

# Outils bio-informatiques pour la protéomique- Comparaison TheGPM/Xcalibur Sequest- Gelprint: génération de gels 2D in-silico

Danielle Moinier, Hiroyuki Ogata et Stéphane Audic  
Institut de Biologie Structurale et Microbiologie  
Information Génomique et Structurale  
31 Chemin Joseph Aiguier  
13402 Marseille cedex 20

# De quoi parlerons nous?

- Présentation de TheGPM
- Mimivirus: un virus géant vivant dans les amibes
- Comparaison TheGPM/Sequest pour l'identification des protéines par spectrométrie de masse MS/MS


The GPM - Mozilla

File Edit View Go Bookmarks Tools Window Help

Back Forward Reload Stop <http://www.thegpm.org/> Search Print

Home Bookmarks People Google annot Tools Debug Actualite bioinfo EMBOSS: The Applic...

---



**The Global Proteome Machine Organization**  
**www.thegpm.org**

---

**Use the GPM**

- [Go](#)
- [FAQ](#)

**Single proteomes**

- [human](#)
- [T. brucei](#)

**Projects**

- [GPM](#)
- [X! Tandem](#)
- [XML](#)
- [LiveCD](#)
- [Quartz](#)
- [TandemCom](#)

**Download**

- [ftp site](#)

**Contact us**

- [email](#)

**Welcome!**

The Global Proteome Machine Organization was set up so that scientists involved in proteomics using tandem mass spectrometry could use that data to analyze proteomes. The projects supported by the GPMO have been selected to improve the quality of analysis, make the results portable and to provide a common platform for testing and validating proteomics results.

**GPM News**

**Two new Projects available: LiveCD and Quartz (2004/4/15)**

The GPMO has added two new projects, LiveCD and Quartz to the site. LiveCD, a project from the University of Michigan NCRR Center for Proteomics, provides a simple method to install a Linux-based version of X! TANDEM and the GPM on a large number of computers for instructional and demonstration purposes. It also includes some software allowing the use of X! TANDEM on clusters of computers running LiveCD.

Quartz is a GPMO staff project. It is a set of annotated spectrum collections, meant to be used for bioinformatics research. The current collections contain > 2000 MS/MS spectra, along with XML-formated annotation files.

**X! TANDEM and the GPM release updates (2004/4/10)**

New releases of both X! TANDEM and the GPM were released today. This is a maintenance release, including fixes for small problems observed with previous versions. The collections of sequences for the GPM have been updated to include the latest sequence releases from ENSEMBL (1/4/2004).

**Probit model published (2004/3/1)**


the GPM site - Mozilla

File Edit View Go Bookmarks Tools Window Help

Back Forward Reload Stop  Search Print

Home Bookmarks People Google annot Tools Debug Actualite bioinfo EMBOSS: The Applic...

### The Global Proteome Machine



advanced [page](#)  
view saved [xml data](#)

what is the [GPM](#)  
powered by [tandem](#)  
send us [email](#)

mirror sites  
north america

[| h003](#) | [h066](#) | [h112](#) |  
[| h319](#) | [h451](#) | [h777](#) |  
[| h874](#) |

**spectra:** DTA, PKL or Matrix Science format only

**taxon:** Select the appropriate species.  
  
Mus musculus  
Rattus norvegicus  
D. melanogaster  
D. rerio  
C. elegans  
A. thalania  
S. cerevisiae

**fragment  $\delta m$ :**   Da  ppm

**output:** log(e) <

**modifications:** format = "m<sub>1</sub>@X<sub>1</sub>,m<sub>2</sub>@X<sub>2</sub>,..." where:  
m<sub>i</sub> = mod. mass (Da) & X<sub>i</sub> = residue.

complete:

potential:

potential:  
(refine)

motif:  
(refine)

**mutations:**  yes  no

**method:** Select device & parent  $\delta m$ .

Quad-TOF (100 ppm)  
Quad-TOF (0.5 Da)  
Ion Trap (4 Da)

Par défaut GPM est configuré pour utiliser les protéomes suivants:

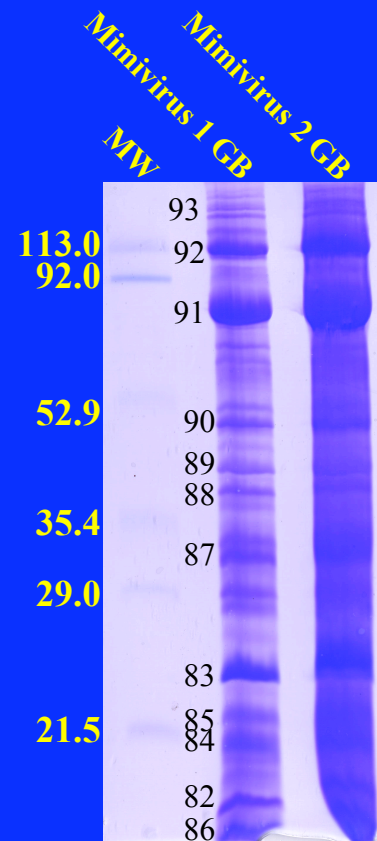
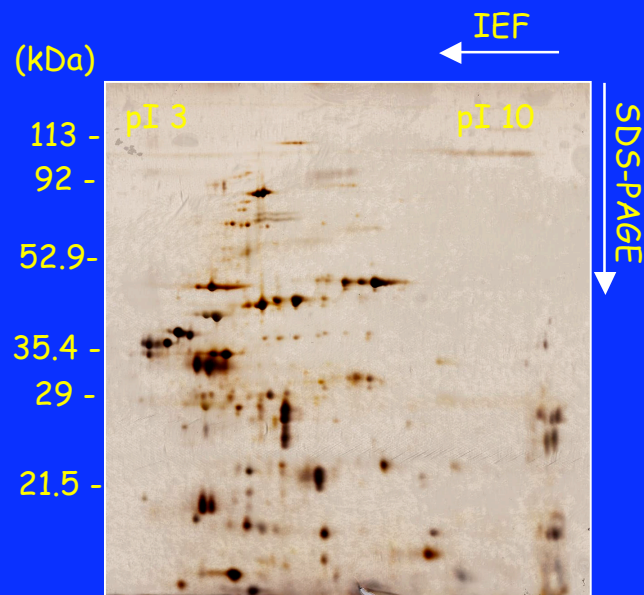
1. Human (34d NCBI 34, 10 May 2004)
2. Mouse (32b NCBI m32, 1 Apr. 2004)
3. Rat (RGSC 3.1b, 9 Feb. 2004)
4. Zebra fish (3b WTSI Zv3, 1 Apr. 2004)
5. Fruit fly (BGDP 3.1, 2 Jul. 2003)
6. *C. elegans* (116a WS 116, 1 Apr. 2004)
7. Yeast (SCD, 15 Nov 2003)
8. *A. thaliana* (ATH1, v. 5.0, Jan 29, 2004)

