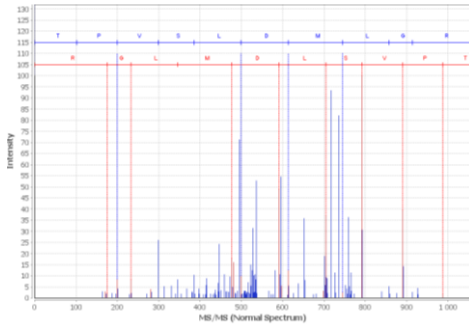


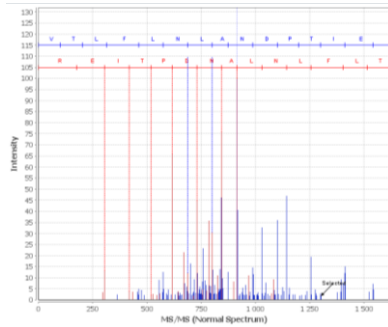
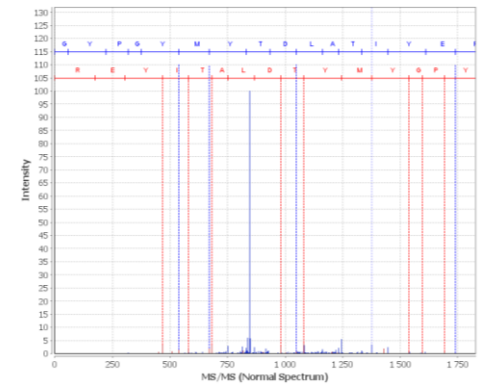


L'identification de protéines séquencées en mode MS/MS : des difficultés encore non résolues.

Dominique Tessier, groupe bioinformatique, plateforme BIBS
Unité BIA, INRA centre Angers-Nantes



Le contexte



4 problématiques majeures en protéomique

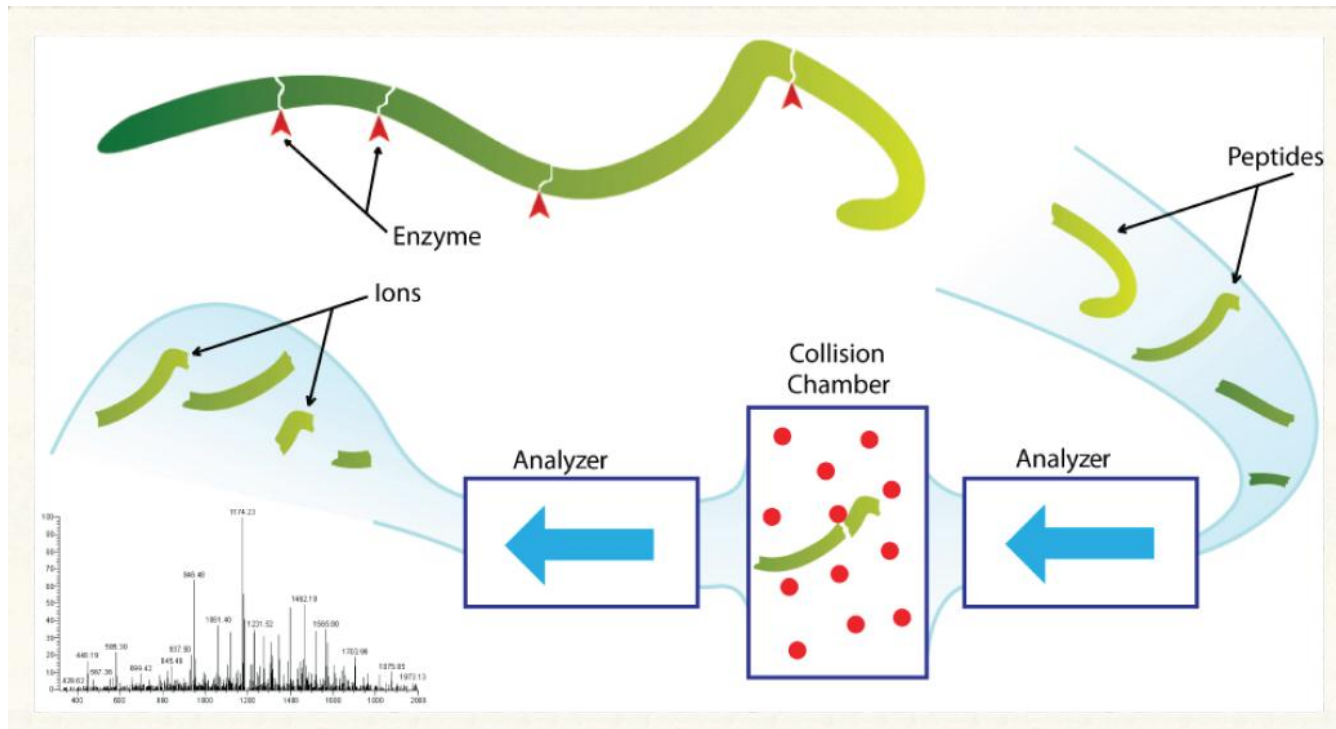
Identification de protéines

Caractérisation de protéines

Quantification de protéines

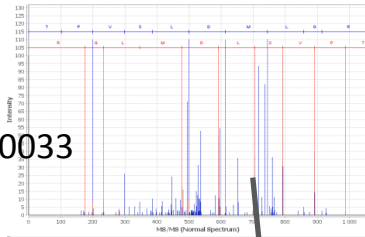
Comparaison d'échantillons

Analyse par spectrométrie de masse en mode MS/MS

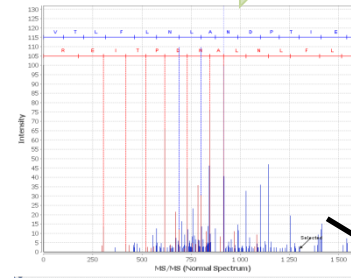




E=0.0033

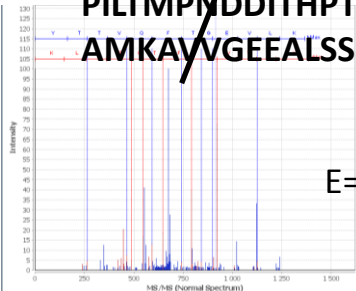


E=0.016

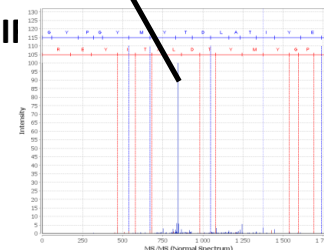


MGLVKDGADMEEGTLEIGMEYRTVSGVAGPLVILDKVKGPKYQEIVNIRLGDGTRR**GQVLEVDGEK**AVVQVFEGTSGID
 NKYTT**VQFTGEVLKTPVSLDMLGR**IFNGSGKPIDNGPPILPEAYLDISGSSINPSERT**YPEEMIQTGISTIDVMNSIARGQKIP**
 LFSAAGLP**INEIAAQICRQAGLVKRLEKGKHAEGGEDDNFAIVFAAMGVNMETAQFFKRDFEENGSMERVTLFLNLAN**
DPTIERIITPRIALTTAEYLAYECGKHVLVILDMSSYADALREVSAAAREEVPGRRGYPGYMYTDLATIYERAGRIEGRTGSITQI
 PILTMPNDDITHPTPDLTGYITEGQIYIDRQLHNR**QIYPPINVLPSLSRLMKS**AIGEGM**TRRDHSDVSNQLYANYAIGKDVQ**
 AMKAV**VGEEALSS**EDLLYLEFLDKFERKFVAQGAYDTRNIFQSLDLAWTLLRII**YSRDATH**

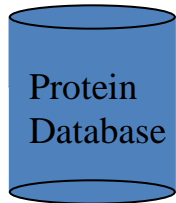
E=0.00037



E=1.4E-005



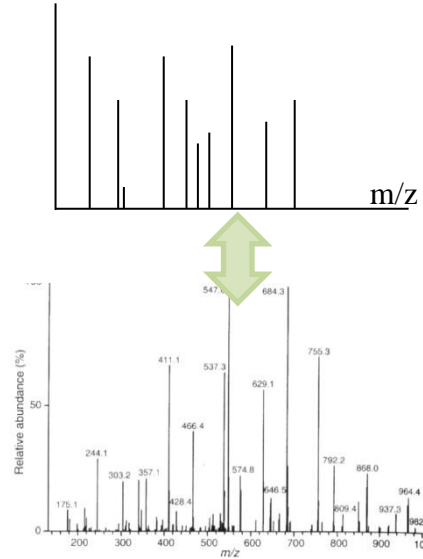
Interprétation : par comparaison...



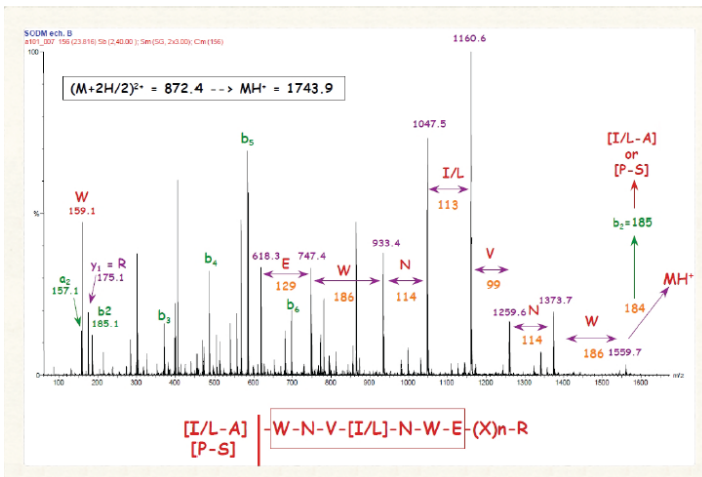
..... • •
NECFLSHK
DDSPDLPK
WVTFISLLLLLFSSAYSR
 • •

Digestion *in silico*

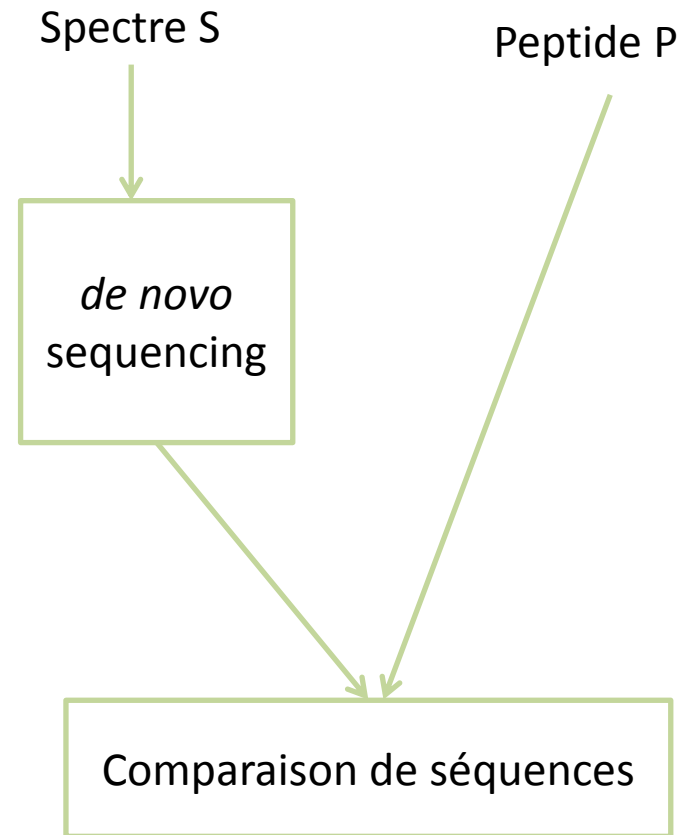
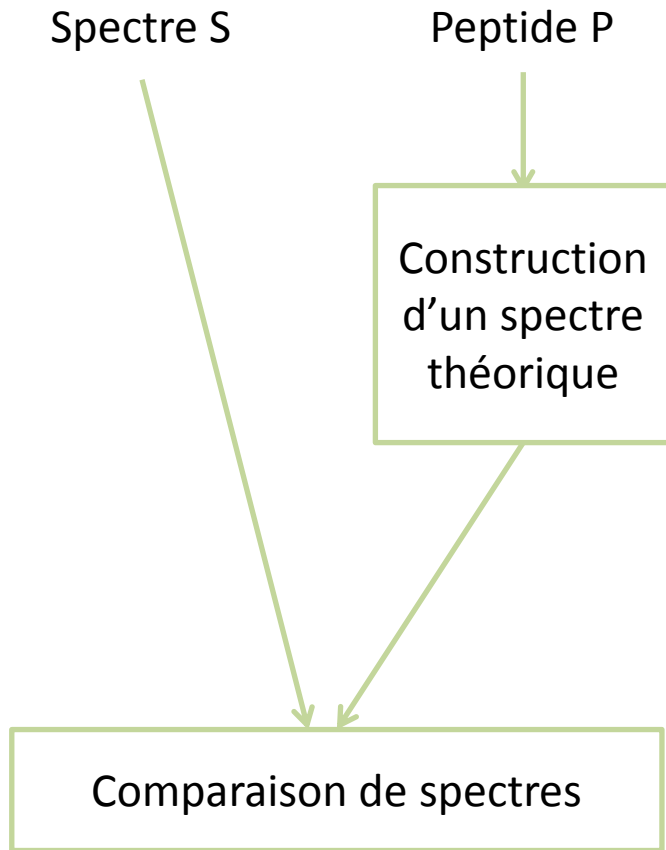
..ou interprétation *de novo*



Comparaison entre les spectres théoriques et les spectres expérimentaux



Interprétation des écarts de masse entre les pics



Plan de la présentation

Le contexte

L'association spectre-peptide par approche comparative (Peptide-Spectrum Match)

L'association spectre-peptide par approche *de novo*

L'identification des protéines

La définition de « pipeline » d'analyse

Identifier et caractériser des protéines du point de vue de l'analyse bioinformatique

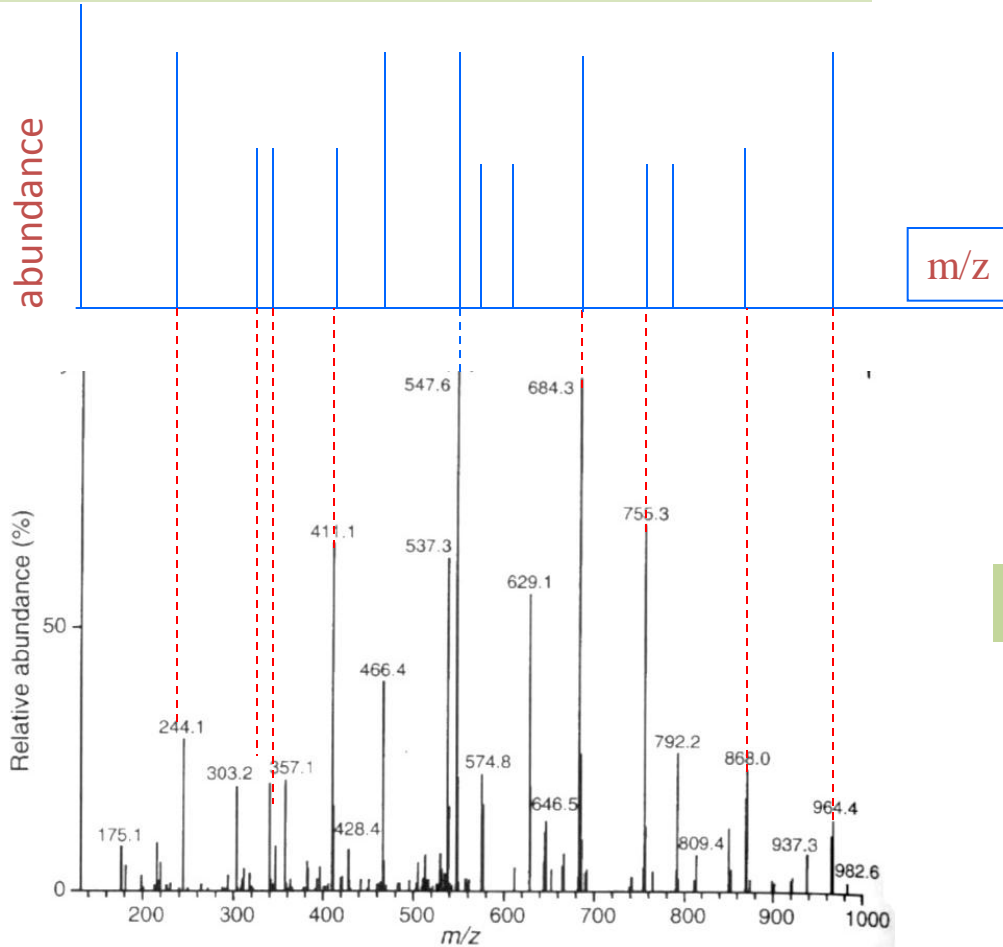
Sélectionner une base de données susceptible de contenir les protéines cherchées.

Comparer les propriétés des protéines cherchées avec les propriétés des protéines de la base de données – les masses des peptides, les propriétés physico-chimiques, l'espèce etc....-

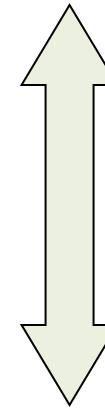
Retenir les meilleures candidates

Calculer les probabilités pour évaluer les chances que les protéines prédites soient les protéines réelles

Comparaison de spectres: Shared peak count



Spectre théorique



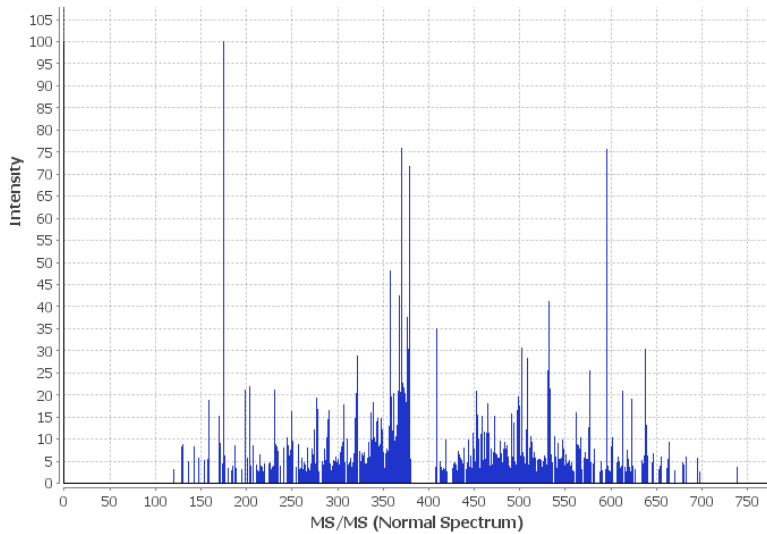
Une fonction de score évalue la similarité entre les 2 spectres

Spectre expérimental

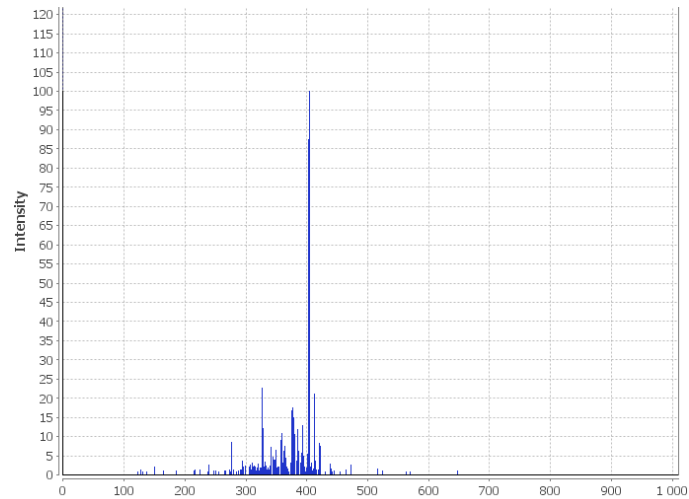
Mascot, X!Tandem, Sequest....

