

Informatik Anwendungen in der Genom- und Proteomforschung

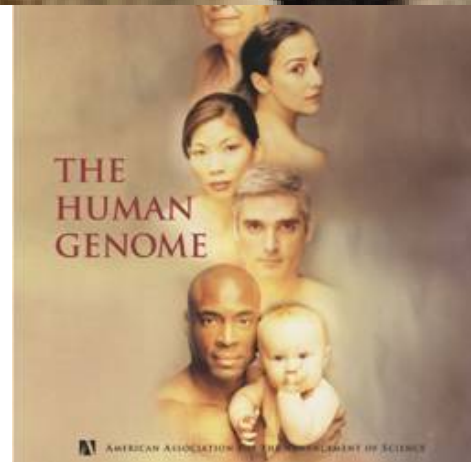
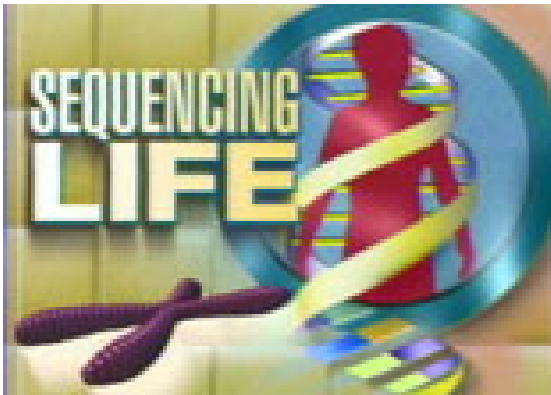
Knut Reinert
FU Berlin

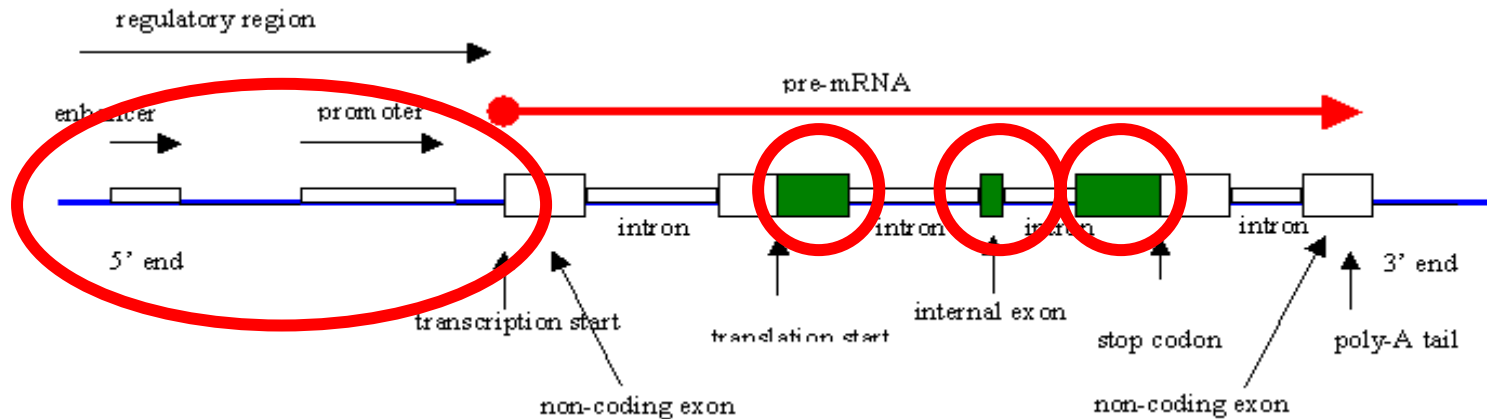
Mai 2006
HU Berlin, Adlershof



funded by the German Federal Ministry for Education and Research





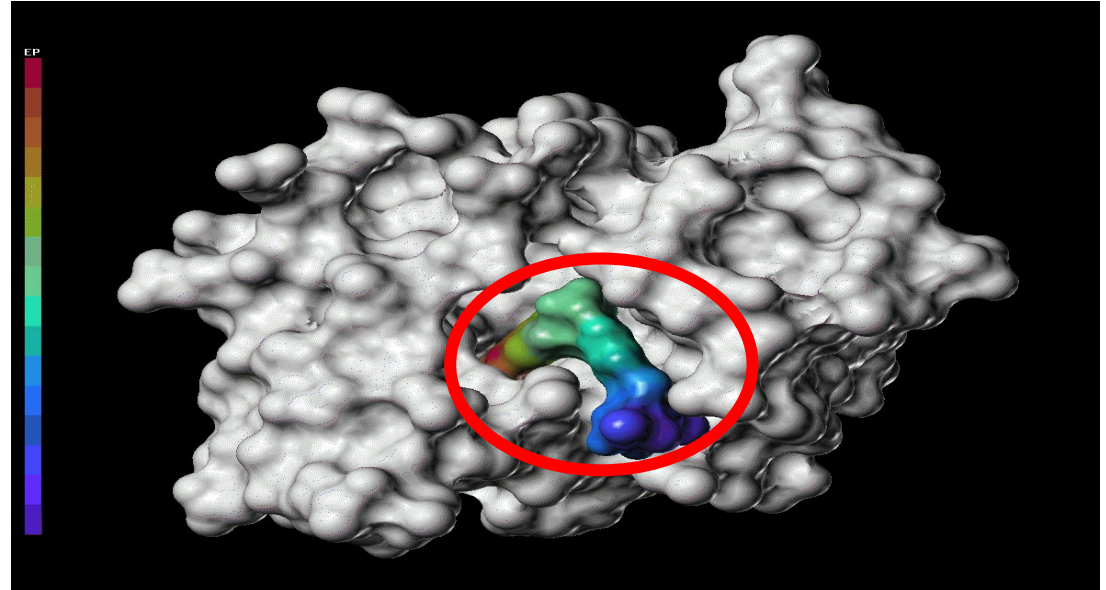


Some pictures with the courtesy of the MPI für Informatik, Saarbrücken

Gene finding:

- Improve prediction of coding and regulatory regions
- Comparing multiple genomes is promising

Molecular therapie
of deseases:
protein-protein
docking



Methotrexate bindet an dihydrofolate reductase

Find drugs that alter or inhibit the function of the target molecule.