Et fonctionne sur une machine standard, sous Windows ou Linux:  
(nous avons teste linux)

1. 2.4 GHz Pentium IV processors;
2. 256 MB memory
3. 40 GB, 7200 rpm EIDE hard drive
4. Apache HTTP server, v. 2
5. ActiveState Perl

Notre expérience est basée sur des spots issu de gels:

- 2D coloré au nitrate d'argent
- 1 D coloré au bleu de Coomassie



Spots analysés sur une trappe à ions LCQ Deca XP plus


the GPM site - Mozilla

File Edit View Go Bookmarks Tools Window Help

Back Forward Reload Stop  Search Print

Home Bookmarks People Google annot Tools Debug Actualite bioinfo EMBOSS: The Applic...

### The Global Proteome Machine



advanced [page](#)  
view saved [xml data](#)

what is the [GPM](#)  
powered by [TANDEM](#)  
send us [email](#)

mirror sites  
north america

[| h003 | h066 | h112 |](#)  
[| h319 | h451 | h777 |](#)  
[| h874 |](#)

**spectra:** DTA, PKL or Matrix Science format only

**taxon:** Select the appropriate species.

- D. melanogaster
- D. rerio
- C. elegans
- A. thalania
- S. cerevisiae
- Mimivirus prots
- Swissprot May 2004
- trEMBL May 2004

**fragment  $\delta m$ :**   Da  ppm

**output:** log(e) <

**modifications:** format = "m<sub>1</sub>@X<sub>1</sub>,m<sub>2</sub>@X<sub>2</sub>,..." where:  
m<sub>i</sub> = mod. mass (Da) & X<sub>i</sub> = residue.

complete:

potential:

refinement:

**mutations:**  yes  no

**method:** Select device & parent  $\delta m$ .

- FTICR (10 ppm)
- Quad-TOF (100 ppm)
- Quad-TOF (0.5 Da)
- Ion Trap (4 Da)




the GPM site - Mozilla

File Edit View Go Bookmarks Tools Window Help

Back Forward Reload Stop [http://10.234.244.194/tandem/thegpm\\_tandem\\_a.html](http://10.234.244.194/tandem/thegpm_tandem_a.html) Search Print

Home Bookmarks People Google annot Tools Debug Actualite bioinfo EMBOSS: The Applic...

---



### The Global Proteome Machine

advanced search page

[simple PAGE](#)  
view saved [xml data](#)

---

what is the [GPM](#)  
powered by [TANDEM](#)  
send us [email](#)

---

mirror sites  
north america

[| h003 | h066 | h112 |](#)  
[| h319 | h451 | h777 |](#)  
[| h874 |](#)

**spectra:** DTA, PKL or Matrix Science format only

**taxon:** [Mimivirus prots](#)

**measurement errors**

Fragment mass error:

Parent mass error: +  -

Isotope error:  yes  no

**output**

Max expect <:

**residue modifications**

modifications:   
m1@X,m2@Y, etc.

potential modifications:   
m1@X,m2@Y, etc.

protein N-terminus:  Da

protein C-terminus:  Da

**cleavage specification**

cleavage C-terminal change:  Da

the GPM site - Mozilla

File Edit View Go Bookmarks Tools Window Help

Back Forward Reload Stop [http://10.234.244.194/tandem/therpm\\_tandem\\_a.html](http://10.234.244.194/tandem/therpm_tandem_a.html) Search Print

Home Bookmarks People Google annot Tools Debug Actualite bioinfo EMBOSS: The Applic...

cleavage site:   
trypsin = [KR]{}P

missed sites:

**model refinement**

refine model:  yes  no

point mutations:  yes  no

potential modifications:   
m1@X,m2@Y, etc.

use these modifications throughout:  yes  no

unanticipated cleaves ((X)[X]):  yes  no

potential N-terminus modifications: 

potential C-terminus modifications:

valid expectation: <

**spectrum synthesis**

spectrum synthesis:  yes  no

**spectrum conditioning**

Noise suppression:  yes  no

Minimum parent M+H:  Da

Minimum fragment m/z:

Total peaks:


Minimum peaks:

GPM - Models from 'bc91.mgf' - Mozilla

File Edit View Go Bookmarks Tools Window Help

Back Forward Reload Stop <http://10.234.244.194/theGPM-cgi/plist.pl?path=/tandem/arcl> Search Print

Home Bookmarks People Google annot Tools Debug Actualite bioinfo EMBOSS: The Applic...

 Models from 'bc91.mgf'

[PERFORMANCE](#) | [parameters](#) | [details](#) | [XML](#) |

[get annotation](#)

| rank | log(e) | accession   |
|------|--------|---|
| #1.  | -125.2 | Contig31_171[63834-62038](REVERSESENSE) <a href="#">homologues</a> <a href="#">protein</a>                  |
| #2.  | -73.8  | Contig31_180[46236-44239](REVERSESENSE) <a href="#">homologues</a> <a href="#">protein</a>                  |
| #3.  | -7.0   | Contig31_115[196187-194196](REVERSESENSE) <a href="#">homologues</a> <a href="#">protein</a>                |
| #4.  | -3.9   | Contig33_76[165386-167494] <a href="#">homologues</a> <a href="#">protein</a>                               |
| #5.  | -2.7   | gi 999627  <a href="#">homologues</a> <a href="#">protein</a><br>Chain B, Porcine E-Trypsin (E.C.3.4.21.4). |
| #6.  | 1.7    | Contig32_164[18449-13218](REVERSESENSE) <a href="#">homologues</a> <a href="#">protein</a>                  |

*plist.pl, v. 2004.03.01*


Major capsid protein

GPM - protein model: Contig31\_171[63834-62038](REVERSESENSE) - Mozilla

File Edit View Go Bookmarks Tools Window Help

Back Forward Reload Stop [http://10.234.244.194/the\\_gpm.cgi/protein.pl?path=/tandem/ε](http://10.234.244.194/the_gpm.cgi/protein.pl?path=/tandem/ε) Search Print

Home Bookmarks People Google annot Tools Debug Actualite bioinfo EMBOSS: The Applic...



protein model: Contig31\_171[63834-62038](REVERSESENSE)

[MODEL](#) | [homologues](#) | [details](#) | [XML](#) |

[ensembl](#)