Une analyse en spectrométrie de masse génère des milliers de spectres

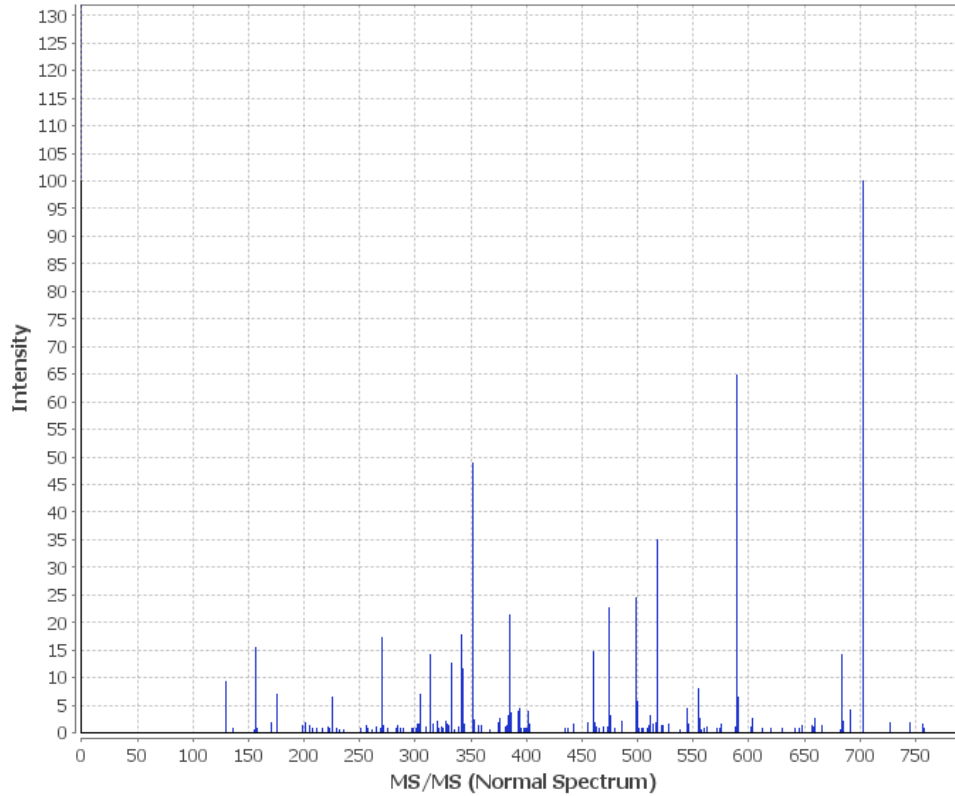
Des spectres très denses



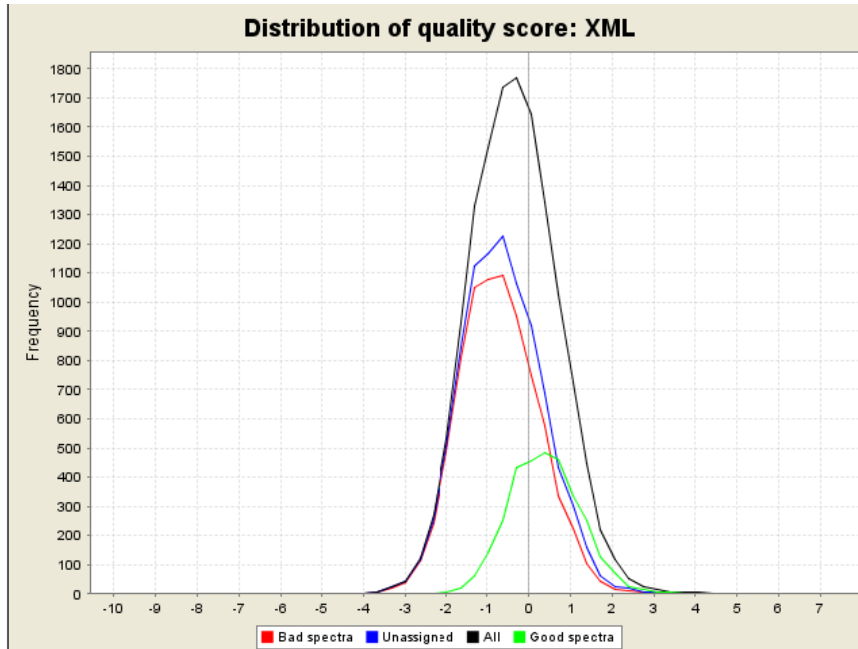
Des spectres peu informatifs



Mais aussi des spectres de bonne qualité:



De nombreux spectres de bonne qualité ne sont pas interprétés.



Qualscore : logiciel d'évaluation de la qualité des spectres

L'interprétation des spectres de masse pour l'identification et la caractérisation des protéines reste un problème ouvert

Nesvizhskii, A.I., Roos, F.F., Grossmann, J., Vogelzang, M., Edes, J.S., Grussem, W., Baginsky, S. and Aebersold, R. (2006) Dynamic spectrum quality assessment and iterative computational analysis of shotgun proteomic data: toward more efficient identification of post-translational modifications, sequence polymorphisms, and novel peptides, Mol Cell Proteomics, 5, 652-670.

Les raisons du faible taux de reconnaissance de spectres

- . Echantillon de mauvaise qualité
- . Bruit
- . Précision de masse insuffisante

- . Des clivages qui ne correspondent pas aux coupures de l'enzyme
- . Du polymorphisme de séquence
- . Des modifications post-traductionnelles
- . Superposition de peptides

On aimerait bien améliorer le taux d'interprétation !

Oui, mais.....

Quel est le coût réel des identifications erronées ?

=> opportunités manquées

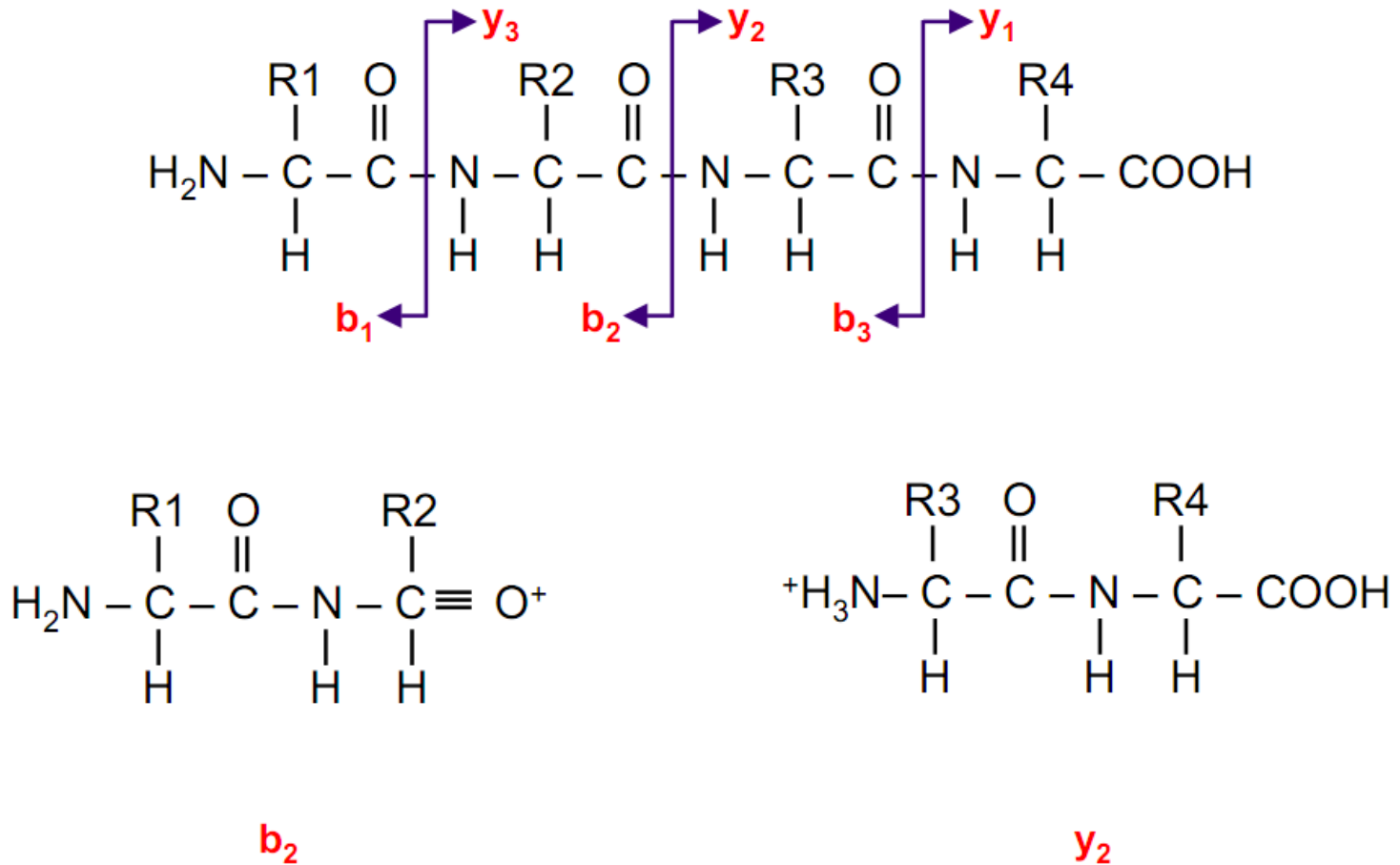
=> du temps perdu à suivre de fausses pistes

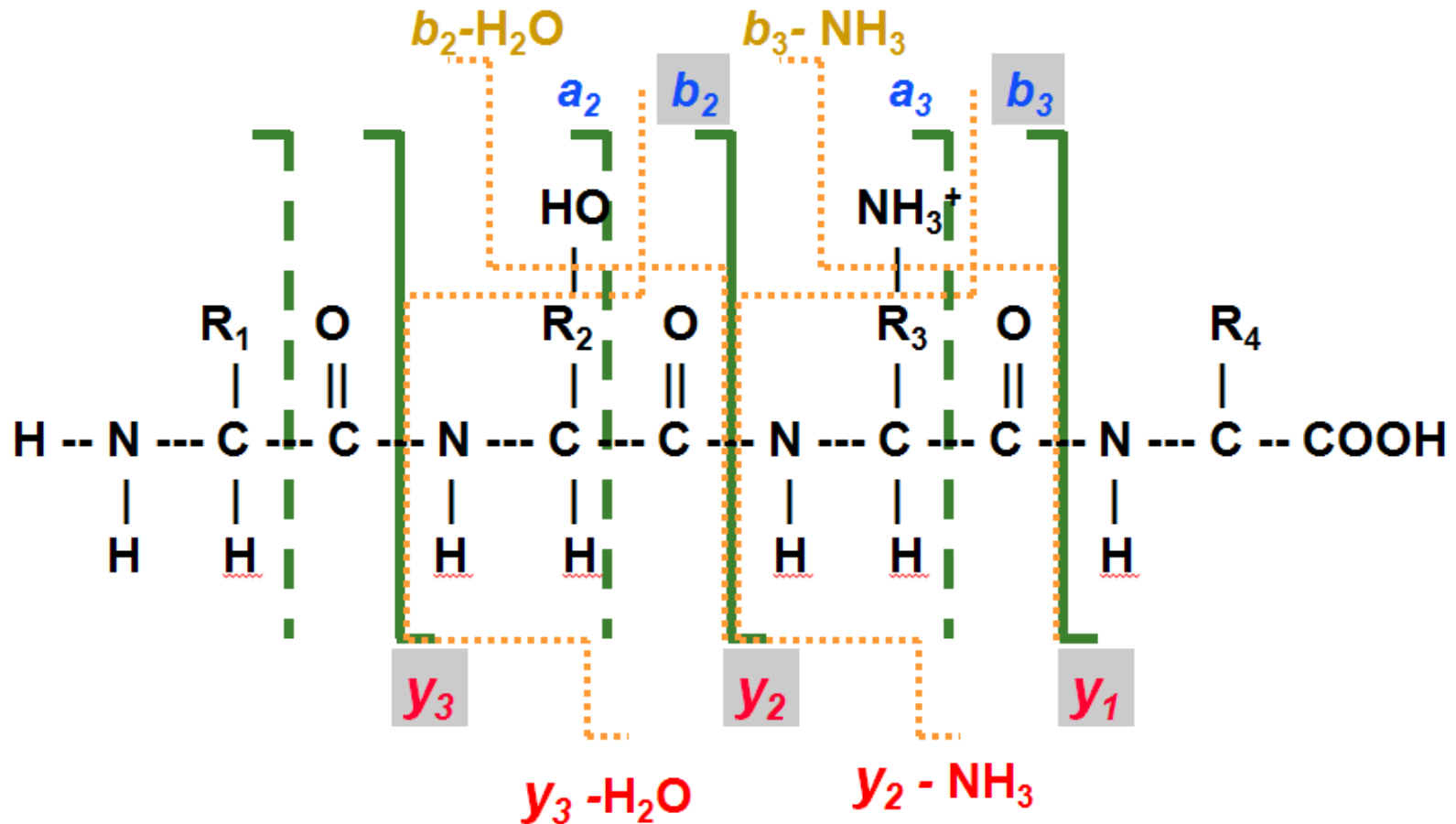
=> la pollution des BD avec des données fausses

=> manque de confiance dans les données de protéomique

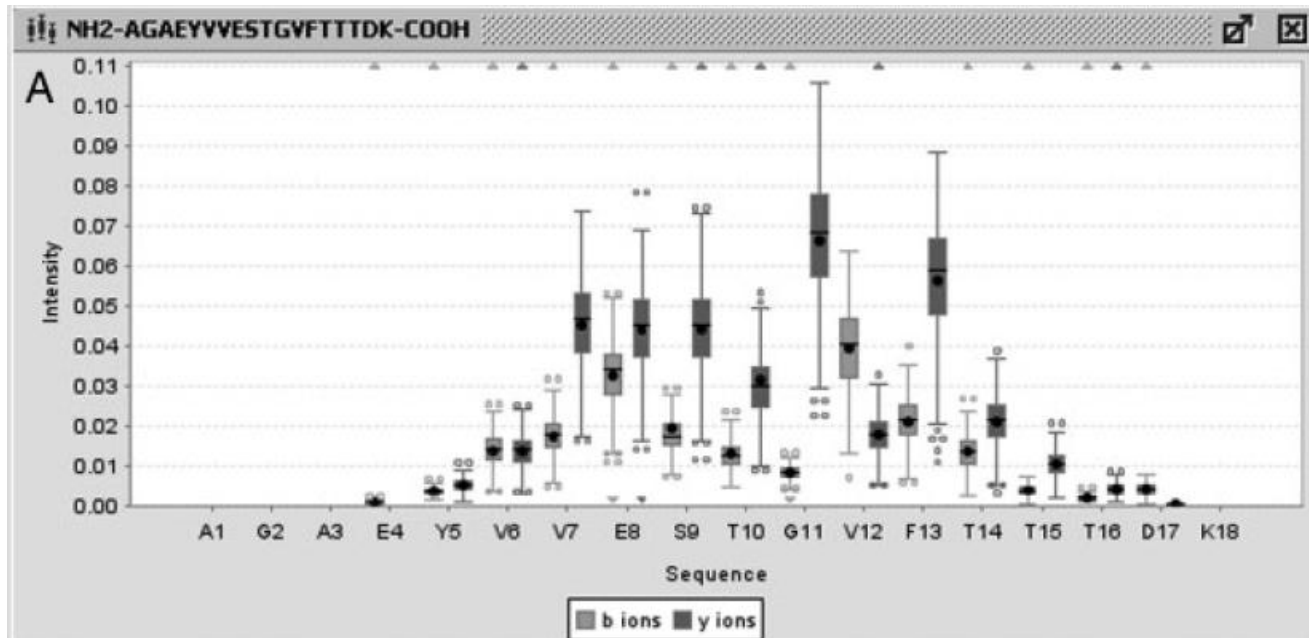
White, F.M. (2011) The potential cost of high- throughput proteomics, *Sci Signal*, **4**, pe8

Fragmentation idéale.....