Searching data bases helps to find suitable candidates and reduce side effects.

Identifying SNPs (or other polymorphisms)

```
GACGTGCACTAAATCGCGCAACTG
TTCGGGTTGGACGTGCACTAACTCG
TTCGGGTTGGACGTGTACTAAATCGCGCAACTG
GGGTTGCACGTGCACTAAATCGCGCAACTG
```

Identifying SNPs allows:

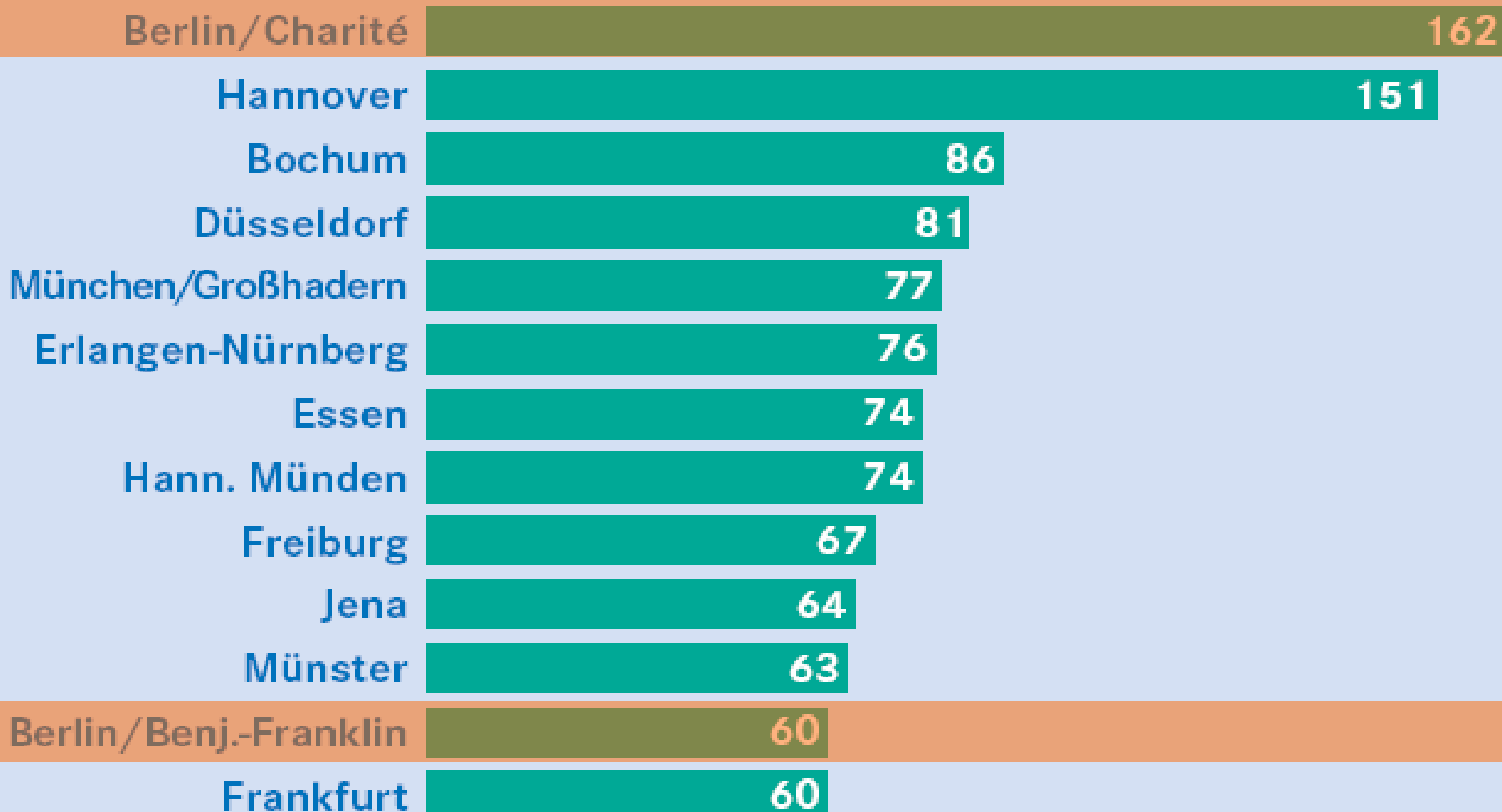
- the association of patterns of genetic diversity to diseases
- the association of genetic patterns to drug tolerance

Deutschland 2004

Nierentransplantationen (ohne Lebendspende)

aufgeschlüsselt nach Transplantationszentren

Anzahl



Problem (Project together with Charité):

Within the first year 8-10% of the patients lose the transplant.

After 10 years about 50% have lost the transplant.