**log(e) = -125.2** Contig31\_171[63834-62038](REVERSESENSE)

```

1 VTSSTTGRLLWVETKFLKLVQSSPLGKLRGNRVFKVVYRRHTNFAVESIEQFFGGNLGFGK 60
61 KSSAEINRSGDLITQVFLKVTLPVRYCGDFTNFHVEFAWVRNIGHAIVEETELEIGGS 120
121 PIDKHYGDWLQIWQDVSSSKDHEKGLAKMLGDVPELTSISTLSWDVPDNTVLKPSYTLV 180
181 PLQFYFNRNNGLALPLIALQYHQVRIYVKFRQADQCYIASDAFKSGCGNLQLDDVSLYVN 240
241 YVFLDTEERRRFAQVSHEYLIEQLQFTGEESAGSSNSAKYKLNFNHPVKAIYWVTKLGNY 300
301 QGGKFMTYDPVCWENARENAAKLLLLAQYDLDDWGYFQEPGGYECEGNDGRSYVGDGCVQ 360
361 YTAVDPSNPSEEPSYIFNDTTAEAFDGSLLIGKLAPCVPLLRNKDLDLKDKEVEGIIRI 420
421 HTDFENDRMKYPEVEKITERNDLTLHDLSPISKYDVDNRVDYIKKFDVTWQHNNFGLLI 480
481 DSGGNPHEAELQLNGQPRQSKRGGIWYDTVNPVHHTKSPRDGVNVFSFALNPBEHQPS 540
541 CTCNFSRIDTAQLNLWFQHTNHNKFAVDFADNDNKVLI FAVNYNVLRLMLSGMAGLAYSN 599

```

| spectrum | log(e) | m+h      | delta | z | sequence                          |
|----------|--------|----------|-------|---|-----------------------------------|
| 129.1    | -8.2   | 2299.116 | 1.621 | 2 | vyrr40HTNFAVESIEQFFGGNLGFGK60kssa |
| 220.1    | -6.1   | 2299.116 | 0.696 | 2 | vyrr40HTNFAVESIEQFFGGNLGFGK60kssa |
| 128.1    | -3.7   | 2300.116 | 1.681 | 3 | vyrr40HTNFAVESIEQFFGGNLGFGK60kssa |
| 210.1    | -2.9   | 2300.116 | 0.677 | 3 | vyrr40HTNFAVESIEQFFGGNLGFGK60kssa |
| 57.1     | -5.6   | 1220.690 | 0.467 | 2 | einr69SGDLITQVFLK79vttp           |
| 1087.1   | -4.7   | 1221.690 | 0.902 | 2 | einr69SGDLITQVFLK79vttp           |



peptide model: 129.1.1 of Contig31\_171[63834-62038](REVERSESENSE)

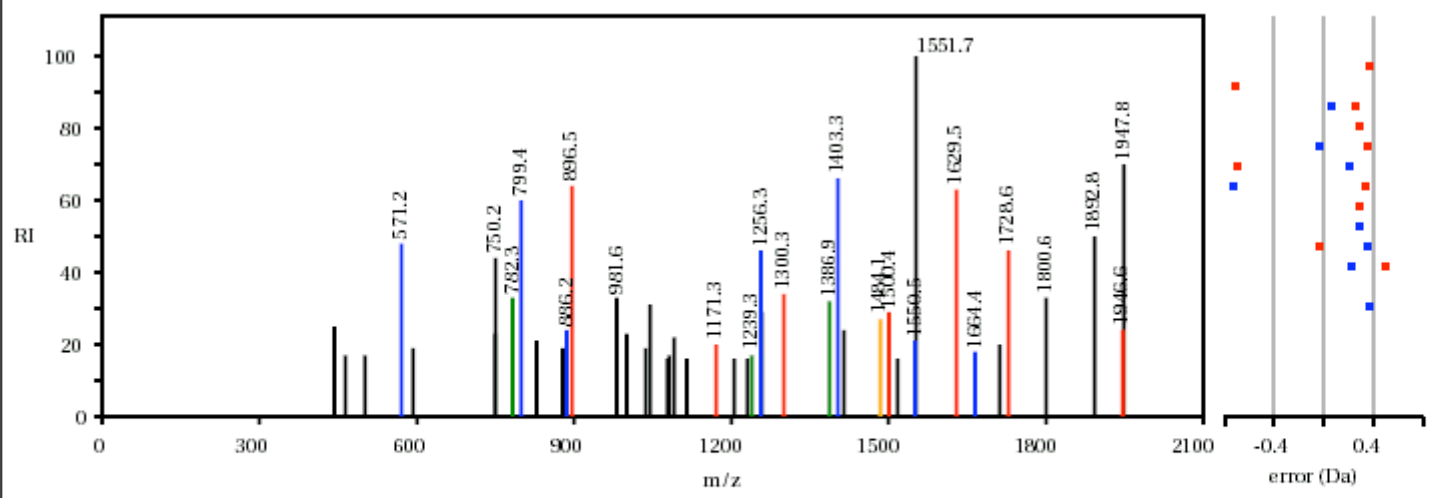
[MODEL](#) | [PROTEIN](#) | [HOMOLOGUES](#) | [XML](#) |

[ADOBE SVG PLUGIN](#) required to view spectrum

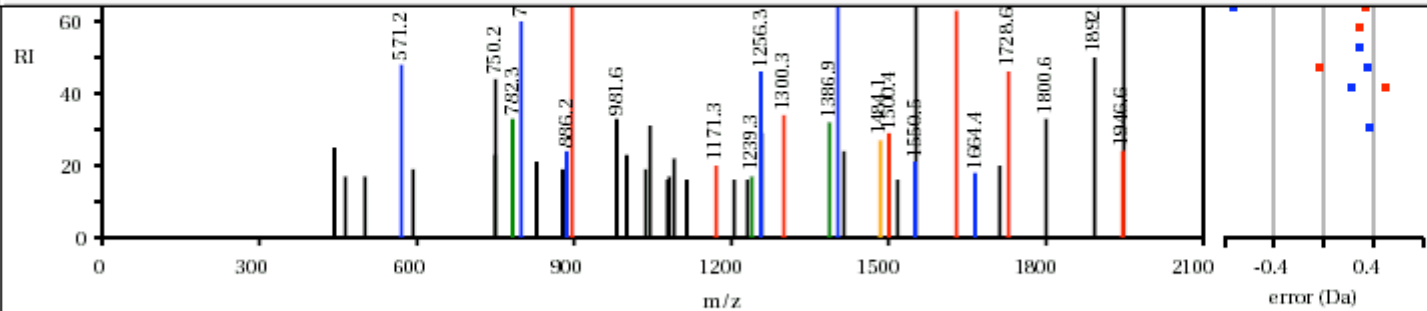
| #   | log(e) | m+h      | delta | z | sequence                  |
|-----|--------|----------|-------|---|---------------------------|
| 129 | -8.2   | 2299.116 | 1.621 | 2 | 40HTNFAVESIEQFFGGNLGFGK60 |