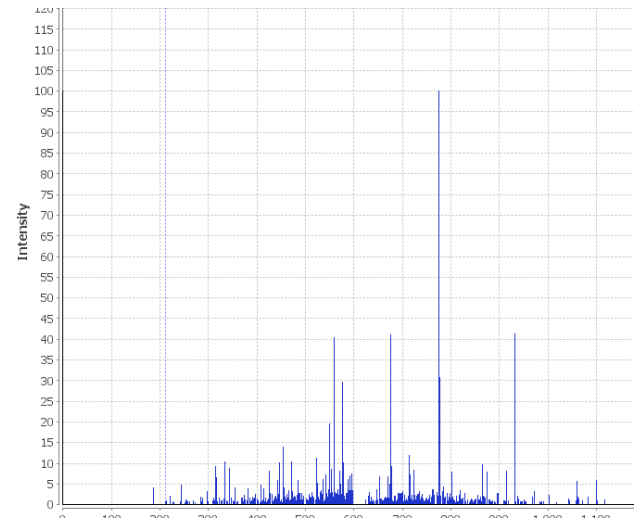
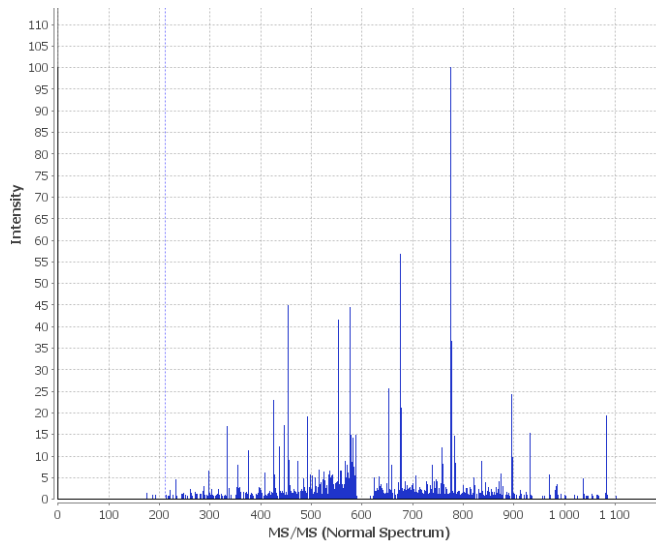
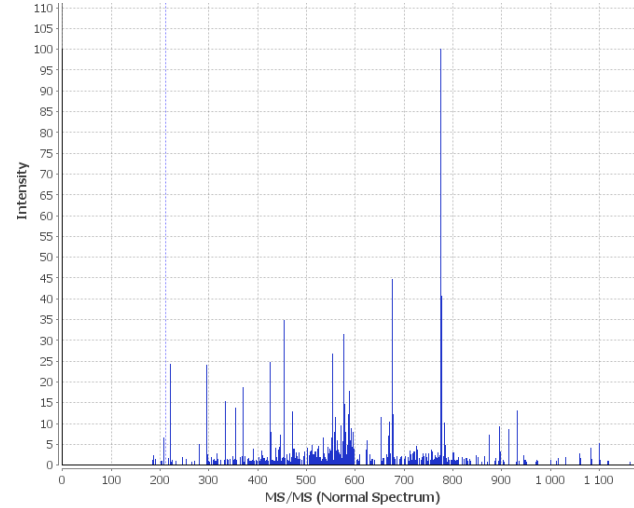
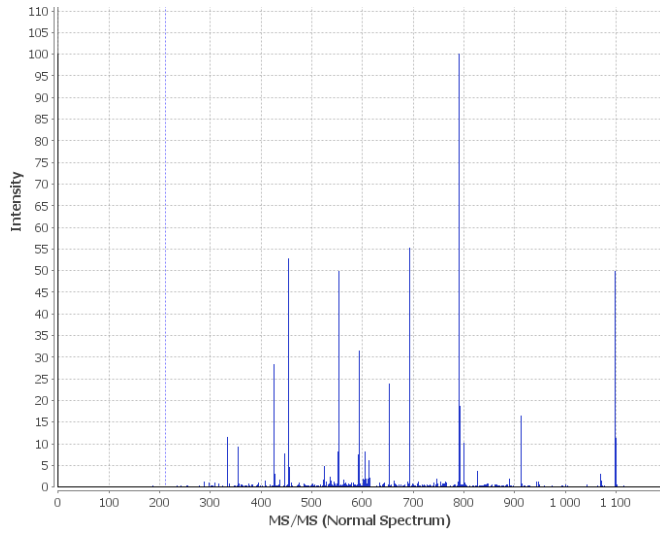


Analyse de la fragmentation de 856 identifications du peptide NH₂-AGAEYVVESTGVFTTTDK-COOH – pics correspondants aux ions y et b (Orbitrap)



FragmentationAnalyser: An open source tool to analyze MS/MS fragmentation data. H. Barsnes et. Al.. Proteomics 2010

Une assez grande variabilité des données générées



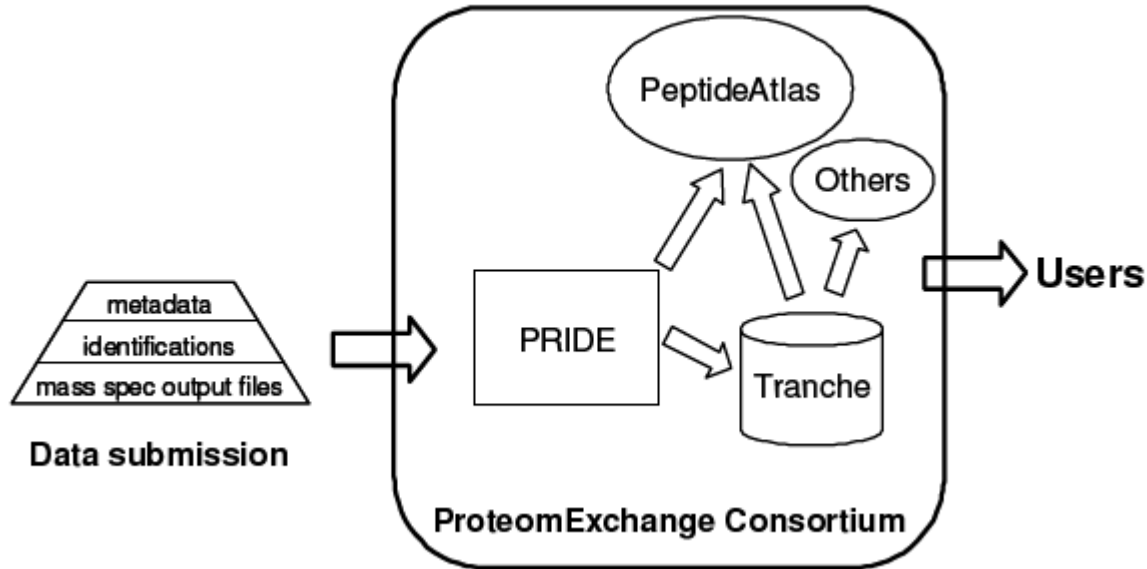
b-ion -H₂O
b-ion -NH₃
y-ion -H₂O
y-ion -NH₃

a
a-H₂O
a-NH₃

bruit,
contaminants

4 spectres issus de l'analyse d'un même échantillon représentant de manière sûre le même peptide EGEGVVMLWK (60% des pics les plus intenses sont affichés, pvalue > 0.99)

Collections non redondantes de spectres obtenus par LC MS/MS



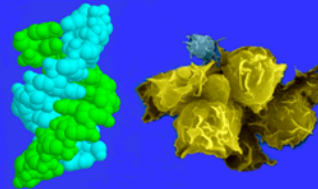
PRIDE Statistics – Novembre 2012

26 589	Experiments
11 380 511	Identified Proteins
64 110 228	Identified Peptides
5 528 099	Unique Peptides
337 882 643	Spectra



The Association
of Biomolecular
Resource Facilities

Research, Technology, Communication, Education



<http://www.abrf.org>

ABRF Annual meeting



The ABRF is an international society dedicated to advancing core and research biotechnology laboratories through research, technology, communication and education. Membership comprises over 700 scientists and administrators, representing over 300 research institutions, who volunteer their efforts on many Research Groups and Committees to provide a unique infrastructure for the benefit of core facilities and their industrial partners.

ABRF Research Groups

benchmark reagents, instruments, software and individual laboratory performance in response to the interest of the members

- ✓ Antibody Technology (ARG)
- ✓ DNA Sequencing (DSRG)
- ✓ Genomic Variation (GVRG)
- ✓ Glycoprotein (gPRG)
- ✓ Light Microscopy (LMRG)
- ✓ Metabolomics (MRG)
- ✓ MicroArray (MARG)
- ✓ Molecular Interactions (MIRG)
- ✓ Nucleic Acids (NARG)
- ✓ Protein Expression (PERG)
- ✓ Protein Sequencing (PSRG)
- ✓ Proteome Informatics (iPRG)
- ✓ Proteomics (PRG)
- ✓ Proteomics Standards (sPRG)

Standardization and Guidelines



Quantitative microscopy



Benchmark the applications of Next Gen sequencers



Use the ABRF infrastructure to

launch your new Research Group, Committee, or Meeting

Core Administrators Network

To promote communication between core personnel and core administrators, and among core administrators. This is coordinated by a new Core Administrators Network – Coordinating Committee (CAN-CC).

ABRF Affiliates & Chapters

Infrastructure assistance to regional interest groups to have meetings that promote common interests. This is facilitated by a new Affiliates and Chapters Committee (ACC)



NERLSCD

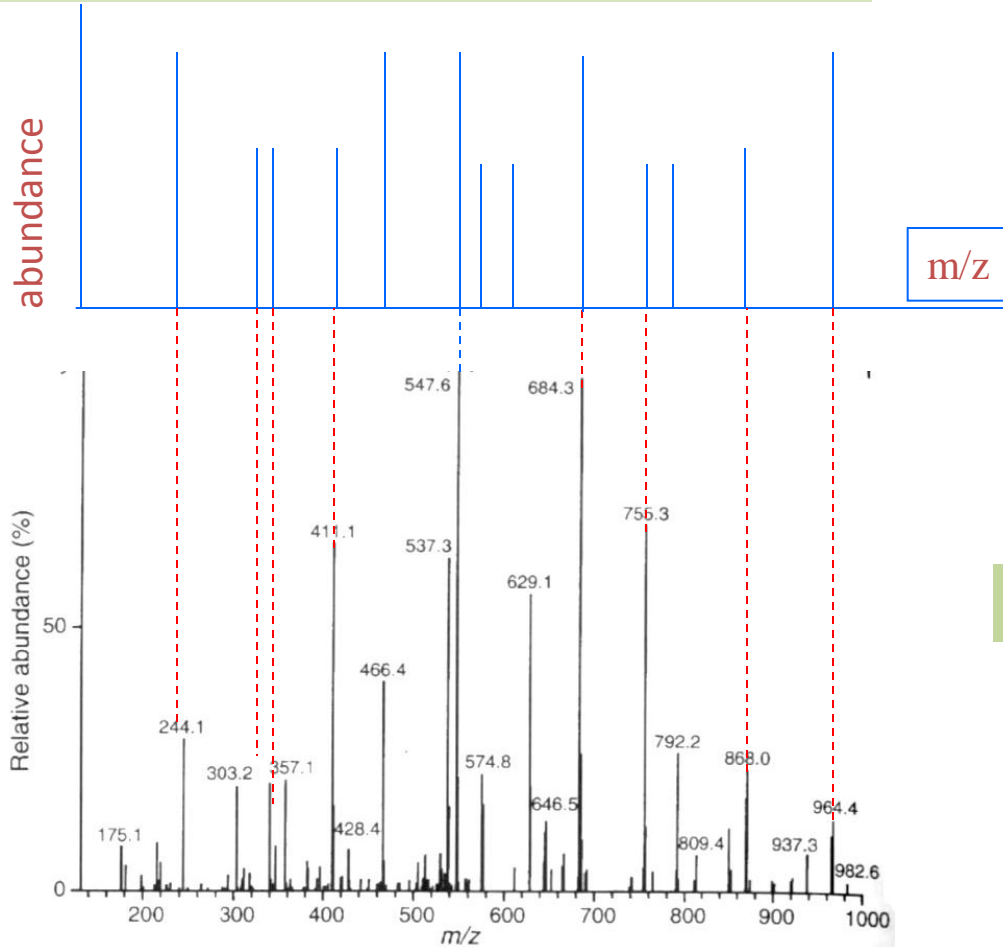
ABRF Committee service is an excellent means of active contribution by our members and a significant networking opportunity

- ✓ Executive Board
- ✓ ABRF Award
- ✓ Affiliates and Chapters
- ✓ Career Development
- ✓ Core Administrators Network - Coordinating Committee
- ✓ Corporate Advisory
- ✓ Corporate Relations
- ✓ Education
- ✓ Finance and Investments
- ✓ Membership
- ✓ Nominations
- ✓ Publications
- ✓ Survey
- ✓ Travel Award
- ✓ Web Site

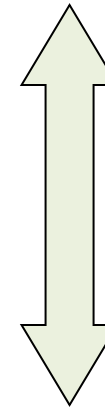
Sponsors



Comparaison de spectres: Shared peak count



Spectre théorique



Une fonction de score évalue la similarité entre les 2 spectres

Spectre expérimental

Mascot, X!Tandem, Sequest....

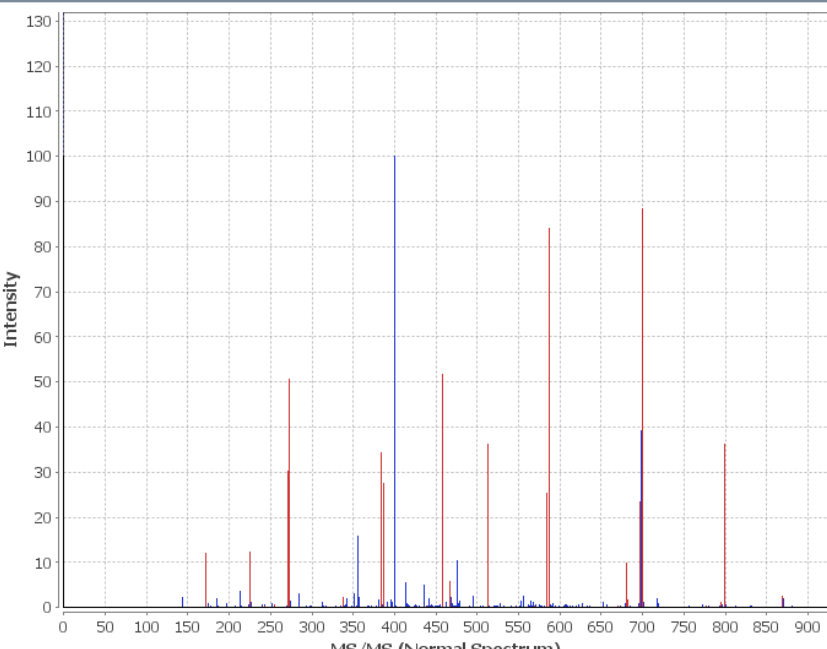
Exemple: un spectre de bonne qualité (qualscore > 3) bien interprété

PacketSpectralAlignment Viewer

File Edit Peptide

Experimental spectrum : Spectrum1830 scans: 3718, Fragment Mass : 967,5469 m/z : 484,7807 Number of experimental peaks : 195

Theoretical sequence : VAVLEANPR Peptide Mass : 967,5483 Mass delta : -0,0015 Number of matched peaks : 27

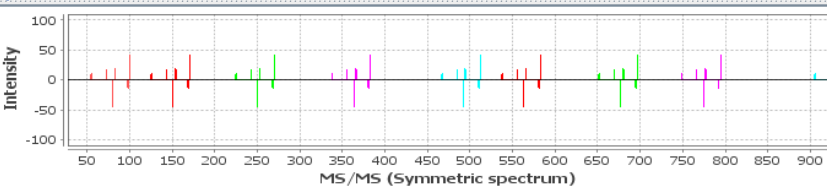


Intensity

MS/MS (Normal Spectrum)

Theoretical peaks :					
Mass	Intensity	Position	Ion	Matched	
968,554		43,5 C-Ter	y9		
951,523		12 C-Ter	y*9		
950,539		13 C-Ter	y-H2O9		
869,486		43,5 C-Ter	y8	yes	
852,455		12 C-Ter	y*8		
851,471		13 C-Ter	y-H2O8		
100,074		41,5 N-Ter	b1		
82,063		19,5 N-Ter	b01		
83,047		18 N-Ter	b*1		
72,07		17 N-Ter	a1		
54,069		8,5 N-Ter	a01		
55,053		10 N-Ter	a*1		
798,449		43,5 C-Ter	y7	yes	
781,418		12 C-Ter	y*7		
780,434		13 C-Ter	y-H2O7		

Experimental peaks :			
Mass	Intensity	Ion	
399.863	1446550.0		
699.361	1276820.0	y6	
586.278	1215920.0	y5	
457.233	748581.0	y4	
272.21	729910.0	y2	
698.352	567597.0		
512.258	525144.0	b5	
798.436	521972.0	y7	
383.277	497281.0	b4	
270.221	438939.0	b3	
587.238	432381.0		
700.337	399486.0		
386.241	398156.0	y3	
583.292	365239.0	b6	
697.362	338005.0	b7	



Intensity

MS/MS (Symmetric spectrum)

Packet	Amino-acid	Position	Score
0-		0	0
1V		99,068	0,87
2A		170,106	1,7
3V		269,174	2,74
4L		382,258	2,14
5E		511,301	2,14
6A		582,338	1,7
7N		696,381	2,9

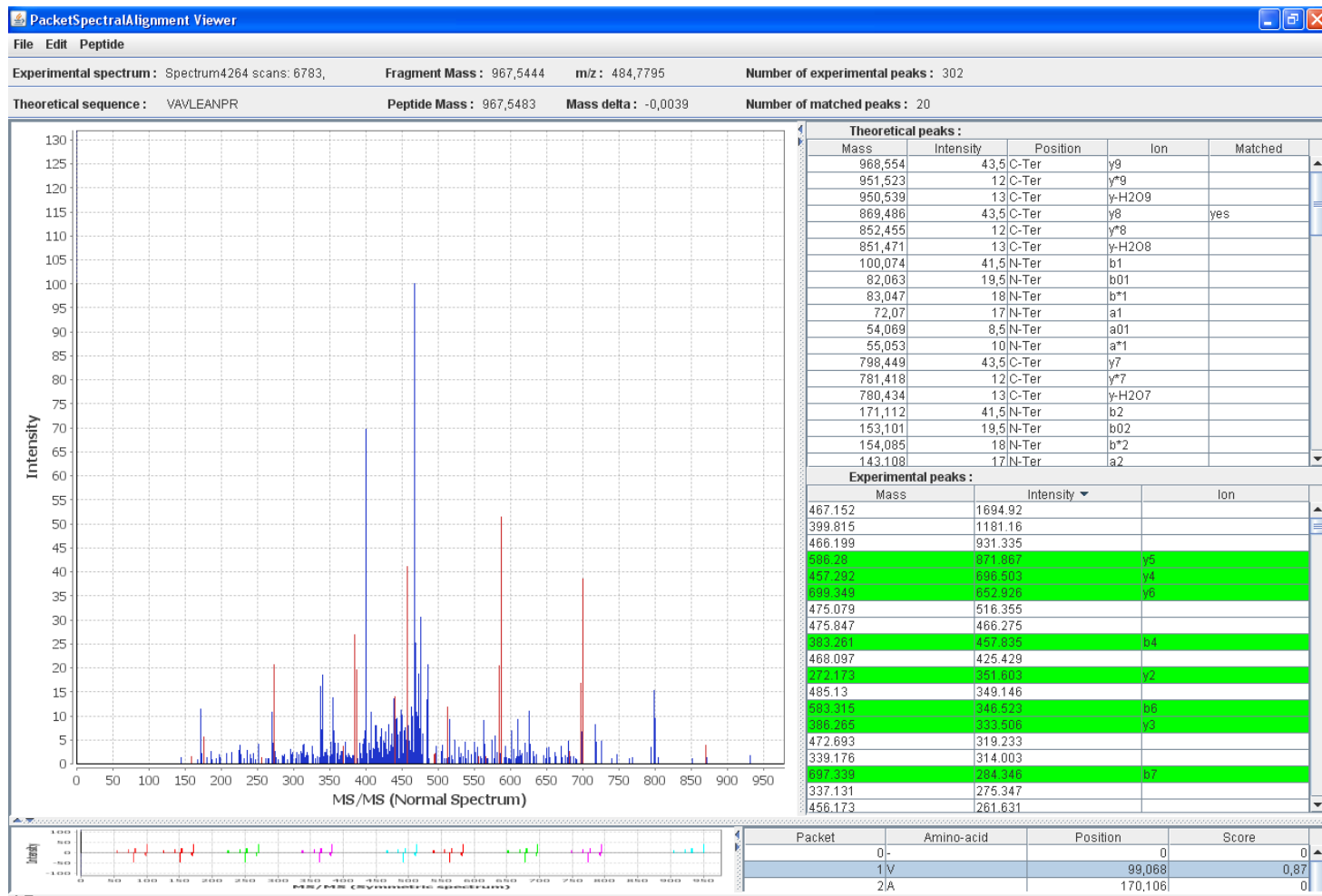
Peptide (type sequence here):

Precision and selection: Accuracy: Threshold:

Scores: Score: Score/Nb-aminoacids:

pvalue=0.9983
 Evalue=0.047
 Quallscore=3.55

Exemple : un spectre bien interprété pourtant de mauvaise qualité



pvalue=0.9924
 evalue=0.091
 Qualscore=-1.09

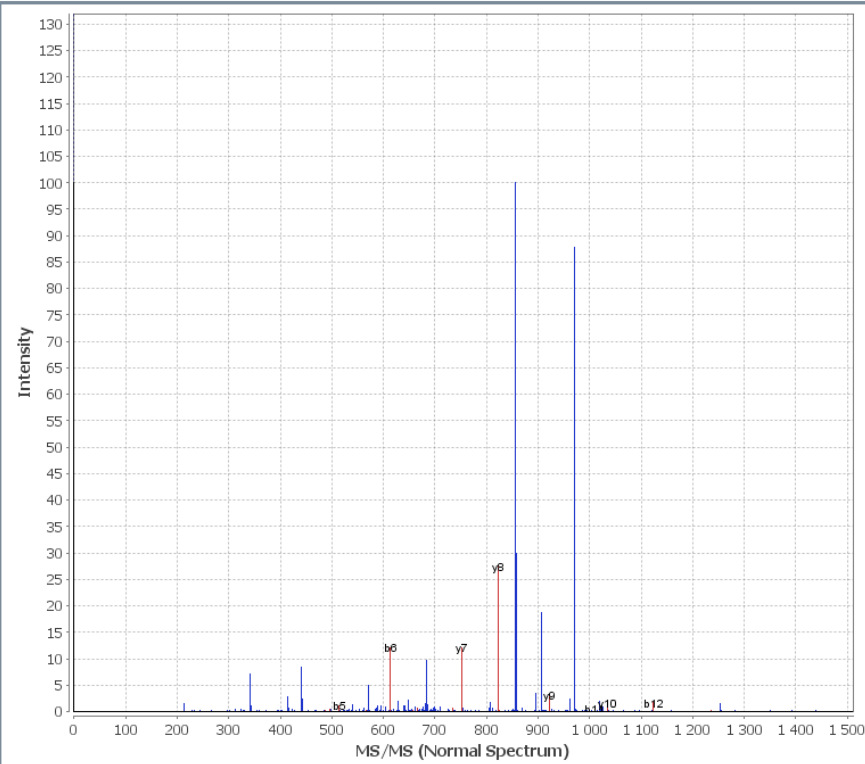
Un spectre de bonne qualité et une interprétation incertaine

PacketSpectralAlignment Viewer

File Edit Peptide

Experimental spectrum : Spectrum3498 scans: 5788, Fragment Mass : 1 434,709 m/z : 718,3618 Number of experimental peaks : 280

Theoretical sequence : EGDVIVAPAGTLMY Peptide Mass : 1 434,7098 Mass delta : -0,0009 Number of matched peaks : 21



Theoretical peaks :					
Mass	Intensity	Position	Ion	Matched	
1 435,716		43,5 C-Ter	y14		
1 418,685		12 C-Ter	y*14		
1 417,701		13 C-Ter	y-H2O14		
1 306,673		43,5 C-Ter	y13		
1 289,642		12 C-Ter	y*13		
1 288,658		13 C-Ter	y-H2O13		
130,049		41,5 N-Ter	b1		
112,038		19,5 N-Ter	b01		
113,022		18 N-Ter	b*1		
102,045		17 N-Ter	a1		
84,044		8,5 N-Ter	a01		
85,028		10 N-Ter	a*1		
1 249,652		43,5 C-Ter	y12		
1 232,621		12 C-Ter	y*12		
1 231,637		13 C-Ter	y-H2O12		
187,07		41,5 N-Ter	b2		
169,059		19,5 N-Ter	b02		
170,043		18 N-Ter	b*2		
159,066		17 N-Ter	a2		

Experimental peaks :		
Mass	Intensity	Ion
856.28	350048.0	
970.355	307092.0	
971.129	124769.0	
857.094	104987.0	
923.43	89197.1	y9
905.801	65747.6	
613.302	42286.5	b6
752.375	42237.5	y7
684.263	34148.1	
441.319	29267.4	
342.308	24853.9	
906.498	22561.7	
571.36	17534.6	
895.566	12381.4	
922.51	10776.5	y9
413.344	9827.62	
961.272	8660.41	
442.303	8236.91	
647.279	7548.61	

Packet	Amino-acid	Position	Score
0-		0	0
1E		129,043	0
2G		186,064	0

Peptide (type sequence here): EGDVIVAPAGTLMY

Precision and selection: Accuracy: 0,06 Threshold: 0,4

Scores: Score: 10,68 / 51,24 Score/Nb-aminoacids: 0,76

pvalue=0.2626
 evalue=0.46
 Qualscore=4.35

Les algorithmes de comparaison de spectres :

Quels sont les types d'ions considérés ?

Comment sont définies les intensités des spectres théoriques ? *

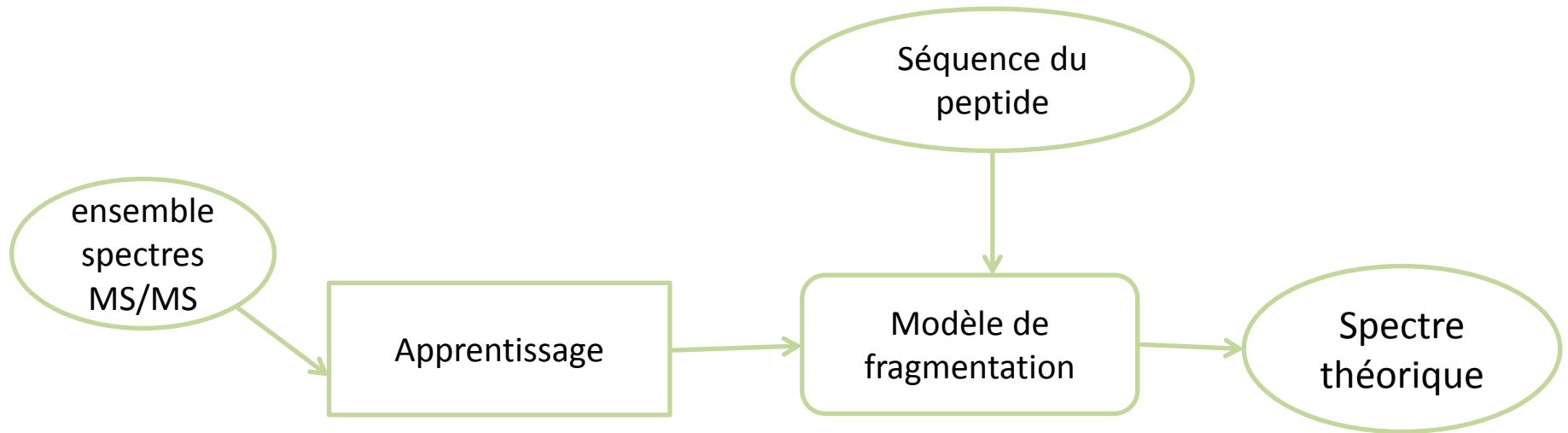
Quel est l'algorithme de comparaison utilisé ?

Quel est le calcul de score ? *

Comment les mutations , les modifications post-traductionnelles sont-elles prises en compte ? *

Comment sont validées les attributions ?

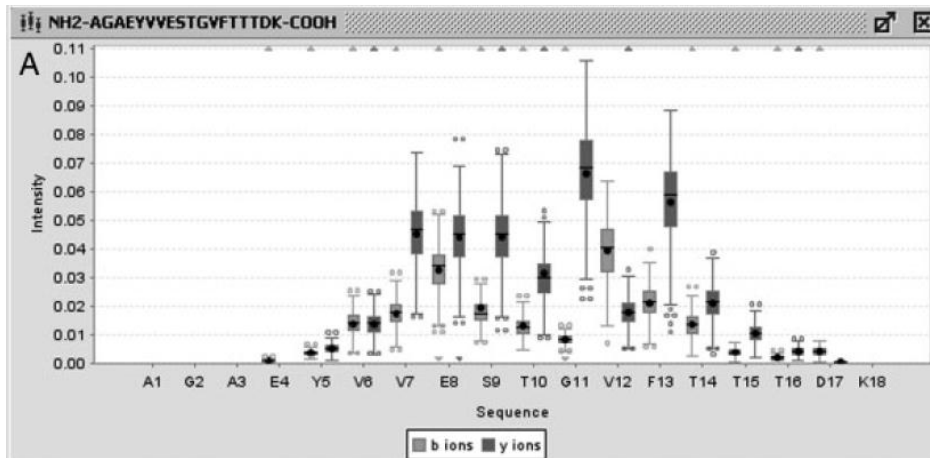
Mieux modéliser un spectre théorique = prédire le modèle de fragmentation



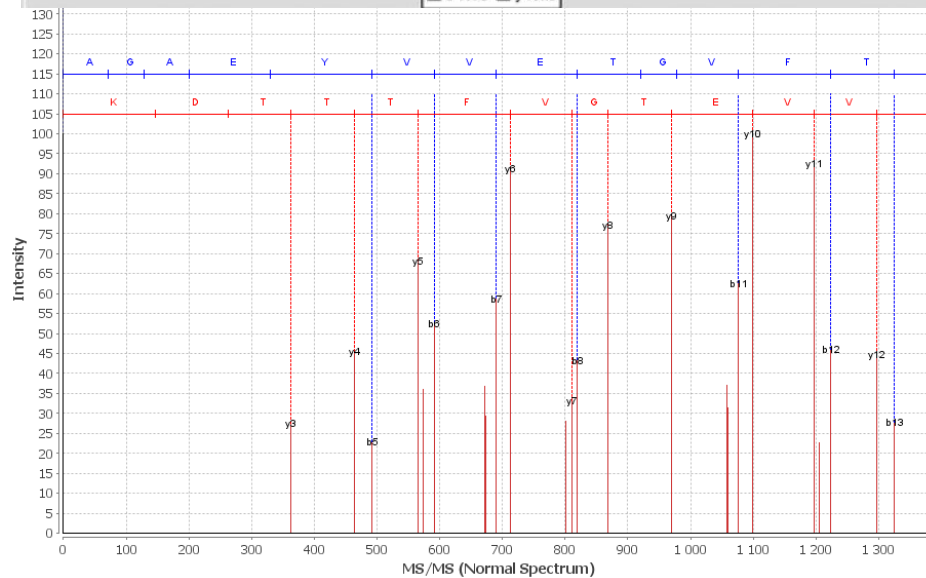
Exemple PepNovo+ : la prédiction ne porte pas sur l'intensité des pics, mais sur l'ordre de l'intensité des pics. Basé sur une méthode RankBoosting

Predicting Intensity Ranks of Peptide Fragment Ions. Frank, A.M. J. Proteome Research, 8:2226-2240, 2009

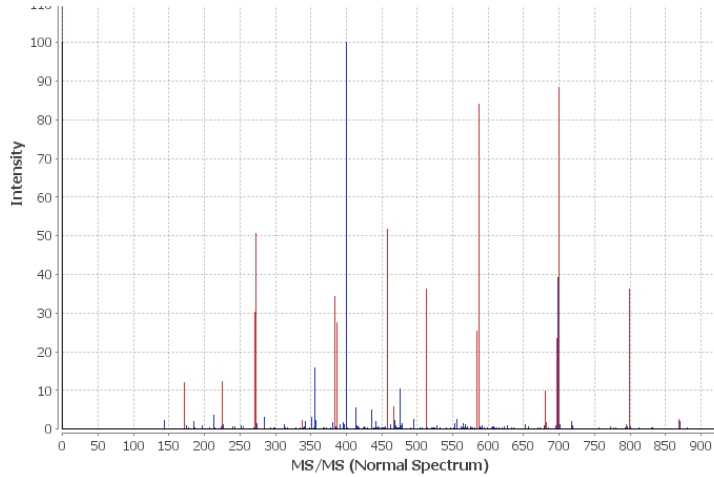
Ranking-Based Scoring Models for Peptide-Spectrum Matches. Frank, A.M. J. Proteome Research, 8:2241-2252, 2009



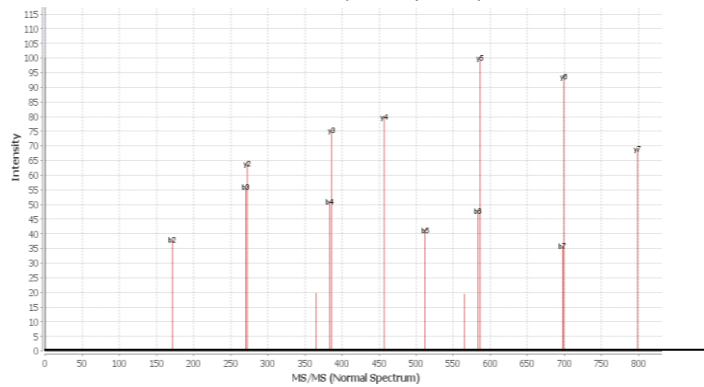
Peptide AGAEYVVETGVFTT DK



Fragmentation prédite
par PepNovo+



Peptide VAVLEANPR



Fragmentation prédite par PepNovo+

Banques de spectres expérimentaux

- ⇒ Des scores plus discriminants qu'en comparaison avec des spectres théoriques
- ⇒ Un espace de recherche plus petit et plus précis
- ⇒ Identification plus facile de spectres sortant de l'ordinaire donc potentiellement porteur de modifications post-traductionnelles
- ⇒ L'intérêt de ces banques décroît lorsque les spectres sont de bonne qualité

NIST : <http://peptide.nist.gov/>

PeptideAtlas : <http://www.peptideatlas.org/speclib/>

PRIDE : <http://www.ebi.ac.uk/pride/>

Les algorithmes de comparaison de spectres :

Quels sont les types d'ions considérés ?

Comment sont définies les intensités des spectres théoriques ? *

Quel est l'algorithme de comparaison utilisé ?

Comment le score est-il calculé ? *

Comment les mutations , les modifications post-traductionnelles sont-elles prises en compte ?

Comment sont validées les attributions ?

Le calcul de score : un exemple simple : X!Tandem

Craig, R. and Beavis, R.C. (2003) A method for reducing the time required to match protein sequences with tandem mass spectra, *Rapid Commun Mass Spectrom*, **17**, 2310-2316.

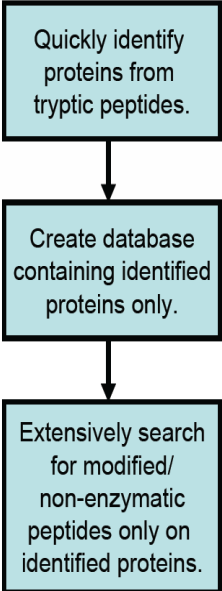
An explanation of the X!Tandem MS/MS spectra search program developed by Craig, R. and Beavis, R.C. (2003) *Rapid Commun. Mass Spectrom.*, **17**, 2310–2316.

Explication complète : www.proteomesoftware.com

X!Tandem_edited.pdf - Adobe Acrobat Pro
Fichier Edition Affichage Fenêtre Aide

PROTEOME SOFTWARE

X!Tandem Workflow



```

graph TD
    A[Quickly identify proteins from tryptic peptides.] --> B[Create database containing identified proteins only.]
    B --> C[Extensively search for modified/non-enzymatic peptides only on identified proteins.]
  
```

X!Tandem, like Mascot and SEQUEST, compares each spectrum to all likely candidate peptides in a protein database.

One of X!Tandem's strengths is its automatic search for modified peptides — but only on proteins it has otherwise identified.

The following pages explain how X!Tandem matches peptides, and how this differs from the way SEQUEST matches them.

Principaux avantages de X!Tandem

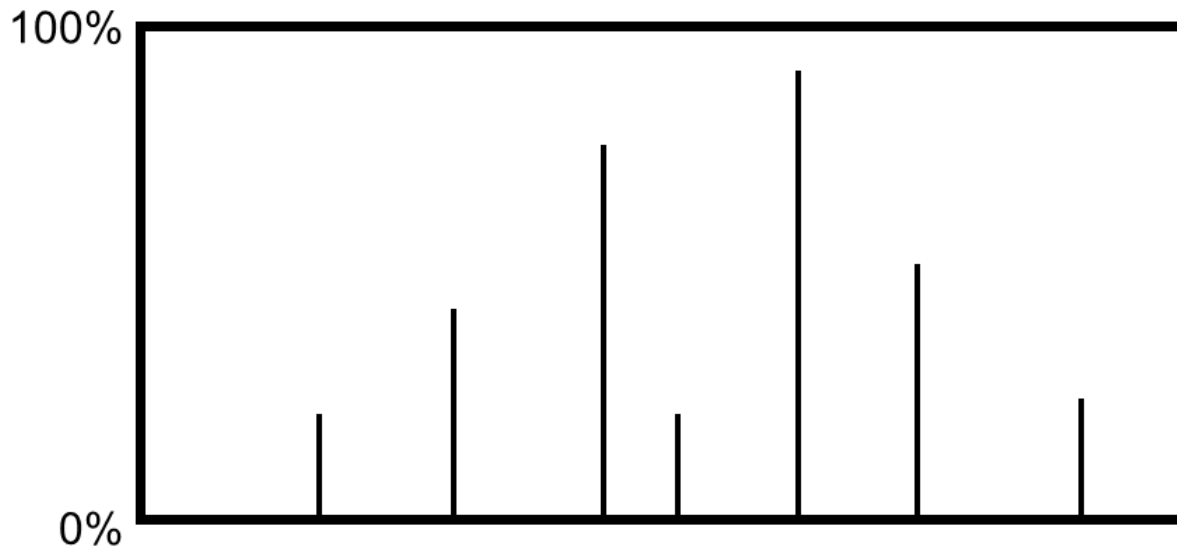
- . Très rapide
- . Prend en compte les peptides semi-tryptiques
- . Prend en compte le polymorphisme
- . Score basé sur des calculs de probabilité

Hyperscore

$$\text{HyperScore} = \left(\sum_{i=0}^n I_i * P_i \right) * N_b! * N_y!$$

spectrum intensities predicted? (1,0)

X!Tandem modifies the preliminary score by multiplying by N factorial for the number of b and y ions assigned. The use of factorials is based on the hypergeometric distribution.