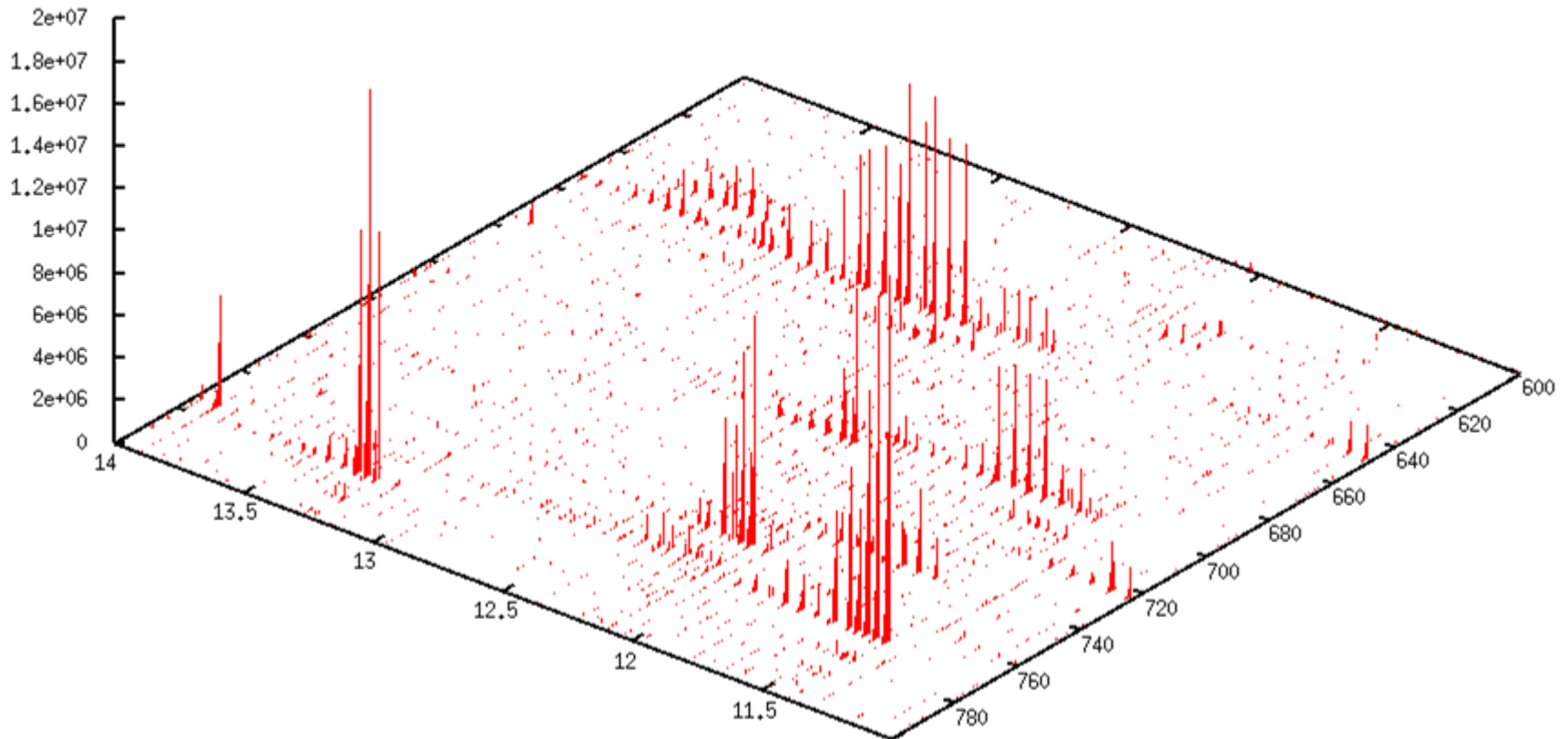
Diagnosis is invasive and leads sometimes to loss of graft.

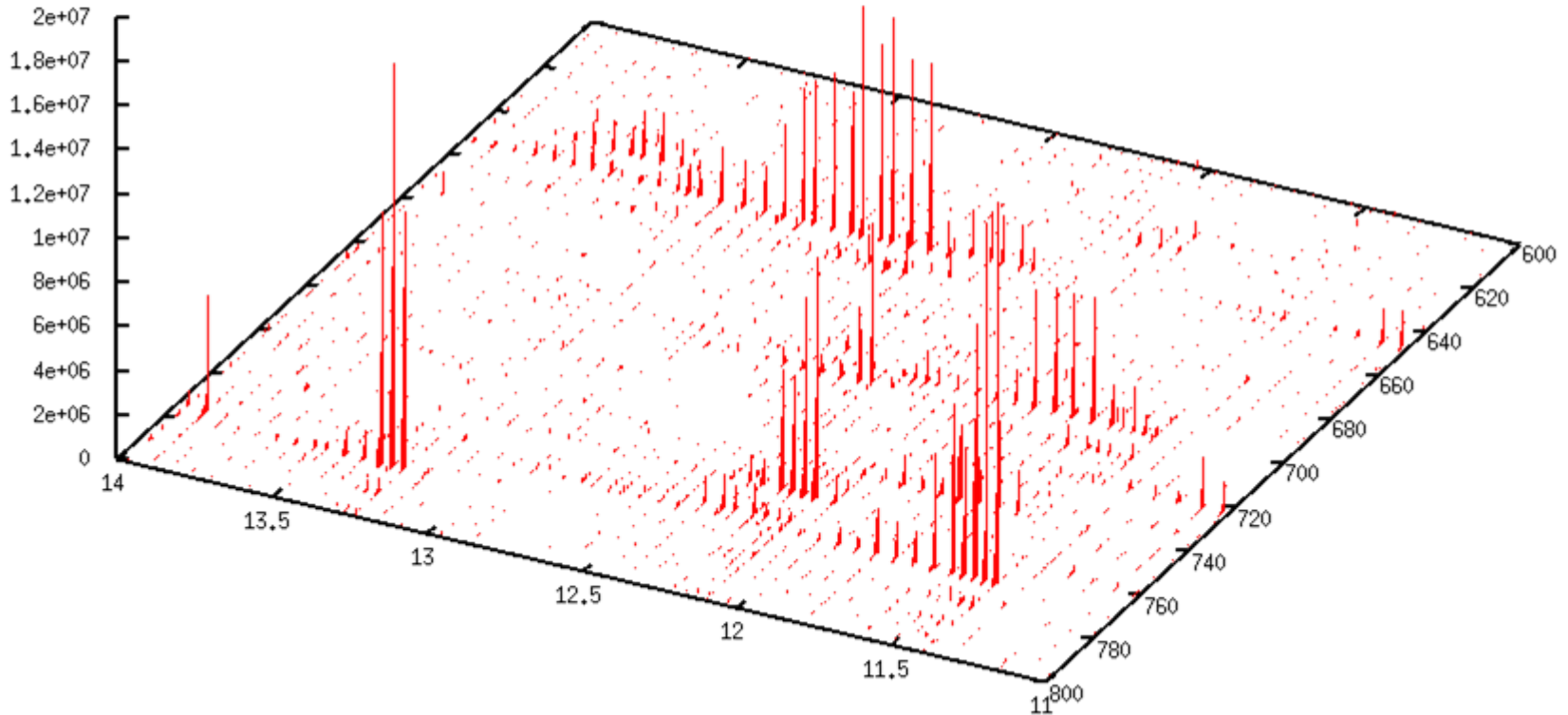
Goal:

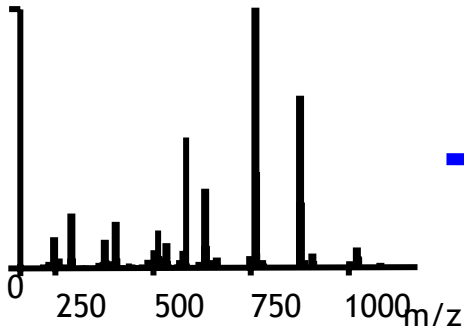
Analyse urine samples of patients and detect as early as possible diagnostic markers to counteract graft loss.