BC91\_01\_010.1248.1250.2.dta

H T N F A V E S I E Q F F G G N L G F G K



| bond | +1y             | +1y*     | +1b     | +1b*    |
|------|-----------------|----------|---------|---------|
| H1   | 2162.055        | 2145.039 | 138.062 | 121.046 |
| T2   | 2061.008        | 2043.992 | 239.110 | 222.094 |
| N3   | <b>1946.965</b> | 1929.949 | 353.153 | 336.137 |
| F4   | 1799.896        | 1782.880 | 500.221 | 483.205 |



| bond | +1y             | +1y <sup>+</sup> | +1b             | +1b <sup>+</sup> |
|------|-----------------|------------------|-----------------|------------------|
| H1   | 2162.055        | 2145.039         | 138.062         | 121.046          |
| T2   | 2061.008        | 2043.992         | 239.110         | 222.094          |
| N3   | <b>1946.965</b> | 1929.949         | 353.153         | 336.137          |
| F4   | 1799.896        | 1782.880         | 500.221         | 483.205          |
| A5   | <b>1728.859</b> | 1711.843         | <b>571.258</b>  | 554.242          |
| V6   | <b>1629.791</b> | 1612.775         | 670.327         | 653.311          |
| E7   | <b>1500.748</b> | <b>1483.732</b>  | <b>799.369</b>  | <b>782.353</b>   |
| S8   | 1413.716        | 1396.700         | <b>886.401</b>  | 869.385          |
| I9   | <b>1300.632</b> | 1283.616         | 999.485         | 982.469          |
| E10  | <b>1171.590</b> | 1154.574         | 1128.528        | 1111.512         |
| Q11  | 1043.531        | 1026.515         | <b>1256.586</b> | <b>1239.570</b>  |
| F12  | <b>896.463</b>  | 879.447          | <b>1403.655</b> | <b>1386.639</b>  |
| F13  | 749.394         | 732.378          | <b>1550.723</b> | 1533.707         |
| G14  | 692.373         | 675.357          | 1607.745        | 1590.729         |
| G15  | 635.351         | 618.335          | <b>1664.766</b> | 1647.750         |
| N16  | 521.308         | 504.292          | 1778.809        | 1761.793         |
| L17  | 408.224         | 391.208          | 1891.893        | 1874.877         |
| G18  | 351.203         | 334.187          | 1948.915        | 1931.899         |
| F19  | 204.134         | 187.118          | 2095.983        | 2078.967         |
| G20  | 147.113         | 130.097          | 2153.005        | 2135.989         |

## main spectra listing

[MODEL](#)

#129,  $e = 5.8e-09$ ,  $M+H = 22300.741.621$ ,  $VYRR40$ HTNFAVESIEQFFGGNLGFGK60KSSA,

**log(E) = -125.2**, [Contig31\\_171](#)

Model protein sequences

Supporting evidence

#128,  $e = 1.9e-04$ ,  $M+H = 32301.81.681$ ,  $VYRR40$ HTNFAVESIEQFFGGNLGFGK60KSSA,

**log(E) = -125.2**, [Contig31\\_171](#)

Model protein sequences

Supporting evidence

#210,  $e = 1.4e-03$ ,  $M+H = 32300.790.677$ ,  $VYRR40$ HTNFAVESIEQFFGGMLGFGK60KSSA,

**log(E) = -125.2**, [Contig31\\_171](#)

Model protein sequences

Supporting evidence

#220,  $e = 7.6e-07$ ,  $M+H = 22299.810.696$ ,  $VYRR40$ HTNFAVESIEQFFGGMLGFGK60KSSA,

**log(E) = -125.2**, [Contig31\\_171](#)

Model protein sequences

Supporting evidence

#1081,  $e = 4.5e-02$ ,  $M+H = 1979.6790.168$ ,  $RINR69$ SGDLITQVF77LKVT,

**log(E) = -125.2**, [Contig31\\_171](#)

Model protein sequences

Supporting evidence

#135,  $e = 9.1e-05$ ,  $M+H = 21222.560.875$ ,  $RINR69$ SGDLITQVFLK79VTLP,

**log(E) = -125.2**, [Contig31\\_171](#)

Model protein sequences

Supporting evidence

#57,  $e = 2.7e-06$ ,  $M+H = 21221.160.467$ ,  $RINR69$ SGDLITQVFLK79VTLP,


**log(E) = -125.2**, [Contig31\\_171](#)

GPM - performance - Mozilla

File Edit View Go Bookmarks Tools Window Help

Back Forward Reload Stop <http://10.234.244.194/thegpm-cgi/perform.pl?path=/tandem/archive/f84b367.x> Search Print

Home Bookmarks People Google annot Tools Debug Actualite bioinfo EMBOSS: The Applic... the GPM site local



modelling performance statistics

[MODEL](#) | [parameters](#) | [XML](#) | [spectra view](#)

**Performance statistics:**

| Parameter                                    | Value                   |
|--|-------------------------|
| list path, sequence source #1:               | ../fasta/mimi.fasta     |
| list path, sequence source #2:               | ../fasta/crap.fasta.pro |
| modelling, spectrum noise suppression ratio: | 0.09                    |
| modelling, total proteins used:              | 1259                    |
| modelling, total spectra assigned:           | 139                     |
| modelling, total spectra used:               | 1114                    |
| process, start time:                         | 2004:05:26:16:51:12     |
| process, version:                            | x! tandem 2004.04.01    |
| refining, # input models:                    | 12                      |
| refining, # input spectra:                   | 899                     |
| refining, # partial cleavage:                | 11                      |
| refining, # potential N-terminii:            | 2                       |
| refining, # unanticipated cleavage:          | 26                      |
| timing, initial modelling total (sec):       | 1.10                    |
| timing, initial modelling/spectrum (sec):    | 0.001                   |
| timing, load sequence models (sec):          | 0.02                    |
| timing, refinement/spectrum (sec):           | 0.003                   |