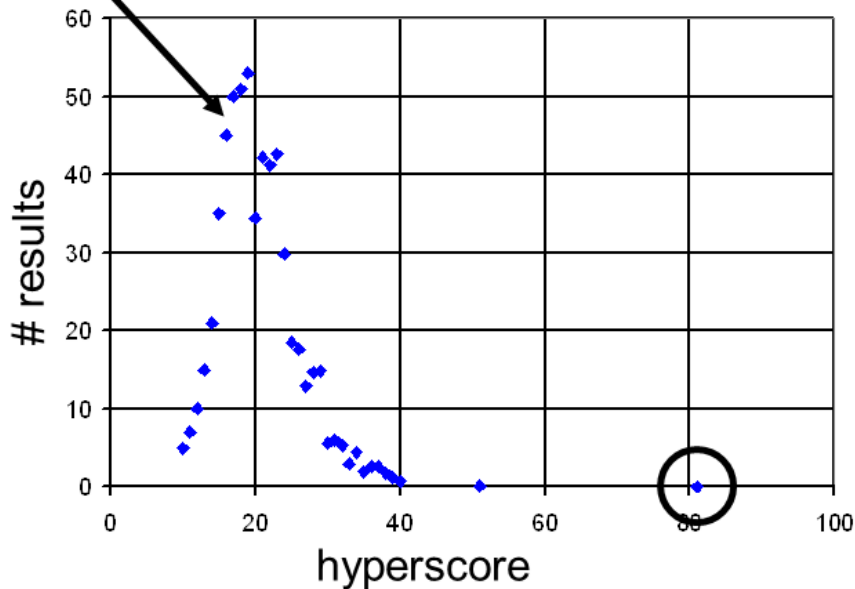


similar
peaks
(y/b ions)



Histogram of hyperscores

incorrect IDs

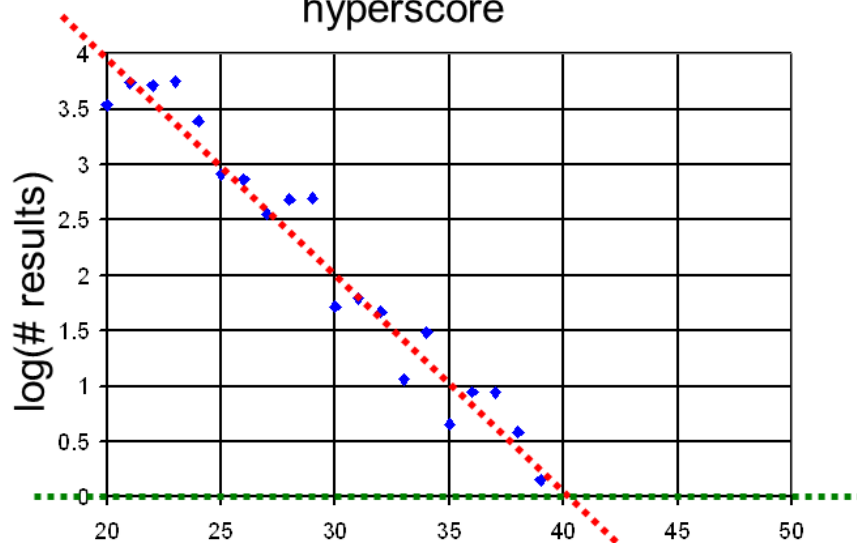
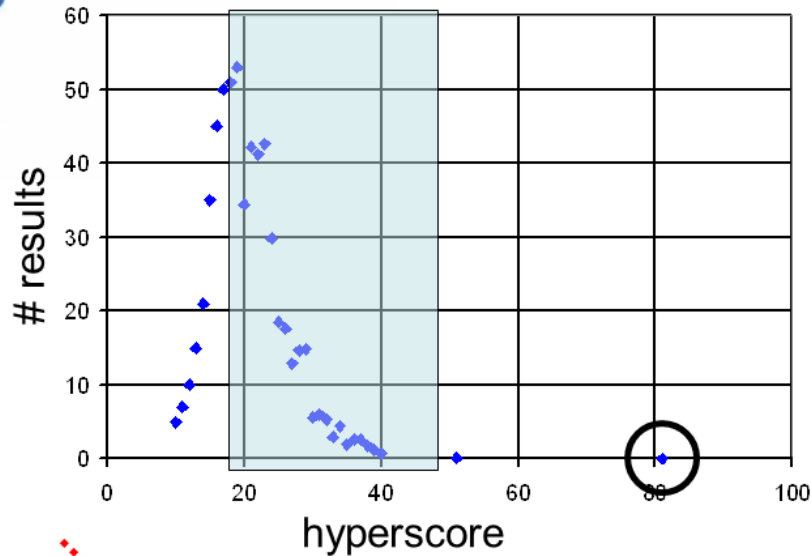


Next, X!Tandem makes a histogram of all the hyperscores for all the peptides in the database that might match this spectrum.

For example, in this figure, 52 peptides were found with a hyperscore of 19, and one peptide with a hyperscore of 83.

X!Tandem assumes that the peptide with the highest hyperscore is correct, and all others are incorrect.

Log histogram

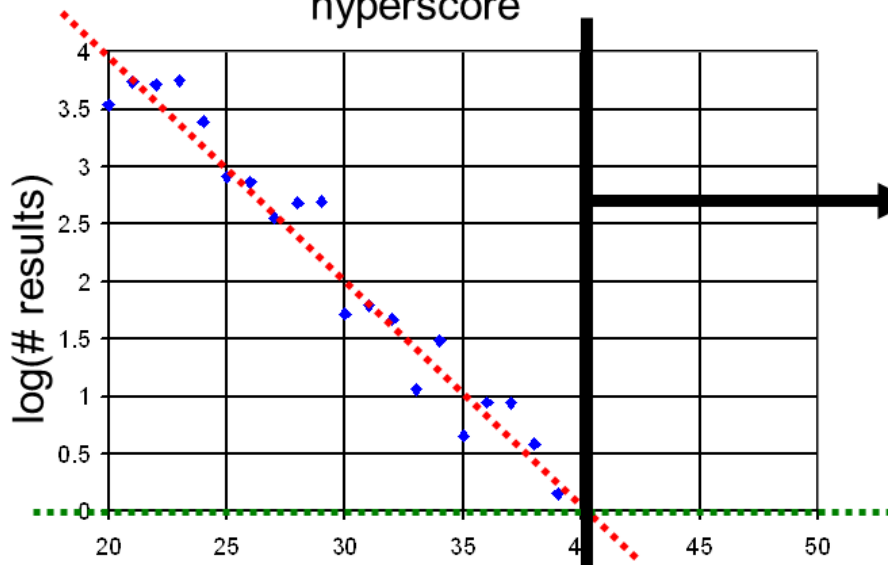
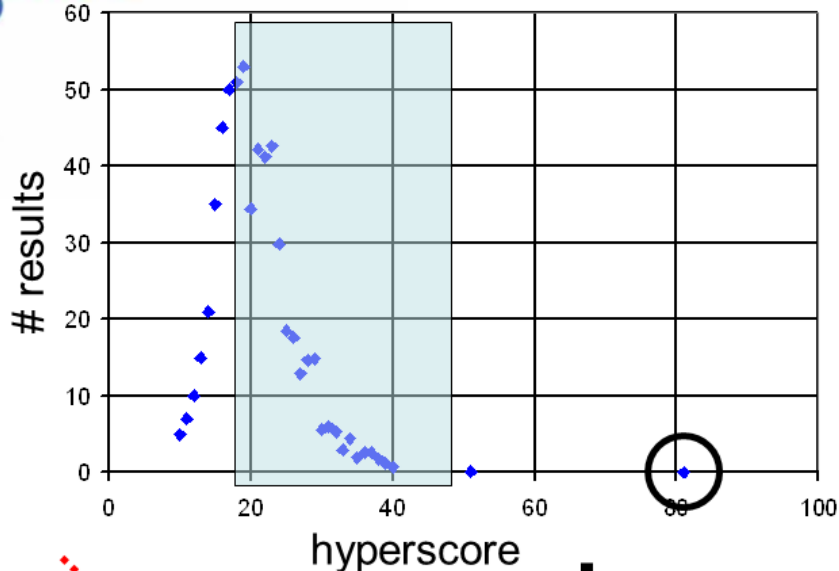


If the data on the right side of the histogram, (colored in upper figure) is taken and log-transformed, the data fall on a straight line.

A straight line is the expected result from a statistical argument that assumes the incorrect results are random.

Note: this histogram is calculated independently for each spectrum.

Significant scores



X!Tandem has already assumed that the top hyperscore is the only possible correct match.

This match is significant if it is greater than the point at which the straight line through the log data intersects the $\log(\# \text{ results})=0$ line.

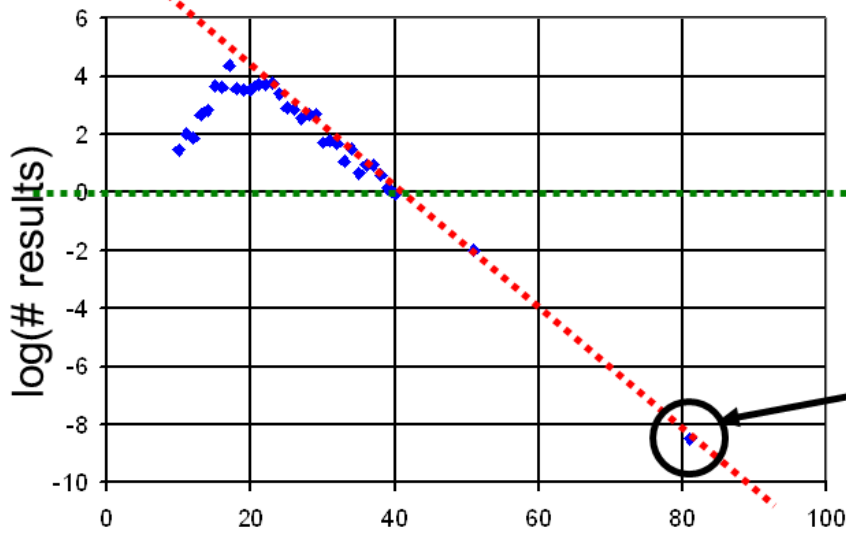
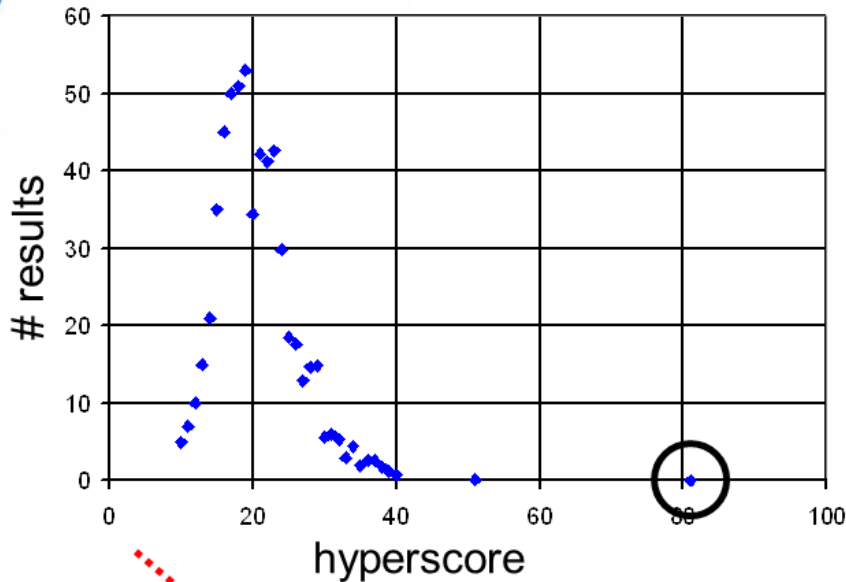
Any hyperscores greater than this are unlikely to have arisen by chance.

E-value

The E-value expresses just how unlikely a greater hyperscore is.

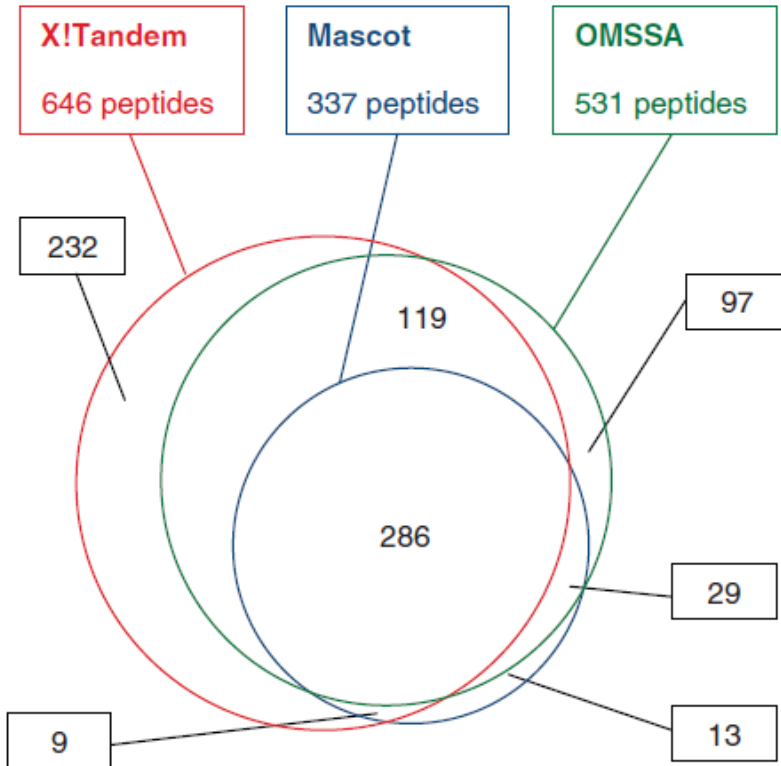
X!Tandem calculates the E-value by extrapolating the red line of the log histogram.

For the example shown, a hyperscore of 83 would occur by chance where the red line crosses 83. The log of this value — the E-value — is -8.2, as shown.



$E\text{-value} = e^{-8.2}$

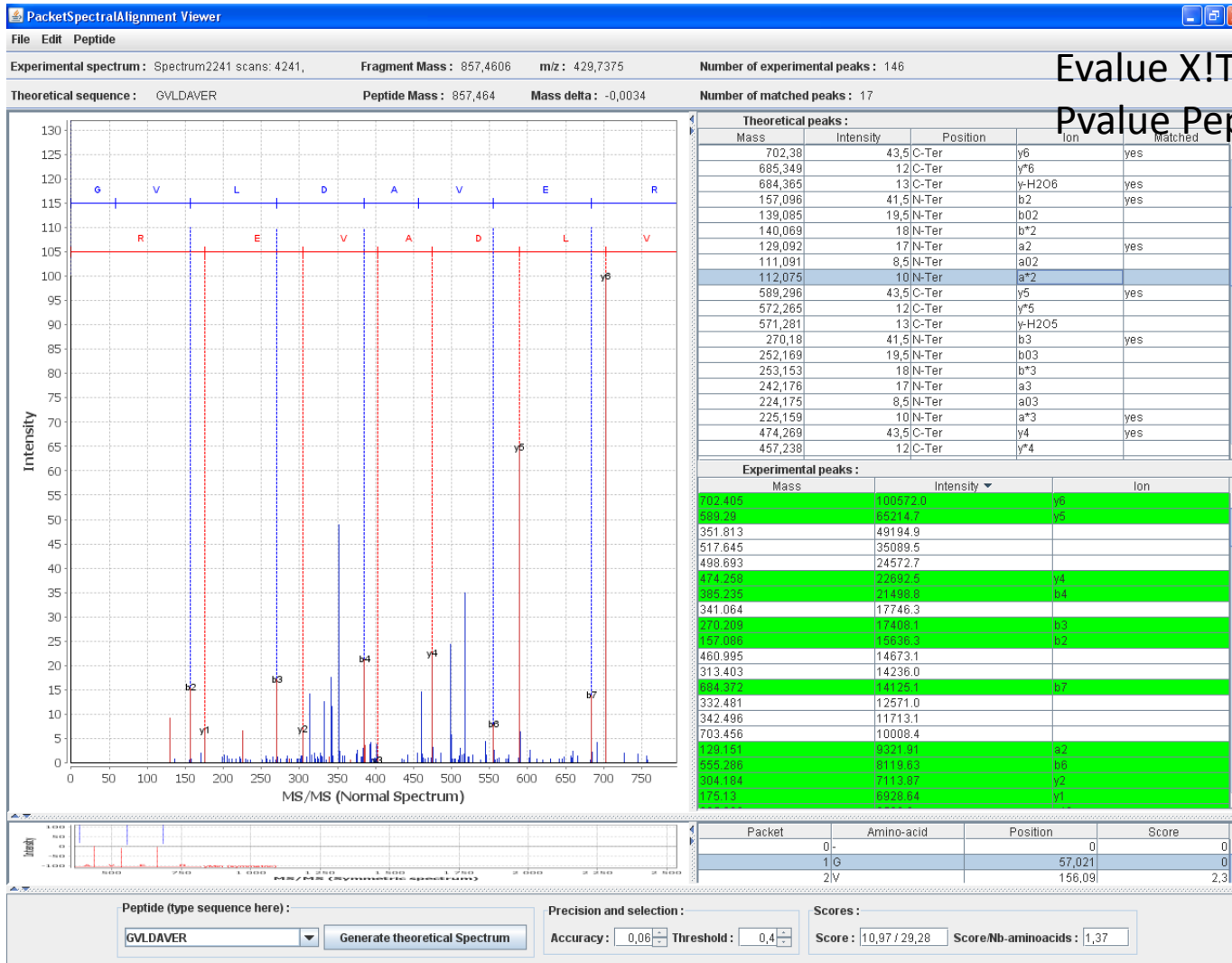
B



Les différents moteurs de recherche sont en accord sur une majorité de peptides, mais chaque moteur permet d'identifier des peptides particuliers

Jeu de données UPS – standard, FDR=1%

Vaudel, M., Burkhart, J.M., Sickmann, A., Martens, L. and Zahedi, R.P. (2011) Peptide identification quality control, *Proteomics*, **11**, 2105-2114.