Automated measurement methods lead to terabytes of data (Prof. Schlüter)

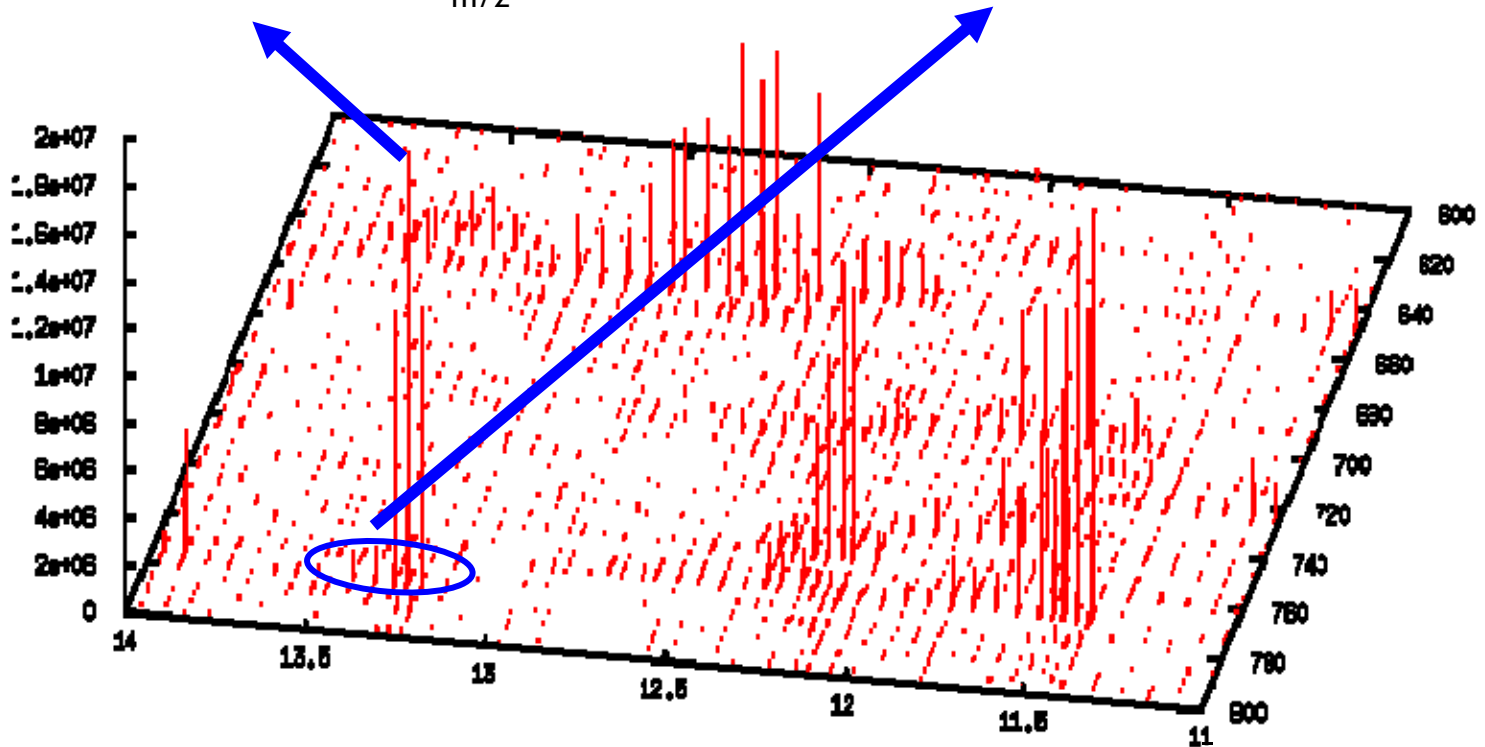




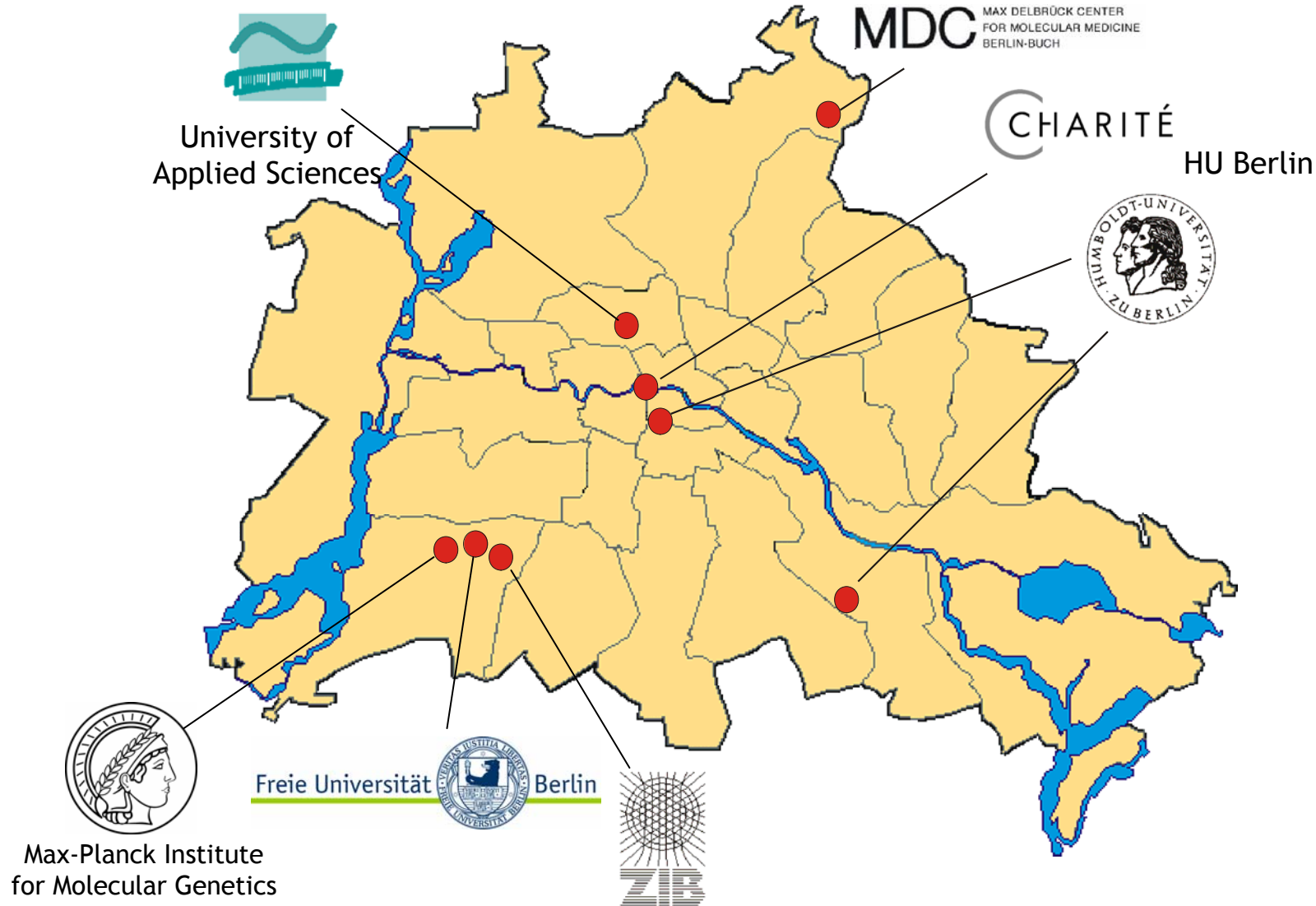


EVAFAQFGSDLASTK

3x increased expression



Bioinformatics in Berlin



- **Motivation for OpenMS**
- **Bioinformatics Issues in quantitative Proteomics**
 - Signal Processing
 - Feature Finding
 - Map Mapping
 - Differential Quantitation
 - Identification, Clustering
 - Software Engineering, Databases





Same genome...

...different proteomes



Genomics

Genome is rather
static

~ **30 k** genes

Established, fully