*perform.pl, v. 2004.03.01*

1200 sequences protéiques  
dans la base de données  
1114 spectres au départ

Temps: 1.10s



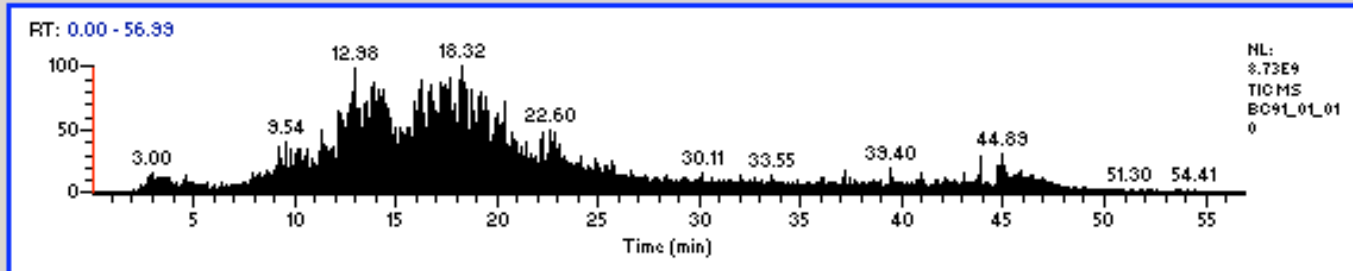




Dta Information: 1114

Scans, Charge

- 2, +1
- 4 - 6, +1
- 12, +1
- 14, +1
- 16, +1
- 24 - 26, +1
- 28, +2
- 28, +3
- 40, +2
- 40, +3
- 74, +1
- 76, +1
- 78, +1
- 80, +1
- 82, +1
- 86, +1
- 88, +2
- 88, +3
- 90, +1
- 92, +1
- 142 - 144, +2



Number of traces: 1

Filter: single, xc (+/- 1,2,3)=1.50,2.00,2.50

|    | Reference   |               |         |        |       | Score         | Accession |   |
|----|---|---------------|---------|--------|-------|---------------|-----------|---|
|    | Scan(s)   | Sequence      | MH+     | Charge | XC    | Delta Cn      | Sp        | F |
| #1 | <b>Contig31_171 [63834 - 62038] (REVERSE SENSE)</b> |               |         |        |       | <b>1168.3</b> |           |   |
|    | 142 - 144   | -.LGNYQGGK.   | 836.43  | 2      | 2.474 | 0.560         | 511.8     |   |
|    | 150 - 152   | -.LGNYQGGK.   | 836.43  | 1      | 2.131 | 0.375         | 273.9     |   |
|    | 316   | -.VDYIKK.-    | 765.45  | 1      | 1.921 | 0.339         | 428.4     |   |
|    | 322   | -.LGNYQGGK.   | 836.43  | 1      | 1.665 | 0.108         | 195.8     |   |
|    | 374 - 376   | -.YPEVEK.-    | 764.38  | 1      | 1.897 | 0.252         | 290.7     |   |
|    | 446 - 448   | -.IHTDFENDR.- | 1146.52 | 2      | 2.627 | 0.663         | 360.4     |   |
|    | 452   | -.IHTDFENDR.- | 1146.52 | 1      | 1.524 | 0.287         | 202.3     |   |
|    | 458   | -.IHTDFENDR.- | 1146.52 | 1      | 2.326 | 0.622         | 458.5     |   |
|    | 462 - 464   | -.VEGIIR.-    | 686.42  | 1      | 1.657 | 0.374         | 297.9     |   |
|    | 466   | -.VEGIIR.-    | 686.42  | 2      | 2.139 | 0.455         | 507.8     |   |
|    | 486   | -.MKYPEVEK.-  | 1023.52 | 3      | 2.584 | 0.278         | 904.1     |   |

Coverage - Contig31\_171 [63834 - 62038] (REVERSE SENSE) □



Reference: Contig31\_171 [63834 - 62038] (REVERSE SENSE) |

Database: mimi.fasta      Monoisotopic Mass: 68033.9      Number of Amino Acids: 599      pI: 5.62

|     | 1-10       | 11-20          | 21-30      | 31-40          | 41-50           | 51-60      | 61-70          | 71-80       | 81-90      | 91-100       |
|-----|------------|----------------|------------|----------------|-----------------|------------|----------------|-------------|------------|--------------|
| 1   | VTSSTTGR   | LVETKFLK       | VQSSPLGK   | LRGMRVFK       | VYRRHTNFAVESIEQ | FFGGNLGFGK | KSSAEINR       | SGDLITQVFLK | TLPEVRYCGD | FTNFGHVEFA   |
| 101 | WVRNIGHAIV | EETELEIGGS     | PIDKHYGDWL | QIWQDVSSSK     | DHEKGLAKML      | GDVPELTSIS | TLSDVPDNT      | VLKPSYTLV   | PLQFYFNR   | NWGLALPLIALQ |
| 201 | YHQVRIYVKF | RQADQCYIAS     | DAFKSGCGNL | QLDDVSLYVM     | YVFLDTEERR      | RFAQVSHEYL | IEQLQFTGEE     | SAGSSNSAKY  | KLNFNHPVKA | IYVWTKLGN    |
| 301 | QGGKFM     | TYDPCWENARENA  | AKLLLLAQYD | LDDWGYFQEP     | GGYECEGNDG      | RSYVGD     | CGVQYTAVDPSNPS | EEPSYIFMDT  | TTAEAFD    | GSLLIGKLA    |
| 401 | LLKRNKD    | VDLKDKVEGIIRI  | HTDFENDRMK | YPEVEK         | ITRMDLTLHDL     | SVPLSKYD   | VDMRVDYIKK     | FDVTVWQHNM  | FGLLIDGSGN | PTHEAELQLNG  |
| 501 | SKRGGI     | WYDTPMPTVHHTKS | PRDGVNV    | VFSFALNPEEHQPS | CTCNFSRIDT      | AQLNLWFQHF | TNHKFADVFA     | DNDNKVLIFA  | VNYNVL     | RMLSGMAGLAYS |

Protein Coverage:

|                                     | Sequence              | MH+     | % by Mass | Position  | % by AA's |
|-------------------------------------|-----------------------|---------|-----------|-----------|-----------|
| <input checked="" type="checkbox"/> | HTNFAVESIEQFFGGNLGFGK | 2299.11 | 3.38      | 40 - 60   | 3.51      |
| <input checked="" type="checkbox"/> | SGDLITQVFLK           | 1220.69 | 1.79      | 69 - 79   | 1.84      |
| <input type="checkbox"/>            | VTLPEVR               | 813.48  | 1.20      | 80 - 86   | 1.17      |
| <input type="checkbox"/>            | YCGDFTNFGHVEFAWVR     | 2047.91 | 3.01      | 87 - 103  | 2.84      |
| <input checked="" type="checkbox"/> | NIGHAIVEETELEIGGSPIDK | 2221.14 | 3.26      | 104 - 124 | 3.51      |
| <input checked="" type="checkbox"/> | NIGHAIVEETELEIGGSPIDK | 2221.14 | 3.26      | 104 - 124 | 3.51      |
| <input checked="" type="checkbox"/> | NIGHAIVEETELEIGGSPIDK | 2221.14 | 3.26      | 104 - 124 | 3.51      |
| <input type="checkbox"/>            | NIGHAIVEETELEIGGSPIDK | 2221.14 | 3.26      | 104 - 124 | 3.51      |
| <input checked="" type="checkbox"/> | HYGDWMLQMWQDVSSSK     | 1948.92 | 2.86      | 125 - 140 | 2.67      |