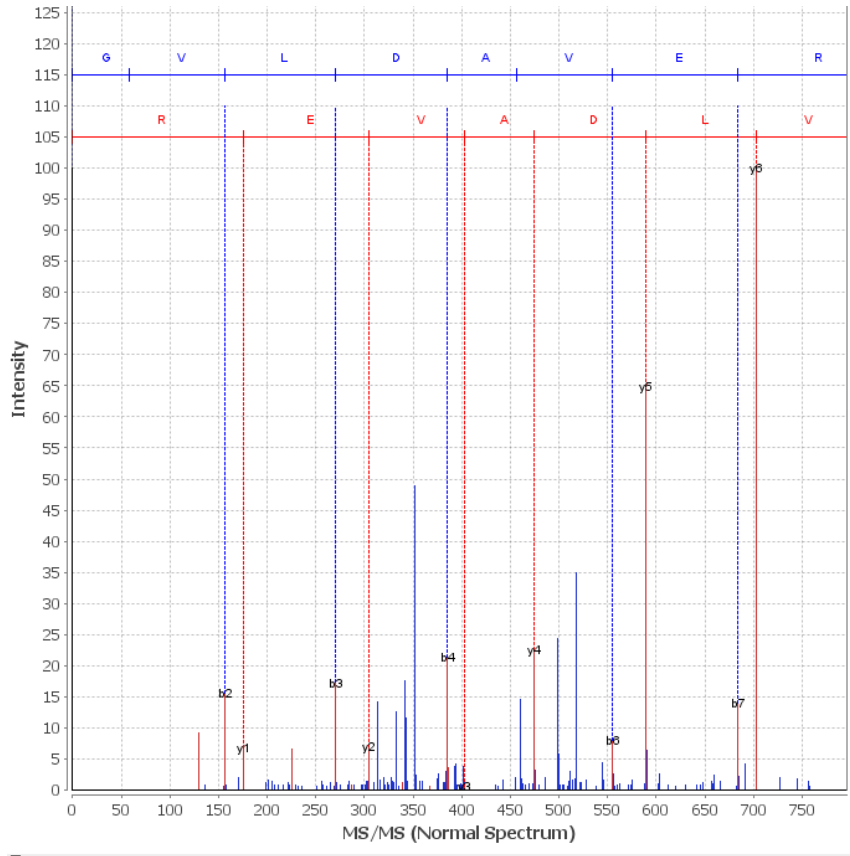
Evalue X!Tandem=4.6

Pvalue PeptideProphet =0.9387

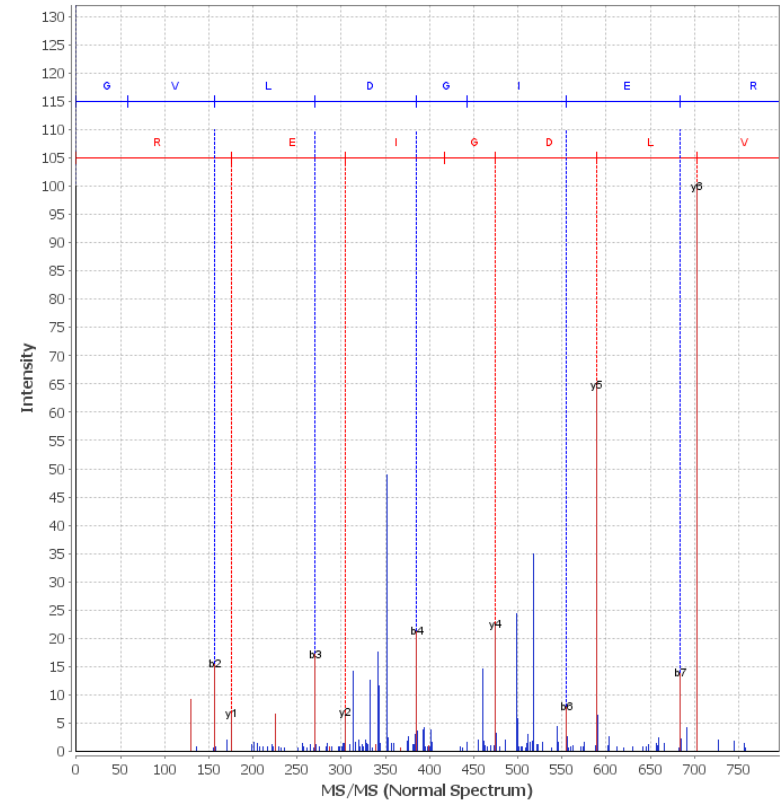
Un beau spectre, une evalue décevante ?

Séquence	Hyperscore	Delta masse	Expect
GVLDAVER	529	-0.0034	4.6
GVLDGIER	525	-0.0034	1
VGGGGGGGGGER	488	-0.9371	2.6
RLATNQR	447	-0.023	7.2
RLAGYLR	408	-0.051	19
RLETLAR	401	-0.048	22
RLAAVKGK	400	-0.118	22

Interprétation GVLD^{AV}ER



Interprétation GVLD^{GI}ER



$$\text{masse (AV)} = \text{masse(GI)}$$

Les algorithmes de comparaison de spectres :

Quels sont les types d'ions considérés ?

Comment sont définies les intensités des spectres théoriques ? *

Quel est l'algorithme de comparaison utilisé ?

Quel est le calcul de score ? *

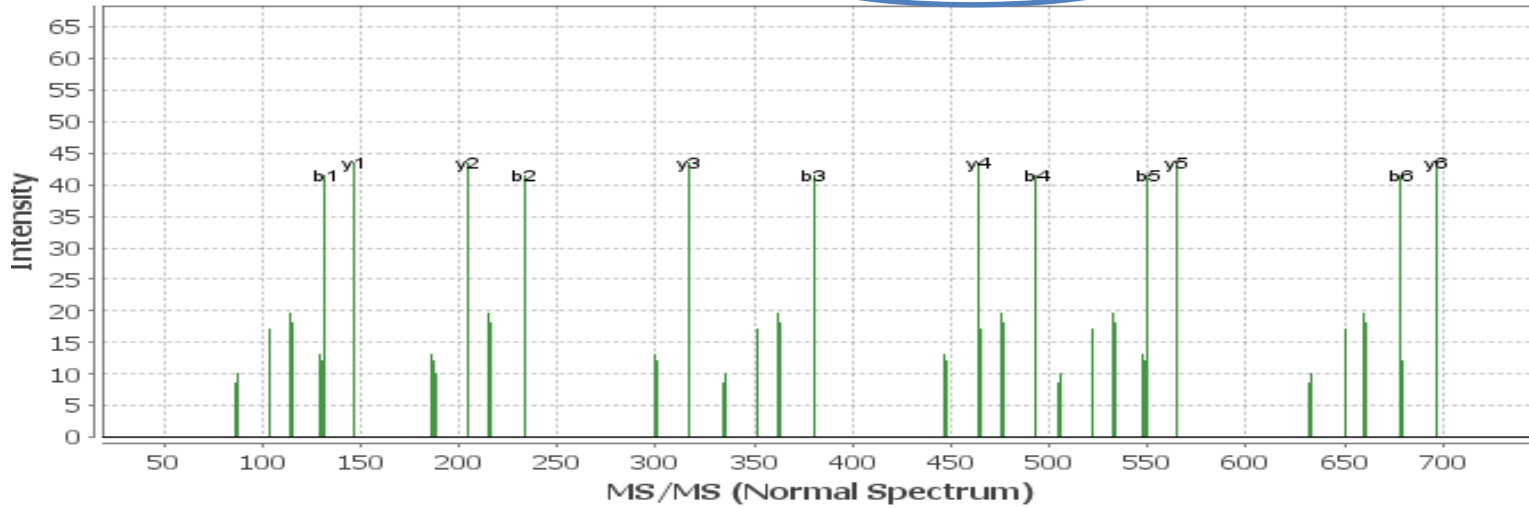
Comment les mutations , les modifications post-traductionnelles sont-elles prises en compte ?

Comment sont validées les attributions ?

theoretical sequence : **MTFLGK**

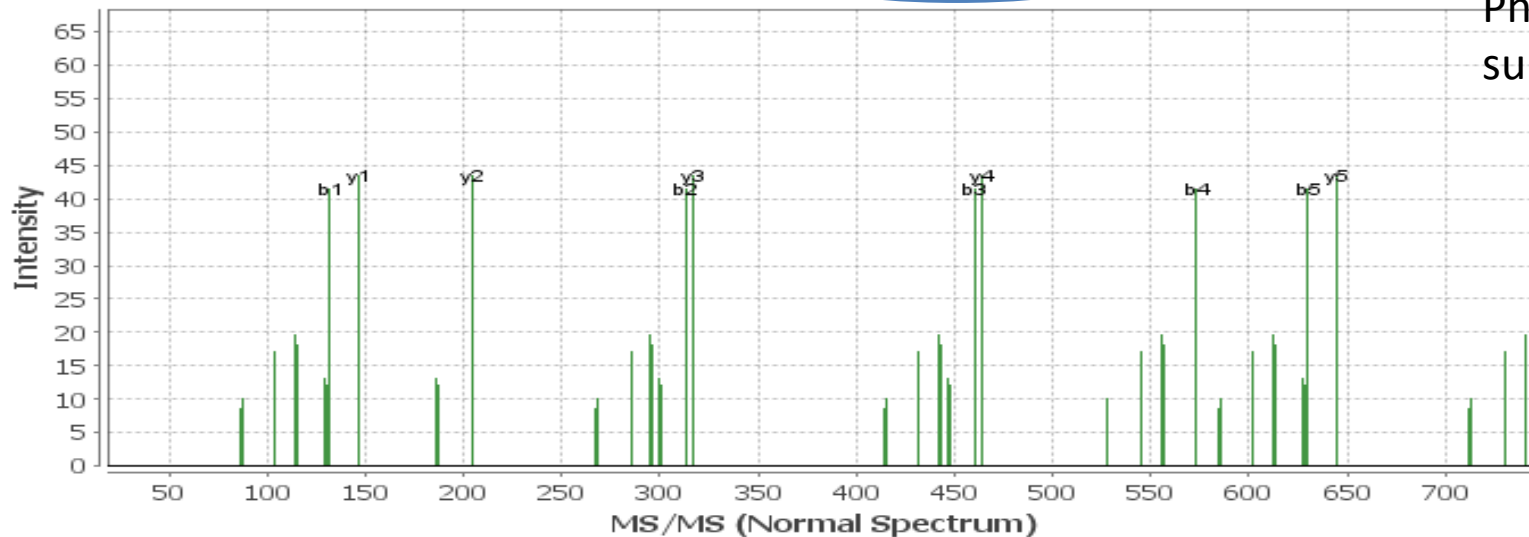
Peptide Mass : 695,371

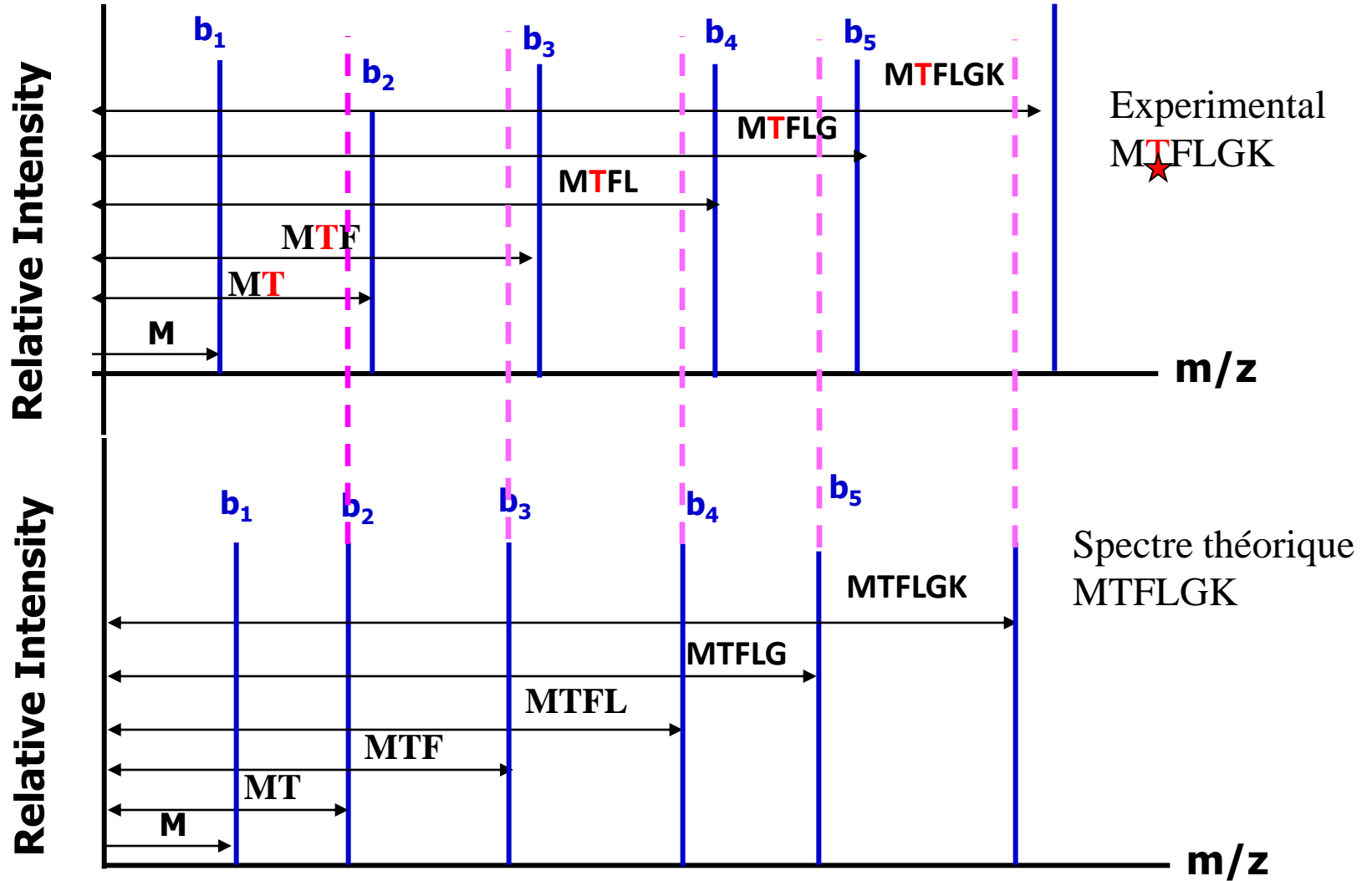
Mass delta :

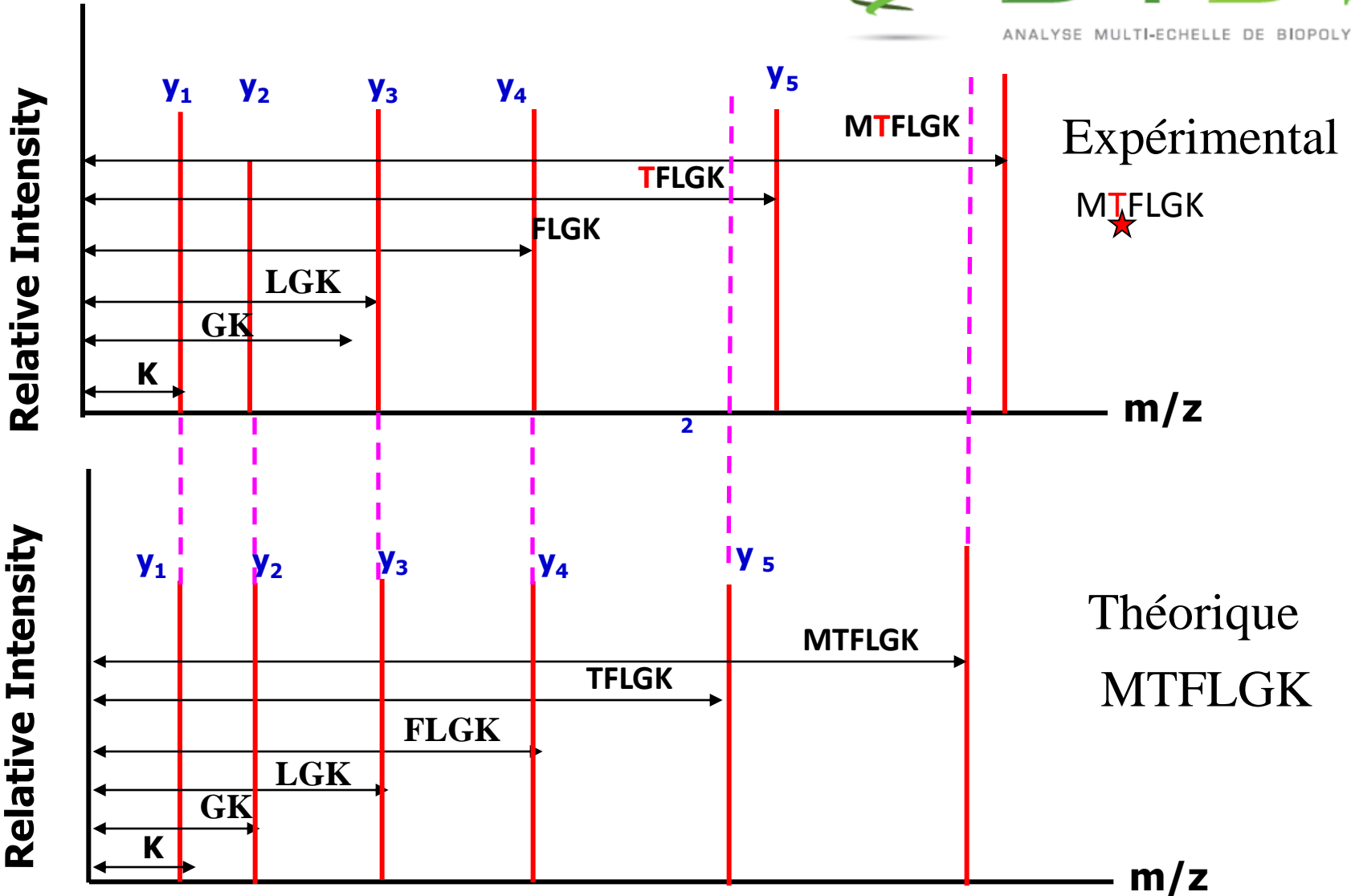

 theoretical sequence : **MT(80)FLGK**

Peptide Mass : 775,371

Mass delta :


 Phosphorylation
sur T (+80 Da)



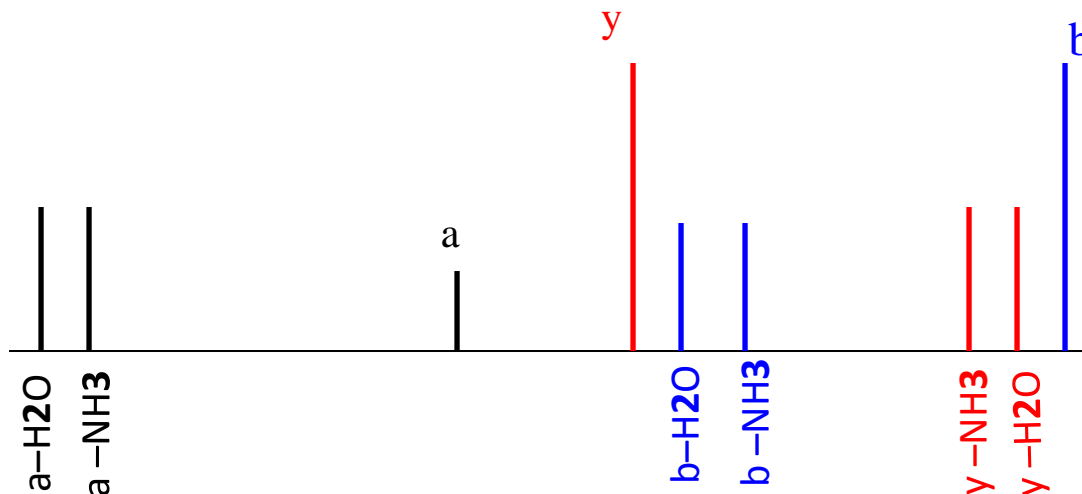


Idée originale de la thèse de Freddy Cliquet (2011): modéliser le spectre sous une forme différente à partir des relations qui lient les ions issus d'une même fragmentation

Relation entre les ions y et b: $\text{masse}(y) = \Sigma \text{masse}(aa) - \text{masse}(b) - 20$
 $\text{masse}(FLGK) = \text{masse}(MTFLGK) - \text{masse}(MT) - 20$