automated technology
(capillary sequencer)

Proteomics

Proteome is **dynamic**
(age, tissue, what you
had for lunch)

Up to **2000 k** Proteins

Emerging technology

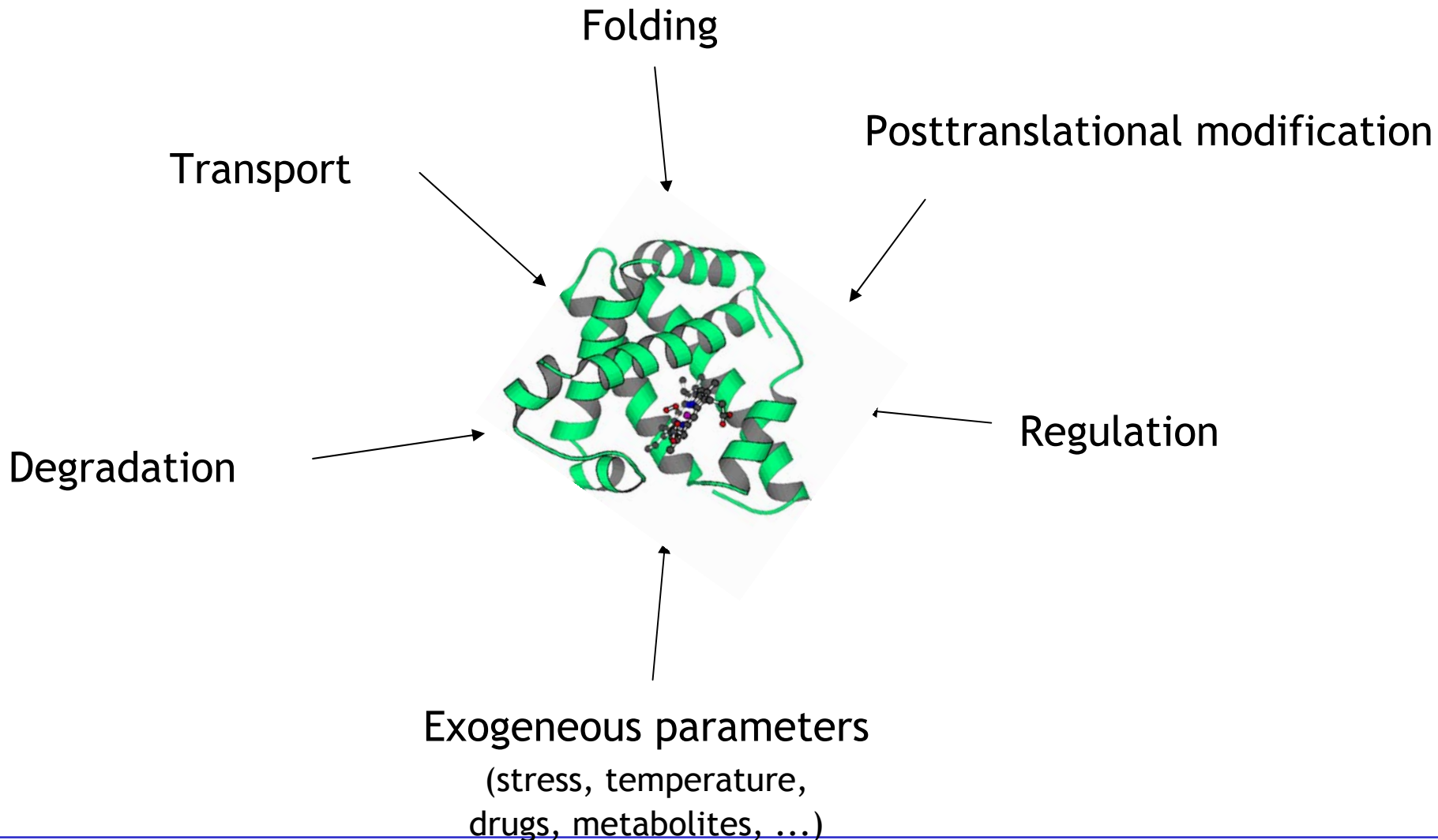
(MS, HPLC/MS,
protein chips)



- Transcription and translation are heavily regulated
- Protein expression levels are not static
- mRNA levels and protein levels often not correlated
- Contradictory results from seemingly similar methods
 - RNA chips
 - DNA chips
 - gene disruption
 - knock out

Anderson et al., Electrophoresis (1998), 19, 1853-61

Proteins = End Products?