Protein Coverage Totals

|              |         |
|--------------|---------|
| by Mass:     | 27075.6 |
| % by Mass:   | 39.80   |
| by Position: | 236     |
| % by AA's:   | 39.40   |

Ready

NUM

DisplayIons - SGDLITQVFLK

Display the ion series for charge state: **+1**

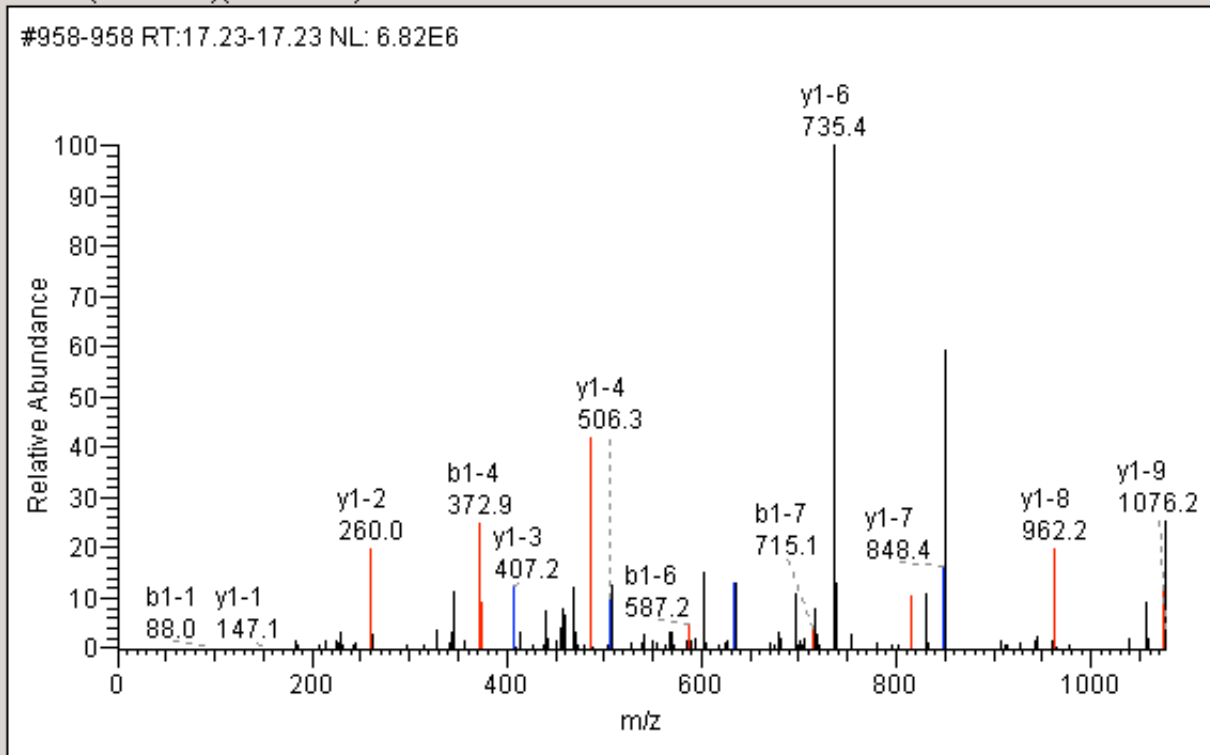
Dta: BC91\_01\_010.958.958.2

Precursor mass: 611.80

Mass type: Monoisotopic

Mod's: (C\* +57.022) (M# +15.999)

|    | AA | A Ions | B Ions  | Y Ions  |    |
|----|----|--------|---------|---------|----|
| 1  | S  |        | 88.04   | -       | 11 |
| 2  | G  |        | 145.06  | 1133.64 | 10 |
| 3  | D  |        | 260.09  | 1076.62 | 9  |
| 4  | L  |        | 373.17  | 961.59  | 8  |
| 5  | I  |        | 486.26  | 848.51  | 7  |
| 6  | T  |        | 587.30  | 735.42  | 6  |
| 7  | Q  |        | 715.36  | 634.37  | 5  |
| 8  | V  |        | 814.43  | 506.32  | 4  |
| 9  | F  |        | 961.50  | 407.25  | 3  |
| 10 | L  |        | 1074.58 | 260.18  | 2  |
| 11 | K  |        | -       | 147.09  | 1  |



S G D L I T Q V F L

Ready

NUM

Sur 10 échantillons,  
8 identifications  
coherentes

| Echantillon | Id match                           | Score gpm                        | Score sequest          |
|-------------|------------------------------------|----------------------------------|------------------------|
| 82          | -                                  | -                                | -                      |
| 83          | 29_13<br>34_197<br>34_273          | -<br>-<br>-22.9                  | 30<br>20<br>-          |
| 84          | 34_129<br>34_156                   | -19.9<br>-16.1                   | 190<br>40              |
| 85          | 34_69<br>34_129                    | -14.7<br>-3.5                    | 100<br>20              |
| 86          | 29_35                              | -13.9                            | 30.2                   |
| 87          | 31_171<br>31_180<br>31_7<br>31_161 | -56.1<br>-26.6<br>-16.1<br>-13.1 | 300<br>110<br>70<br>20 |
| 88          | 31_161<br>34_103<br>34_288         | -49<br>-23.7<br>-13.8            | 348<br>72<br>30        |
| 89          | 31_161<br>34_152                   | -50<br>-5.7                      | 210<br>30              |
| 90          | 31_161                             | -111                             | 1140                   |
| 91          | 31_171<br>31_180<br>31_115         | -125.2<br>-73.8<br>-7            | 1168<br>290<br>30      |

## Conclusion:

Une différence considérable en temps de calcul:

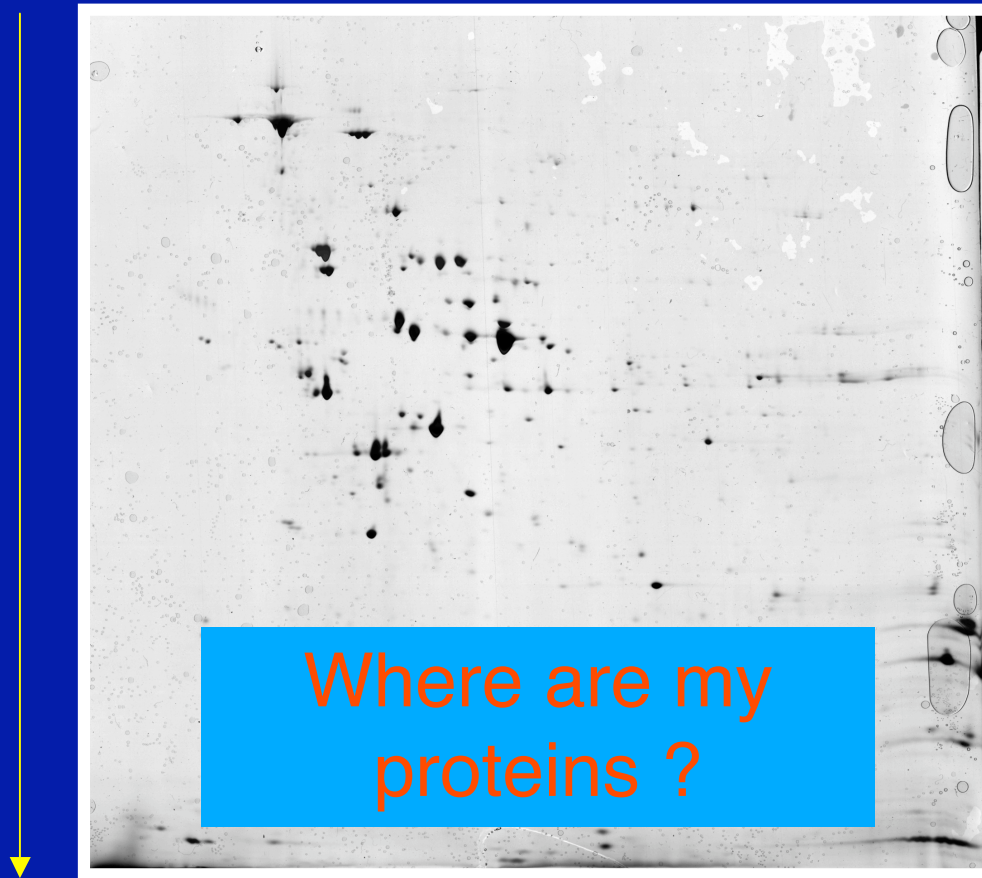
theGPM: 1 s

Sequest : 10 mn

Pour des résultats similaires

## 2D-gel of *Rickettsia conorii*

Molecular weight



Where are my  
proteins ?

Isoelectric Point (pI)

2D-gel provided by P. Renest et al.

Gelprint: In Silico 2D-PAGE - Konqueror

Location Edit View Go Bookmarks Tools Settings Window Help

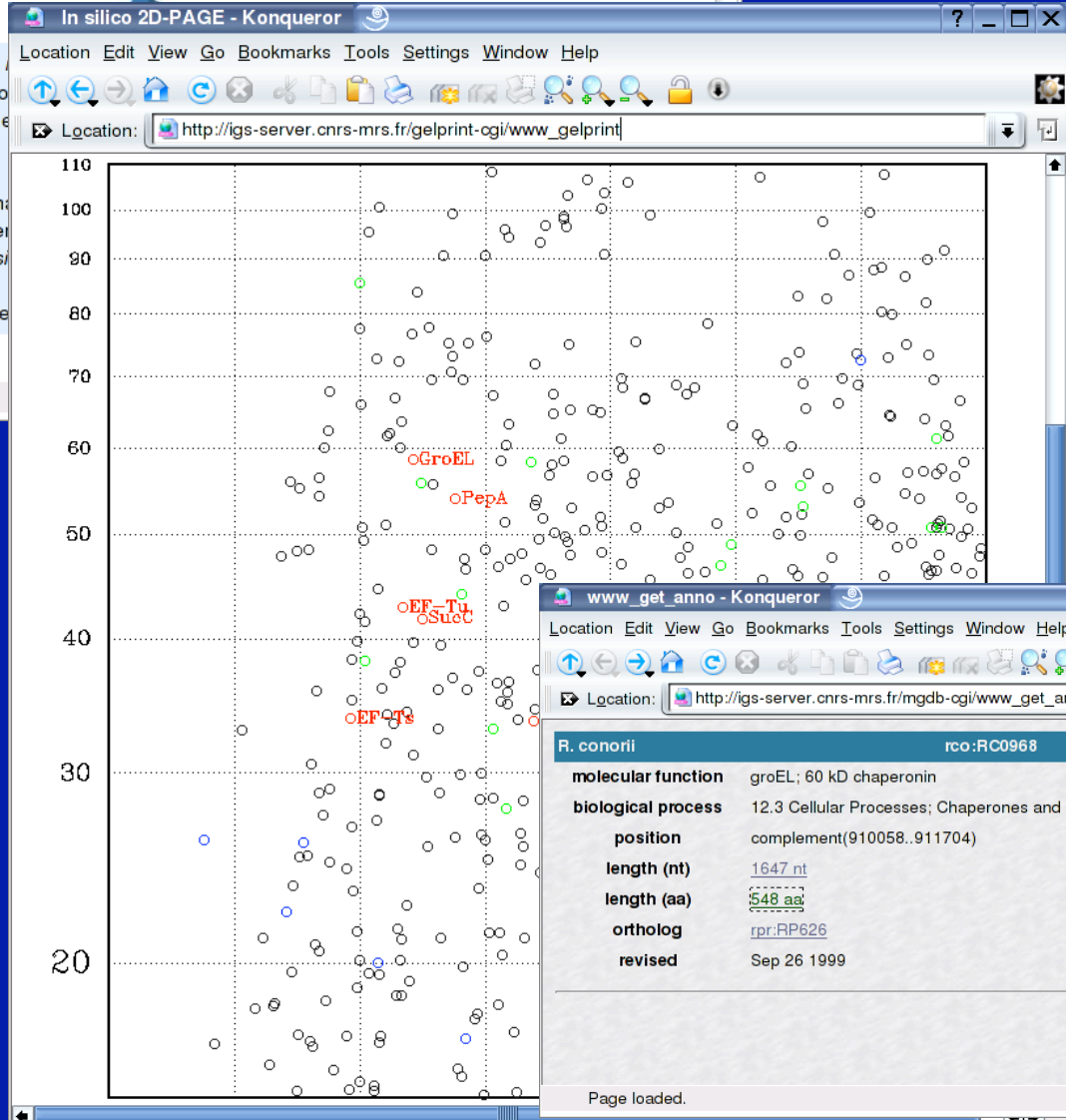
Location: <http://igs-server.cnrs-mrs.fr/gelprint/>

# Gelprint

**Gelprint** is a tool that generates an *in silico* 2D-gel (molecular weight) data for a set of proteins. Compare two-dimensional polyacrylamide gel electrophoresis (2D-PAGE) data of interest.