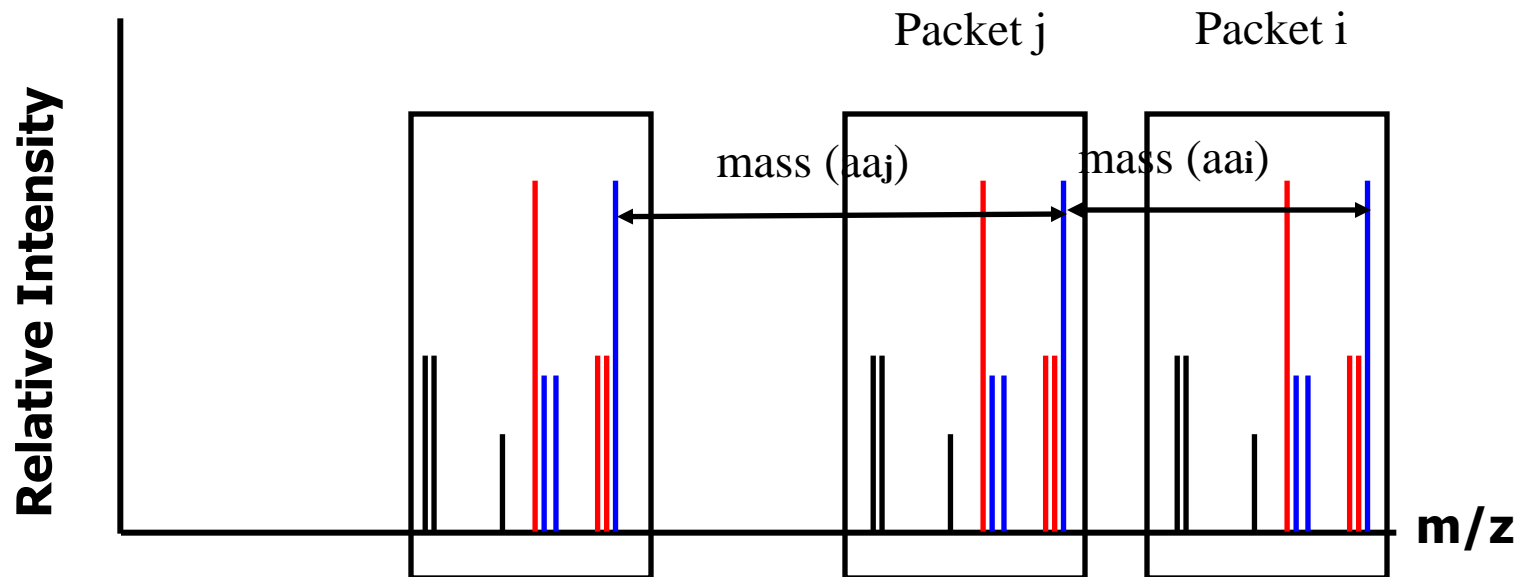
Nous pouvons regrouper les ions issus d'une même fragmentation en changeant le position des ions y

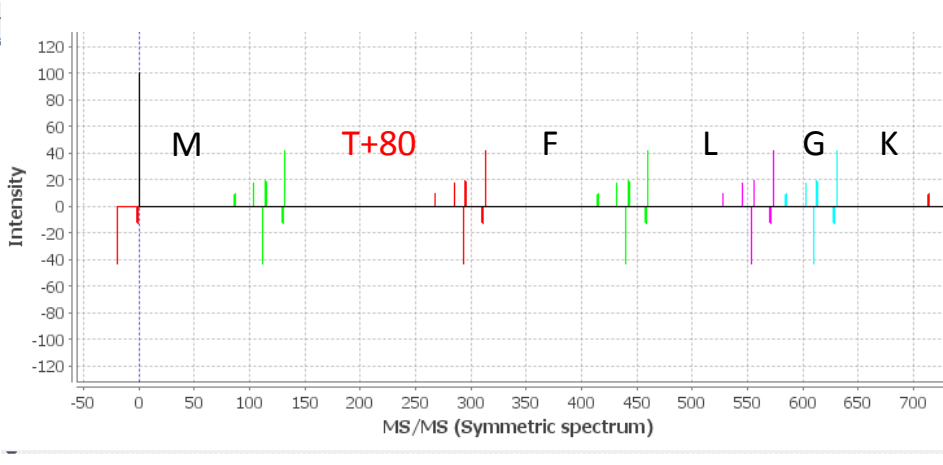
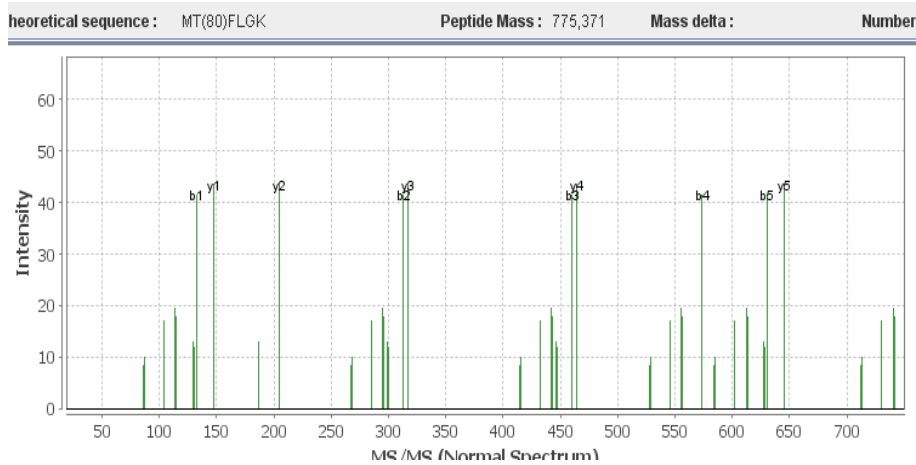
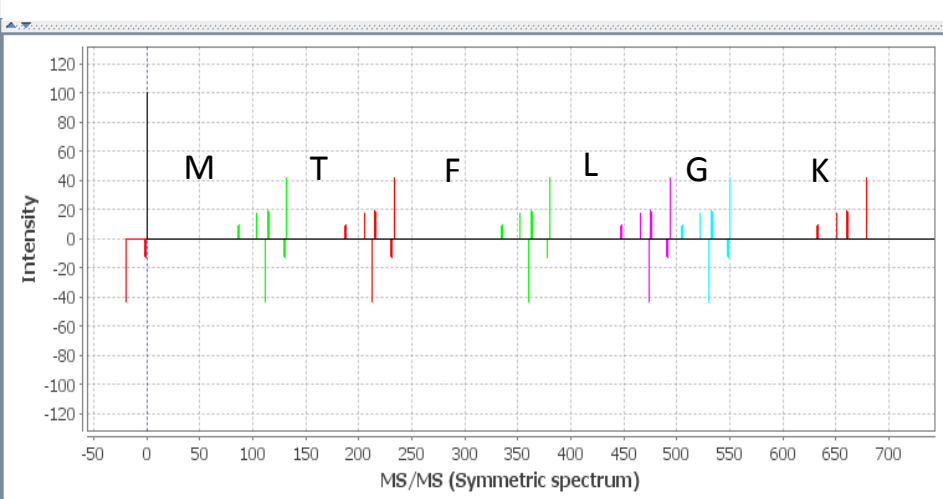
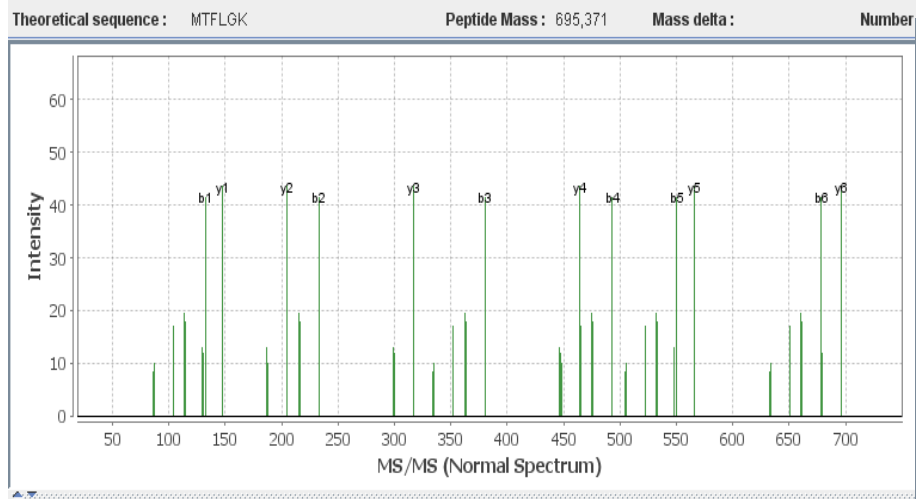
Position(y-ion) = $\Sigma \text{mass}(aa) - \text{mass}(b\text{-ion})$

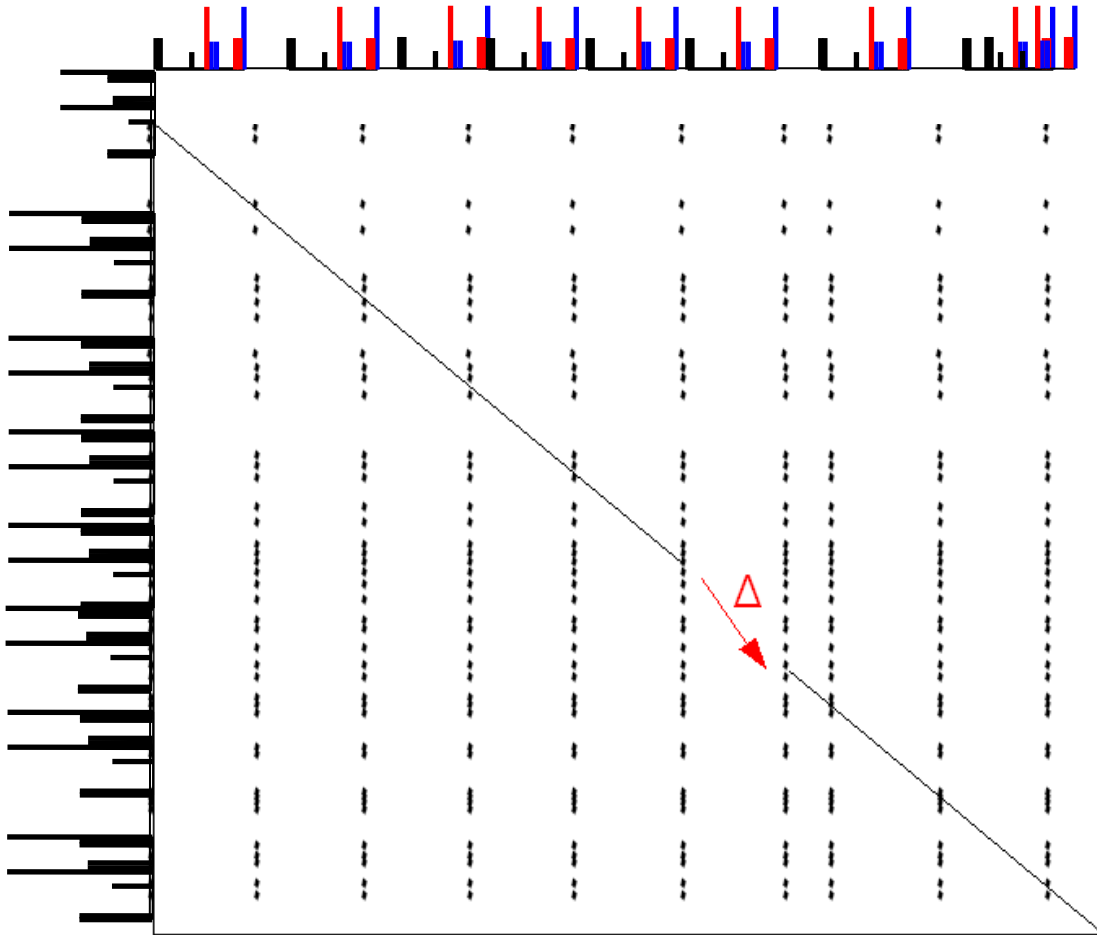


Les pics issus d'une même fragmentation sont regroupés dans un paquet

Un spectre théorique est alors représenté par une liste de paquets







L'algorithme PSA (Packet Spectral Alignment) s'appuie alors sur la programmation dynamique pour aligner les paquets du spectre théorique sur le spectre expérimental comme l'algorithme SA développé par Pevzner et al.

La transformation du spectre théorique est simple !

La transformation du spectre expérimental est beaucoup plus complexe car on ne sait pas distinguer les différents types d'ions.

PSA génère donc un 'pic complémentaire' pour chaque pic, mais attention à ne pas générer de bruit (nombre de pics * 2!)

Les ions b du spectre théorique – resp. les ions y - ne seront alignés qu'avec les pics d'origine – resp. les pics complémentaires - du spectre expérimental.



Temps d'exécution longs :

AVG ~1 minute par spectrum sur un serveur linux et l'utilisation de 7 coeurs sur une banque d'environ ~500 000 peptides

L'algorithme ne peut être appliqué que sur

- . Des spectres de bonne qualité
- . Des spectres non-interprétés par des approches plus rapides

Cliquet, F., Fertin, G., Rusu, I. and Tessier, D. (2009) Comparison of Spectra in Unsequenced Species. . In (LNBI), L.N.i.B. (ed), *4th Brazilian Symposium on Bioinformatics (BSB 2009)*. Porto Alegre, Brazil.

Cliquet, F., Fertin, G., Rusu, I. and Tessier, D. (2010) Proper alignment of MS/MS spectra from unsequenced species. *Proc. 11th International Conference on Bioinformatics and Computational Biology (BIOCOMP 2010)*. CSREA Press, 766-772.

Plan de la présentation

Le contexte

L'association spectre-peptide par approche comparative (Peptide-Spectrum Match)

L'association spectre-peptide par approche *de novo*

L'identification des protéines

La définition de « pipeline » d'analyse

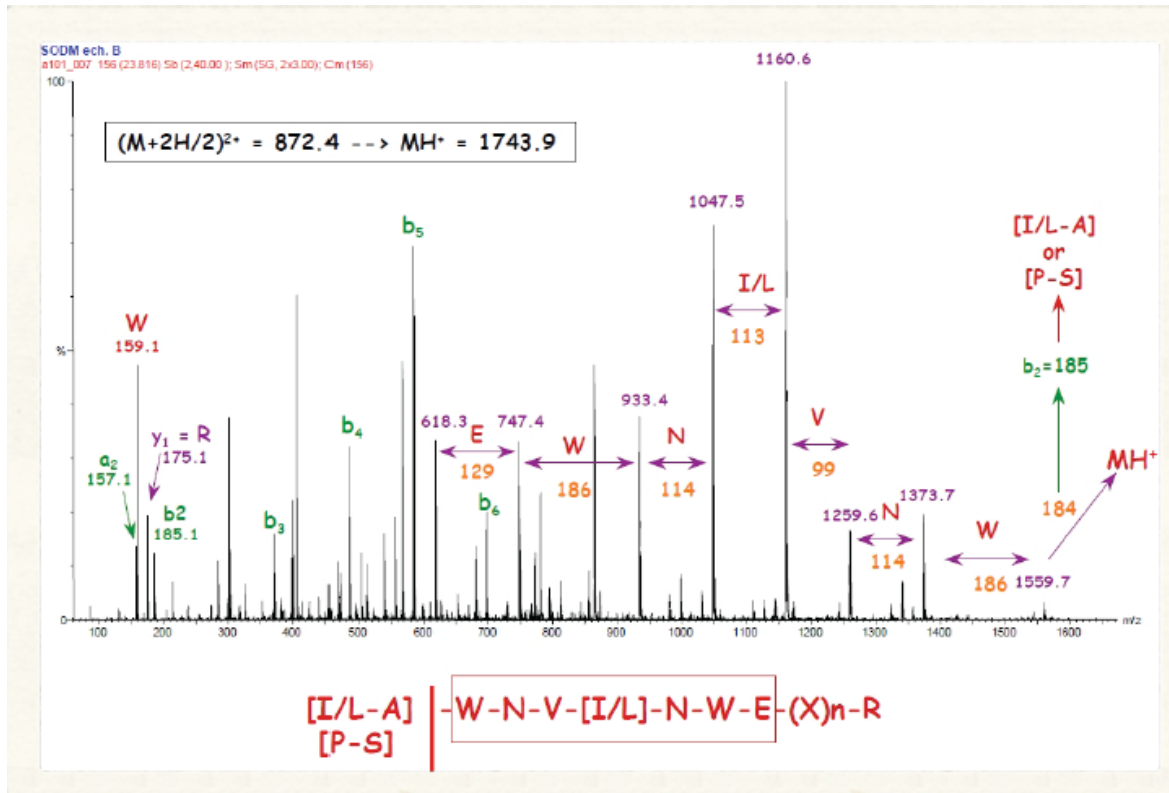
Plan de la présentation

Le contexte

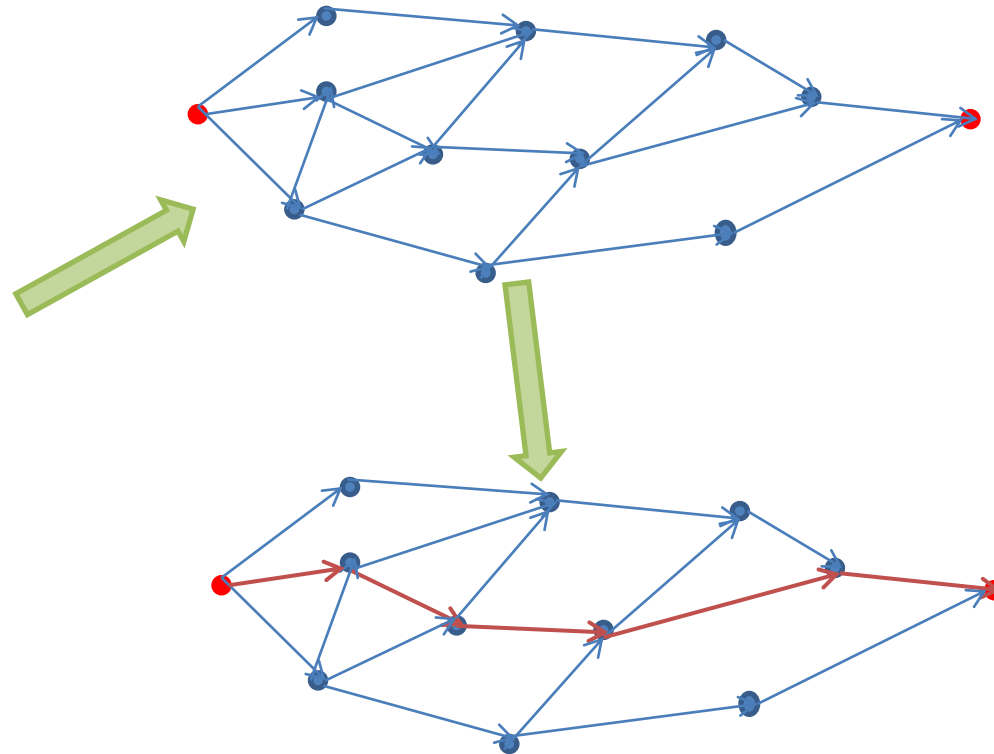
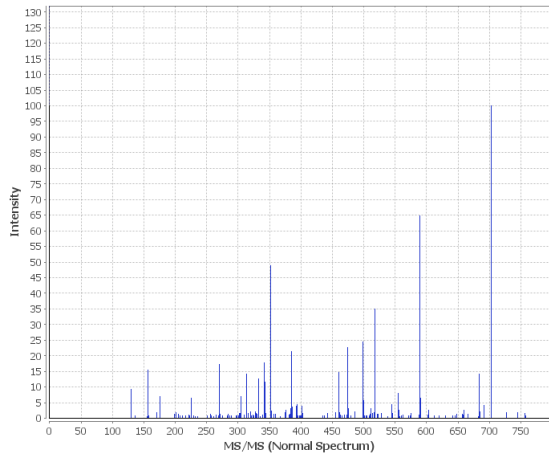
L'association spectre-peptide par approche comparative