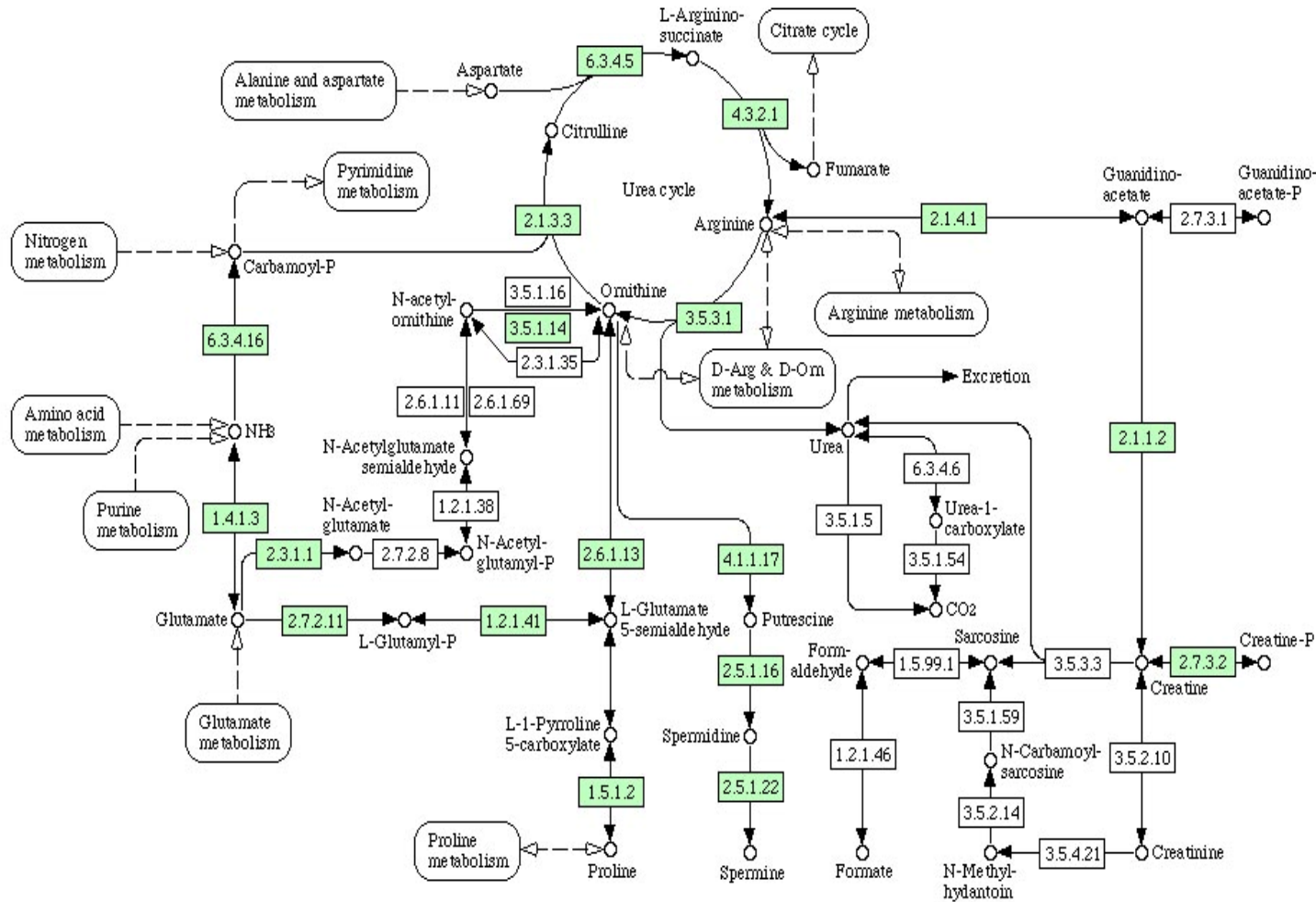
- ~~• Proteomics [n]: the branch of genetics that studies the full set of proteins encoded by a genome~~

(www.hyperdictionary.com)

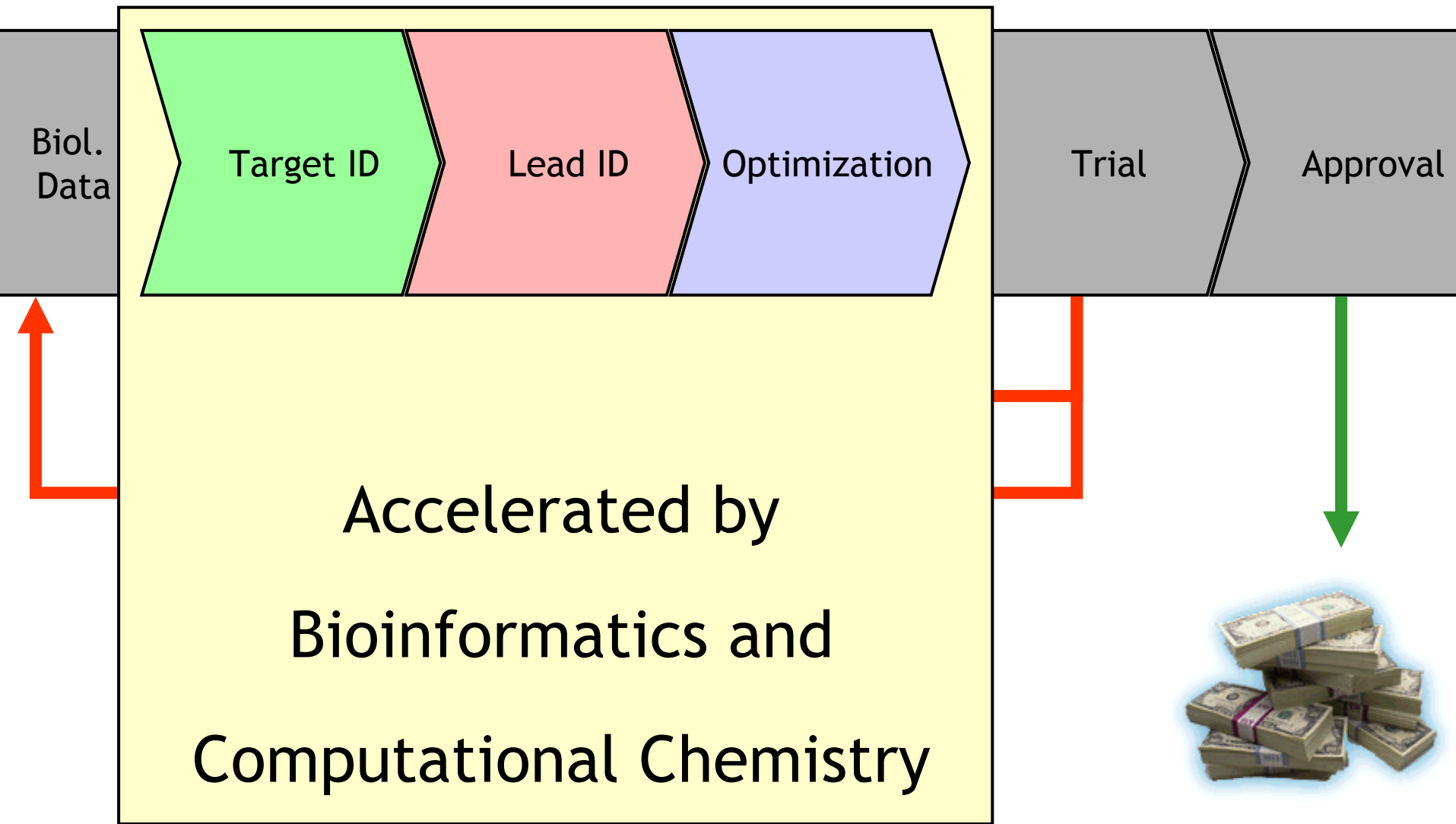
- Proteomics can be defined as *the qualitative and quantitative comparison of proteomes under different conditions to further unravel biological processes.*