The *in silico* 2D-gel in printable format can be used to overlay and compare the two different gels. Gelprint also outputs a clickable *in silico* 2D-gel.

Click links here to see example pages.



www\_get\_anno - Konqueror

Location Edit View Go Bookmarks Tools Settings Window Help

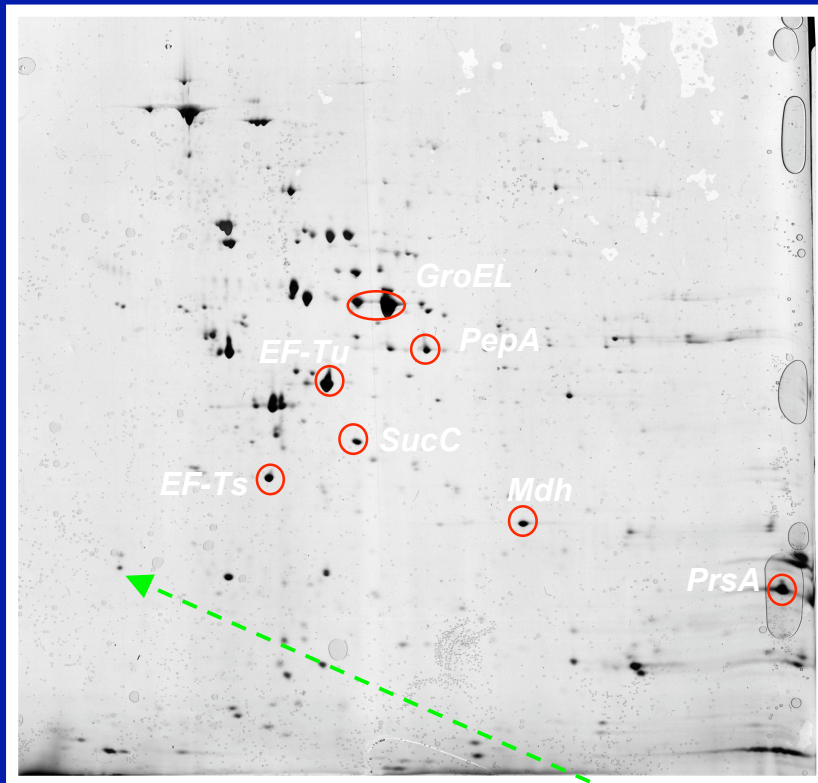
Location: [http://igs-server.cnrs-mrs.fr/mgdb/cgi/www\\_get\\_anno?ric.ann+rco:RC0968](http://igs-server.cnrs-mrs.fr/mgdb/cgi/www_get_anno?ric.ann+rco:RC0968)

| R. conorii                | rco:RC0968  | CDS |
|---------------------------|---|-----|
| <b>molecular function</b> | groEL; 60 kD chaperonin   |     |
| <b>biological process</b> | 12.3 Cellular Processes; Chaperones and stress-induced proteins |     |
| <b>position</b>           | complement(910058..911704)                                      |     |
| <b>length (nt)</b>        | 1647 nt   |     |
| <b>length (aa)</b>        | 548 aa  |     |
| <b>ortholog</b>           | rpr:RP626   |     |
| <b>revised</b>            | Sep 26 1999   |     |

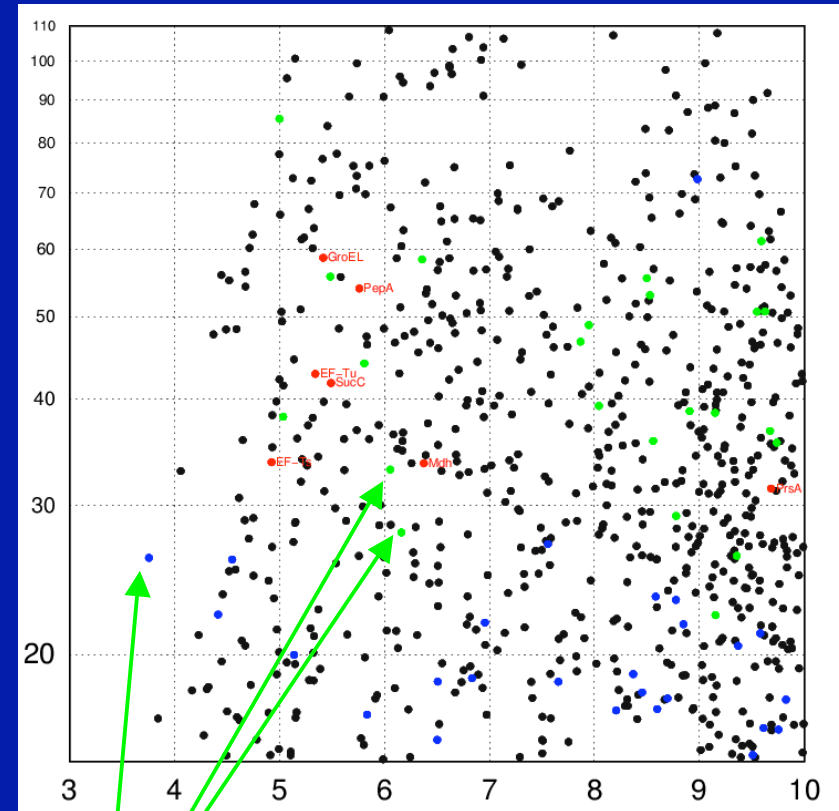
Page loaded.



## 2D gel (experimental)



## 2D gel (*in silico*)



?

**Proteins of interest !**  
"Split genes and Repeat-containing proteins"

# Input data for Gelprint

- List of theoretical (pre-computed) pI and Mw for a set of proteins
- Ranges of pI/Mw to show in the in silico 2D-gel
- Database address to make links from in silico 2D-gel to the database
- Gelprint generates in silico 2D-gels in Postscript/PDF that fit to your “real” gel
  - Size of your 2D-gel (in centimeter)
  - Some references points for fitting (i.e. positions of some identified spots on the real 2D-gel (in centimeter))

## *Gelprint*

<http://igs-server.cnrs-mrs.fr/gelprint/>

## *Acknowledgements*

Patricia Renesto

Didier Raoult

(Unité des Rickettsies, CNRS UMR6020, Marseille)