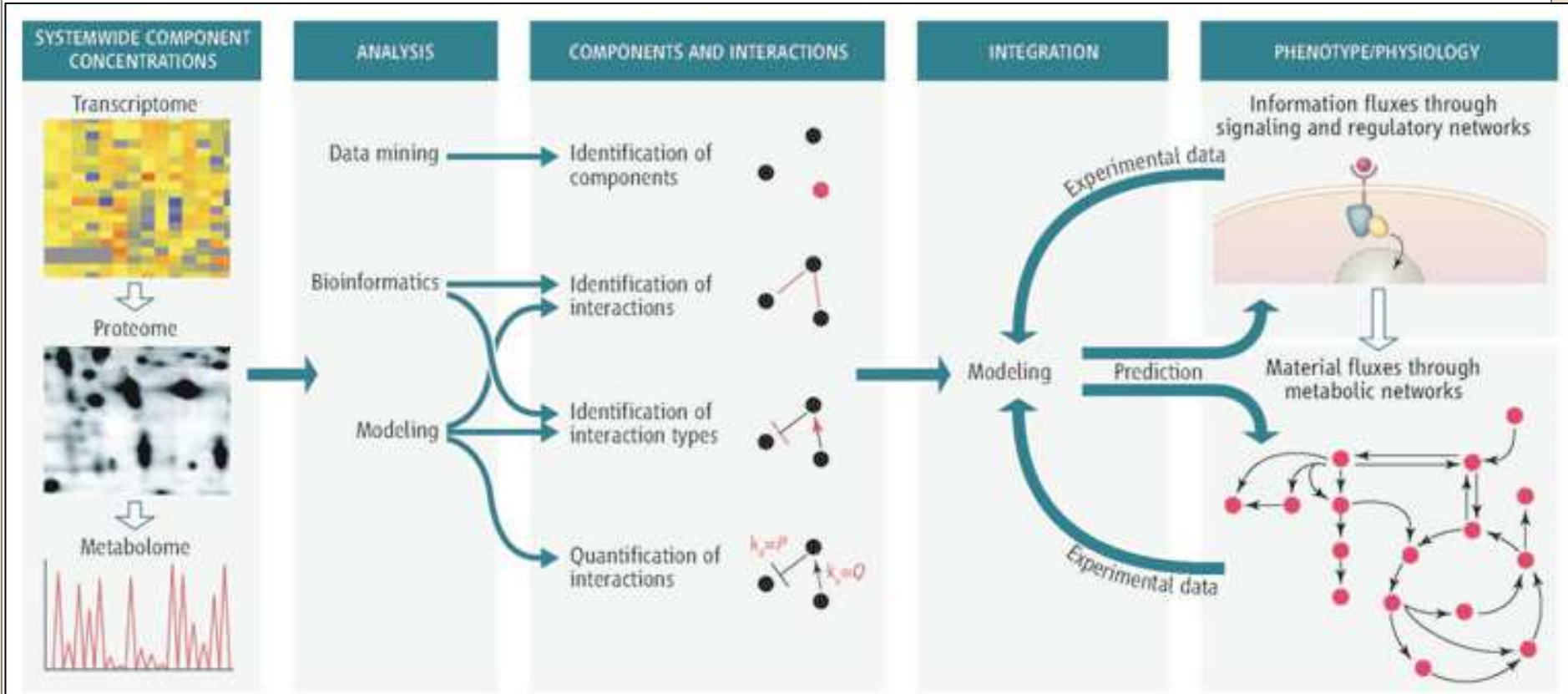
Integration of OMICS data



Proteomics: detection of proteins with high expression levels

Transcriptomics: differential expression between normal and lymphoma cells (GEO database; Basso set)

Systems Biology



→ G. Launay, Axe 2, 29/11

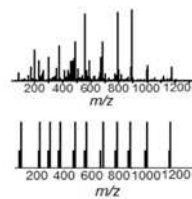
→ Axe 2, 30/11

→ L. Tichit, Axe 3, 29/11

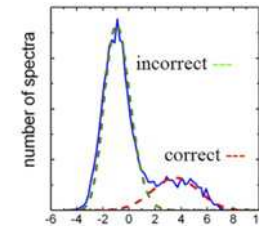
Schematic overview of a typical workflow of the proteomics informatics processing of a data set



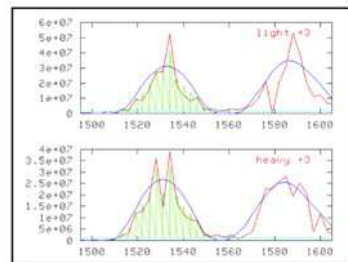
1. Conversion to and use of open data formats



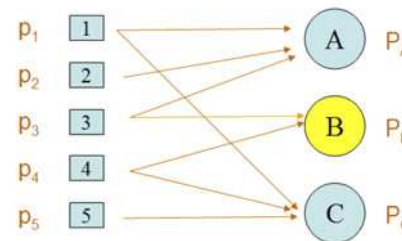
2. Spectrum identification with a search engine



3. Validation of identifications



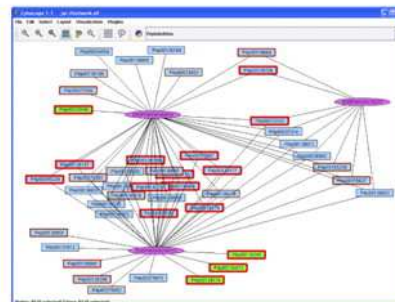
5. Quantification



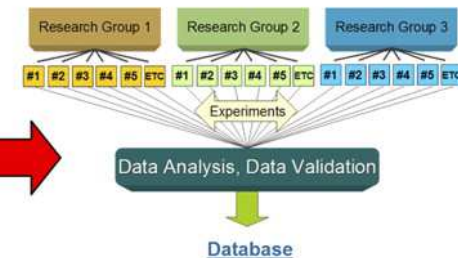
4. Protein inference



6. Organization in local data management systems



7. Interpretation of the protein lists



8. Transfer to public data repositories

Some resources and references

- <http://www.proteomicstutorials.org>
- Cottrell JS. Protein identification using MS/MS data. *J Proteomics*. 2011 Sep 6;74(10):1842-51.
- Beck M, Claassen M, Aebersold R. Comprehensive proteomics. *Curr Opin Biotechnol*. 2011 Feb;22(1):3-8.
- Deutsch EW, Lam H, Aebersold R. Data analysis and bioinformatics tools for tandem mass spectrometry in proteomics. *Physiol Genomics*. 2008 Mar 14;33(1):18-25.
- Walther TC, Mann M. Mass spectrometry-based proteomics in cell biology. *J Cell Biol*. 2010 Aug 23;190(4):491-500.