L'association spectre-peptide par approche *de novo*

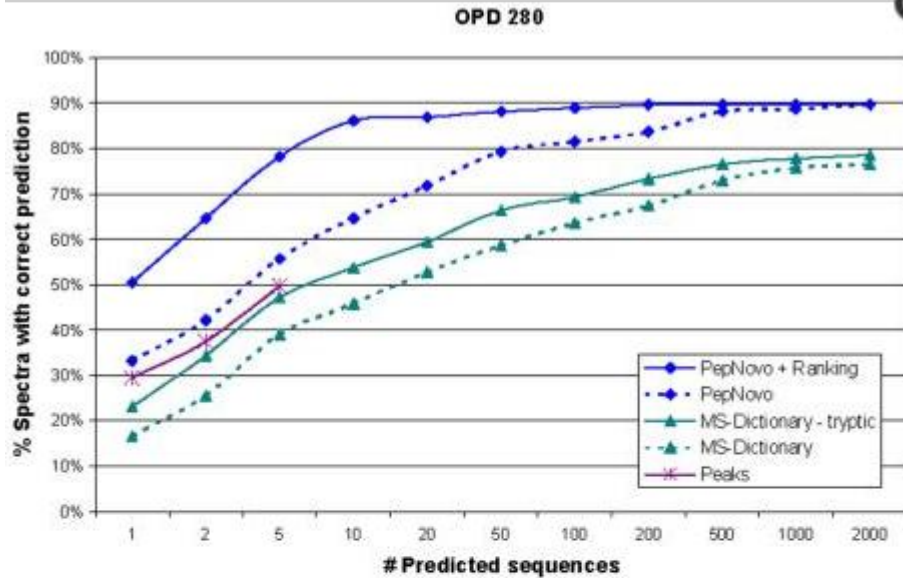
L'identification de protéines



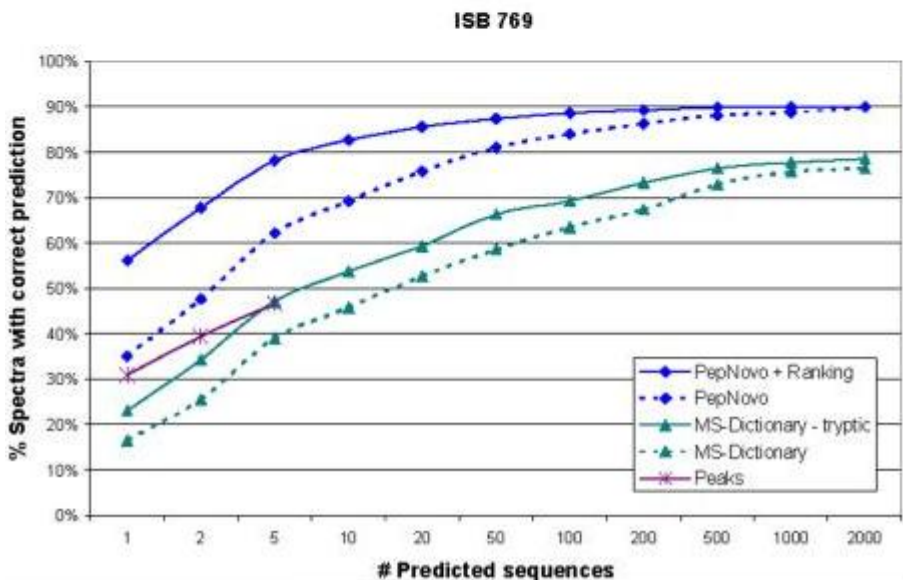
Utilisation de graphes de spectres (introduits par Bartels – 1990 -)



Recherche d'un chemin optimal



Le jeu d'évaluation correspond à une sélection de spectres : en moyenne, ils contiennent 47.7% des ions b et 51.2% des ions y



Ranking-Based Scoring Models for Peptide-Spectrum Matches. Frank, A.M. J. Proteome Research, 8:2241-2252, 2009

En résumé

Trouver une association entre un spectre et un peptide ? Facile !!!! **Oui mais est-elle correcte ?**

Erreurs:

des spectres de mauvaise qualité, mais pas toujours. Même des spectres de bonne qualité peuvent fournir des associations incorrectes si le peptide analysé n'est pas dans la banque de protéines.

Plan de la présentation

Le contexte

L'association spectre-peptide par approche comparative (Peptide-Spectrum Match)

L'association spectre-peptide par approche *de novo*

L'identification des protéines

La définition de « pipeline » d'analyse

Difficulté à assembler les peptides : Protein inference problem

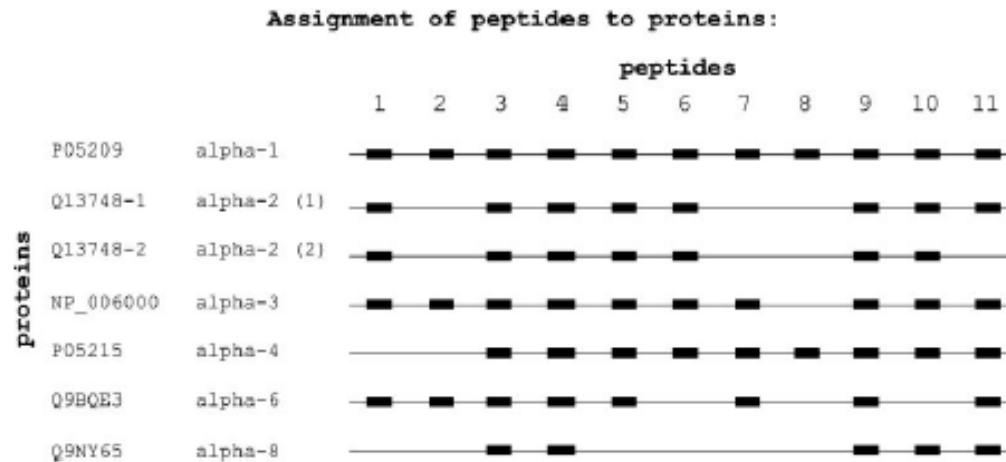


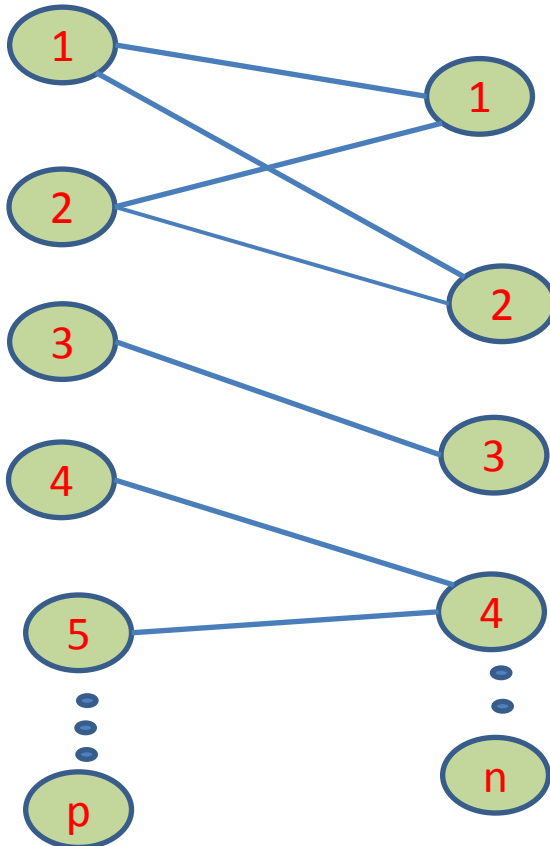
FIG. 3. **An example of a protein family.** Eleven tryptic peptides are identified that are shared between the members of the α -tubulin family. None of the proteins is identified by a peptide that is unique to it, thus making it impossible to determine which particular member(s) of the family is present in the sample.

Nesvizhskii, A.I. and Aebersold, R. (2005) Interpretation of shotgun proteomic data: the protein inference problem, *Mol Cell Proteomics*, **4**, 1419-1440.

Le problème est en général formalisé sous la forme d'un graphe biparti

Peptides

Protéines



m = nombre max de protéines

n = nombre min de protéines

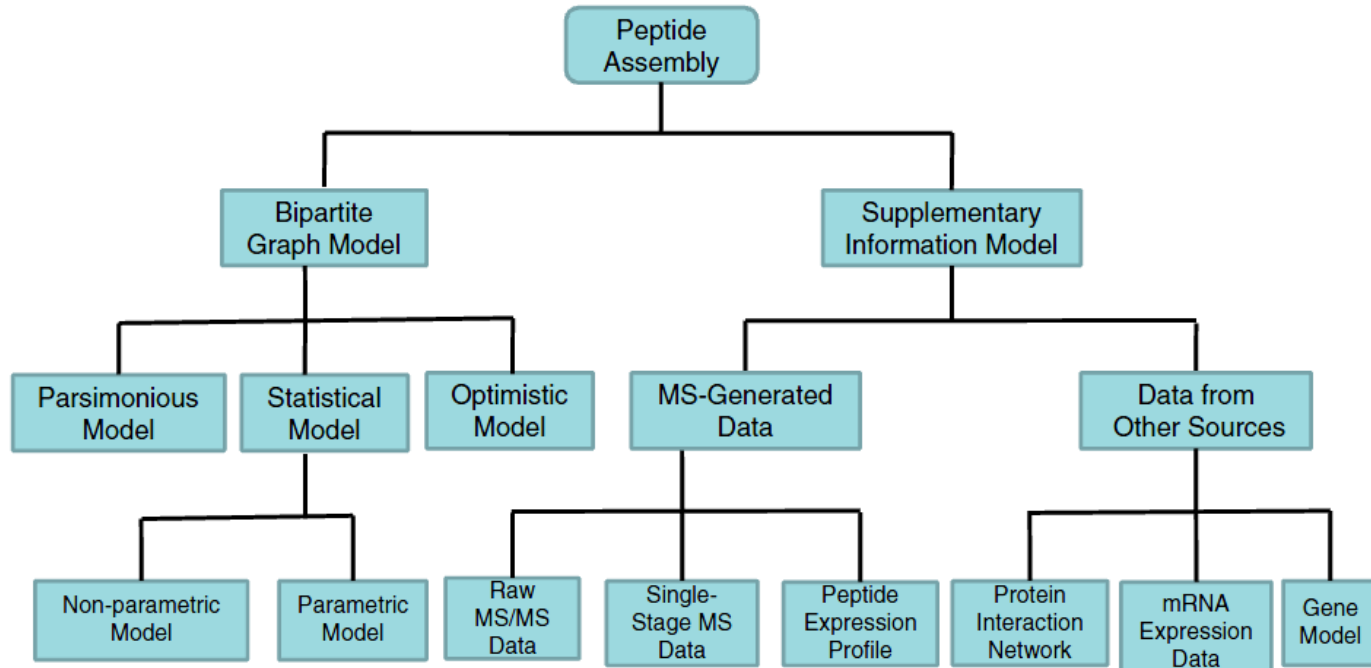
⇒ Set coverage problem (NP-complet)

⇒ Principe de parcimonie

Le problème des protéines avec 1 seul peptide

Gupta, N. and Pevzner, P.A. (2009) False discovery rates of protein identifications: a strike against the two-peptide rule, *J Proteome Res*, **8**, 4173-4181.

Les différents algorithmes

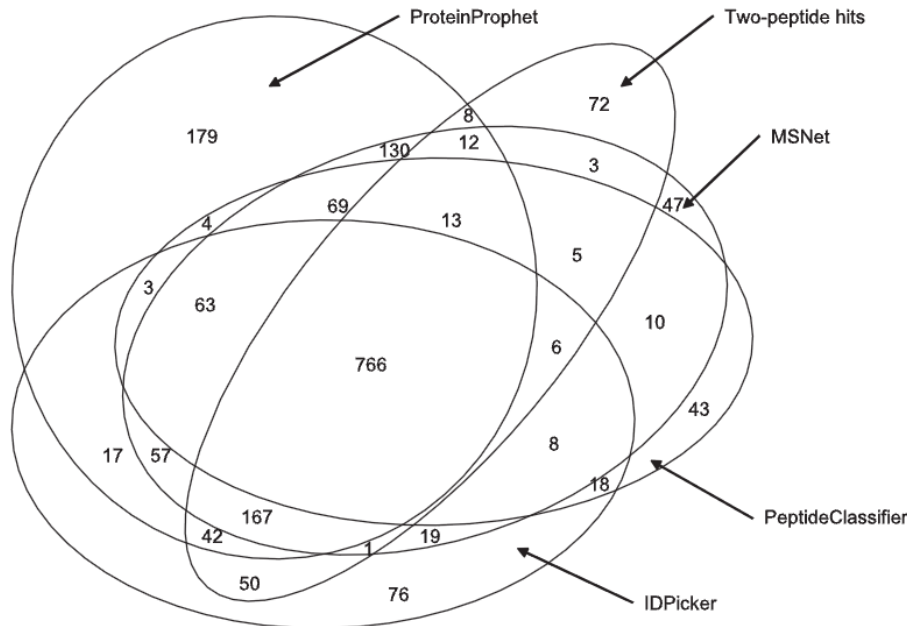


Huang, T., Wang, J., Yu, W. and He, Z. (2012) Protein inference: a review, *Brief Bioinform*, **13**, 586-614.

Table 2: The performance of five methods

Data	IDPicker	ProteinProphet	MSNet	PeptideClassifier	Two-peptide hits
Yeast.D1	1293 (1152)	1530 (1254)	1376 (1319)	1008 (977)	1112 (978)
Yeast.D2	154 (116)	322 (162)	291 (166)	232 (145)	317 (133)

Two protein samples undergo MS/MS analysis to generate ten lists of proteins identified by five peptide assembly algorithms. The figures in brackets indicate the number of proteins present in the reference set.



Huang, T., Wang, J., Yu, W. and He, Z. (2012) Protein inference: a review, *Brief Bioinform*, **13**, 586-614.

PSM	UPS ₂	Yeast YPD	iPRGog(Red)	iPRGog(Yellow)
Total MS/MS spectra observed	74,602	240,781	69,416	70,970
SEQUEST	32,651 (87)	57,955 (268)	9,524 (98)	9,492 (83)
X!Tandem	27,264 (210)	74,244 (332)	15,147 (117)	15,366 (112)
MyriMatch	26,262 (79)	41,179 (106)	9,706 (88)	9,134 (46)
InsPecT	25,618 (64)	69,341 (414)	12,691 (202)	13,295 (216)
Union	40,829 (434)	95,315 (1053)	21,764 (505)	21,684 (455)
MSblender	39,273 (336)	99,814 (1011)	23,580 (153)	23,717 (177)
1 engine	4,043 (190)	10,441 (100)	2,138 (38)	2,073 (52)
2 engines	7,389 (89)	16,861 (546)	3,768 (76)	3,878 (74)
3 engines	5,560 (35)	32,111 (203)	6,820 (24)	6,816 (21)
4 engines	22,202 (24)	38,257 (18)	10,830 (3)	10,826 (4)
Protein	UPS ₂	Yeast YPD	iPRGog(Red)	iPRGog(Yellow)
Total proteins	48	6,698	4,417	4,417
SEQUEST	38	1,391	757	749
X!Tandem	38	1,459	870	847
MyriMatch	36	1,241	722	657
InsPecT	29	1,527	877	902
Union	44	1,873	999	1,024
MSblender	42	1,911	1,185	1,147

Kwon, T., Choi, H., Vogel, C., Nesvizhskii, A.I. and Marcotte, E.M. (2011) MSblender: A probabilistic approach for integrating peptide identifications from multiple database search engines, *J Proteome Res*, **10**, 2949-2958.

← Résultat à 0.5% de FDR

Plan de la présentation

Le contexte

L'association spectre-peptide par approche comparative (Peptide-Spectrum Match)

L'association spectre-peptide par approche *de novo*

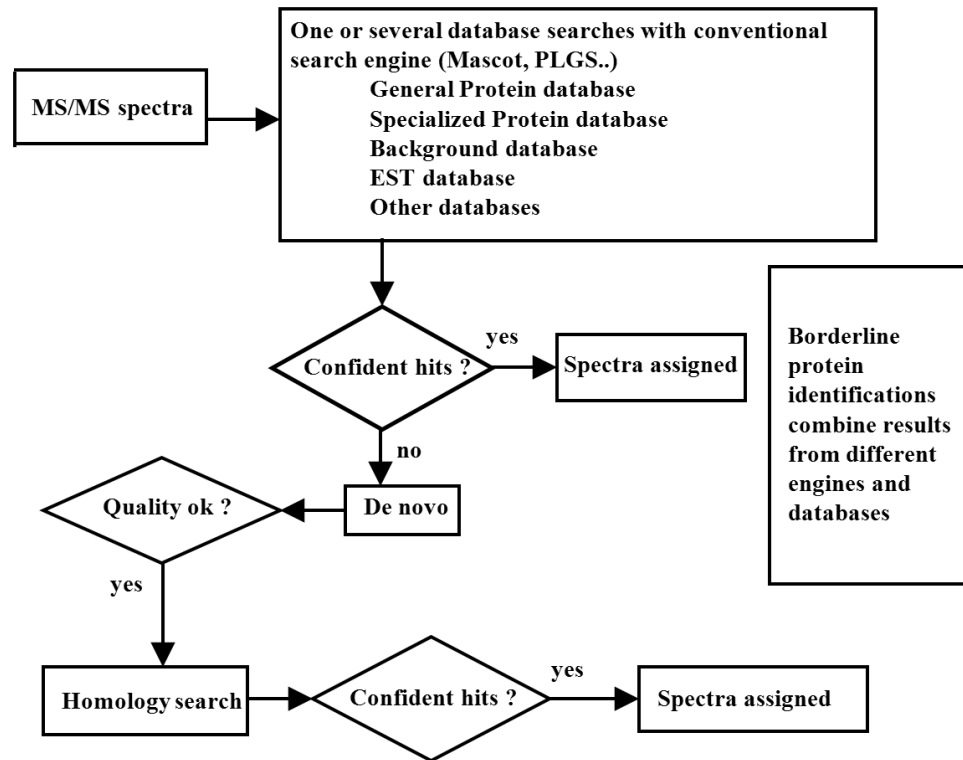
L'identification des protéines

La définition de « pipeline » d'analyse

Développement de processus d'analyse adaptés aux questions de recherche



Développement du logiciel OVNlp



Remerciements

La composante de spectrométrie de masse de la plateforme BIBS, en particulier Hélène Rogniaux, Audrey Geatron

Freddy Pliquet (doctorant), Guillaume Merceron (master 2)

L'équipe PomBi du LINA, en particulier Guillaume Fertin, Irena Rusu

Polette Larré (unité INRA-BIA)

Le groupe bioinformatique de BIBS : Virginie Lollier, Stéphane Bansard, Marc Vasseur