(www.expasy.org)

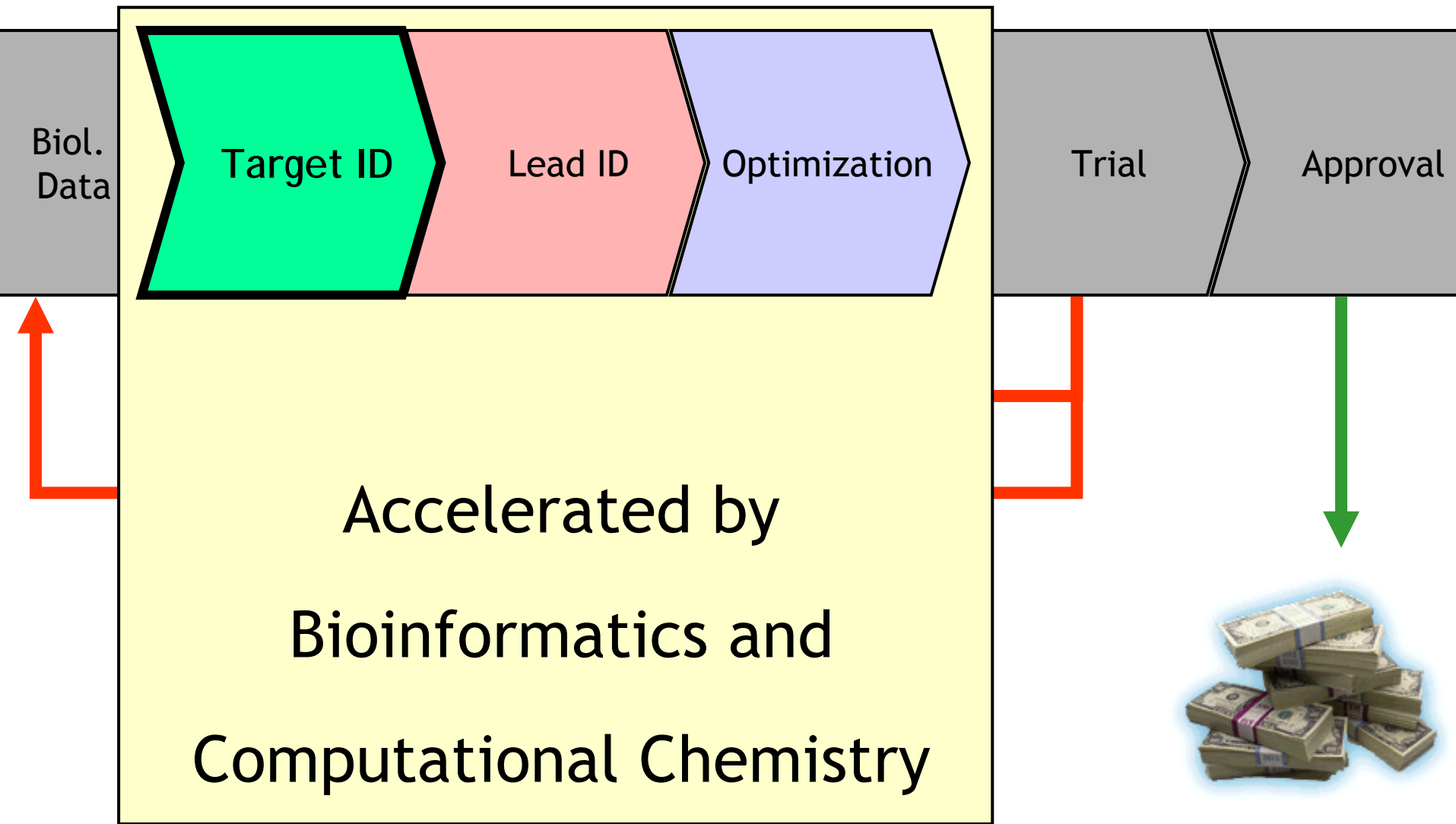
Target ID and biological networks



Drug Discovery Pipeline



Drug Discovery Pipeline



- **Diagnostics:** Find relevant patterns in one- or two-dimensional LC measurements
- **Time series:** Analyze the temporal behaviour in a time series experiment
- **Quantitative Measurements:** Determine absolute content of peptides using additive method (Myoglobin, Gliadin)

All techniques share the following steps:

1. **Separation** of proteins/peptides.

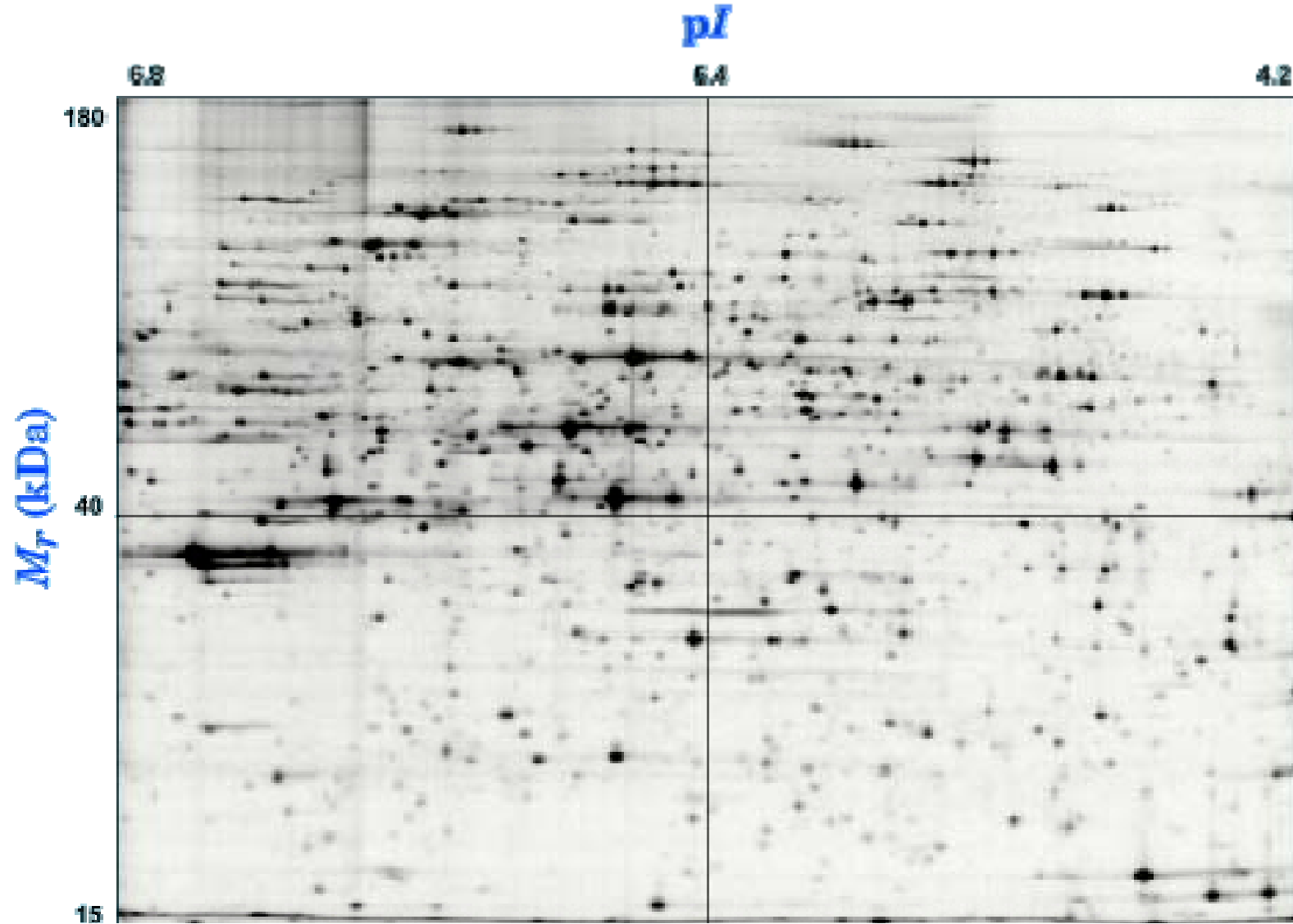
Typical techniques: 2D PAGE, HPLC, DC, etc.

2. **Assignment** of proteins/peptides for *relative* quantitation

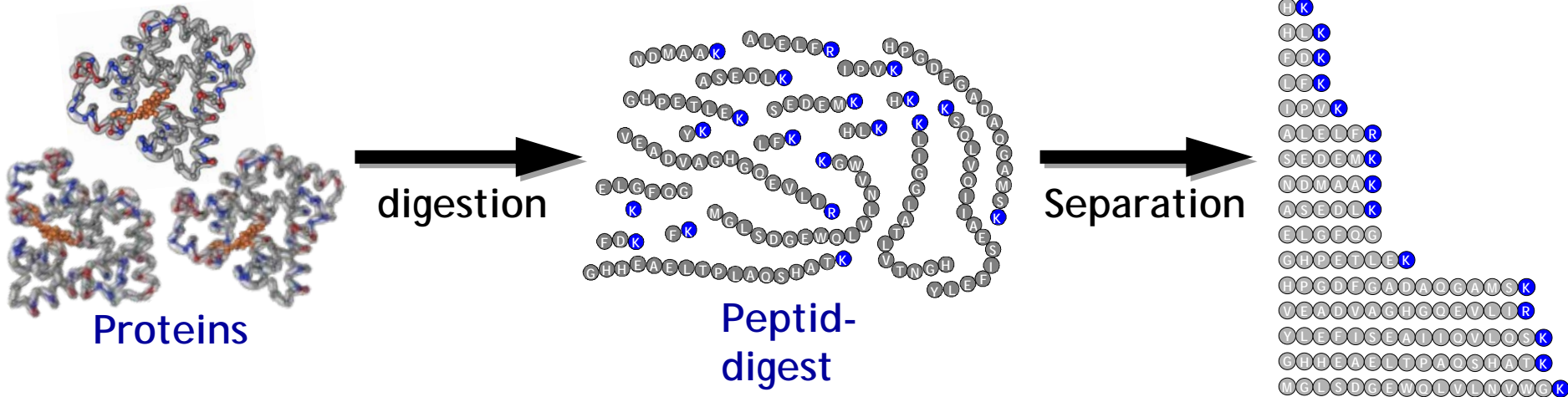
3. **(Relative) Quantitation**

4. **Identification** of peptides/proteins

Separation (2D-PAGE)



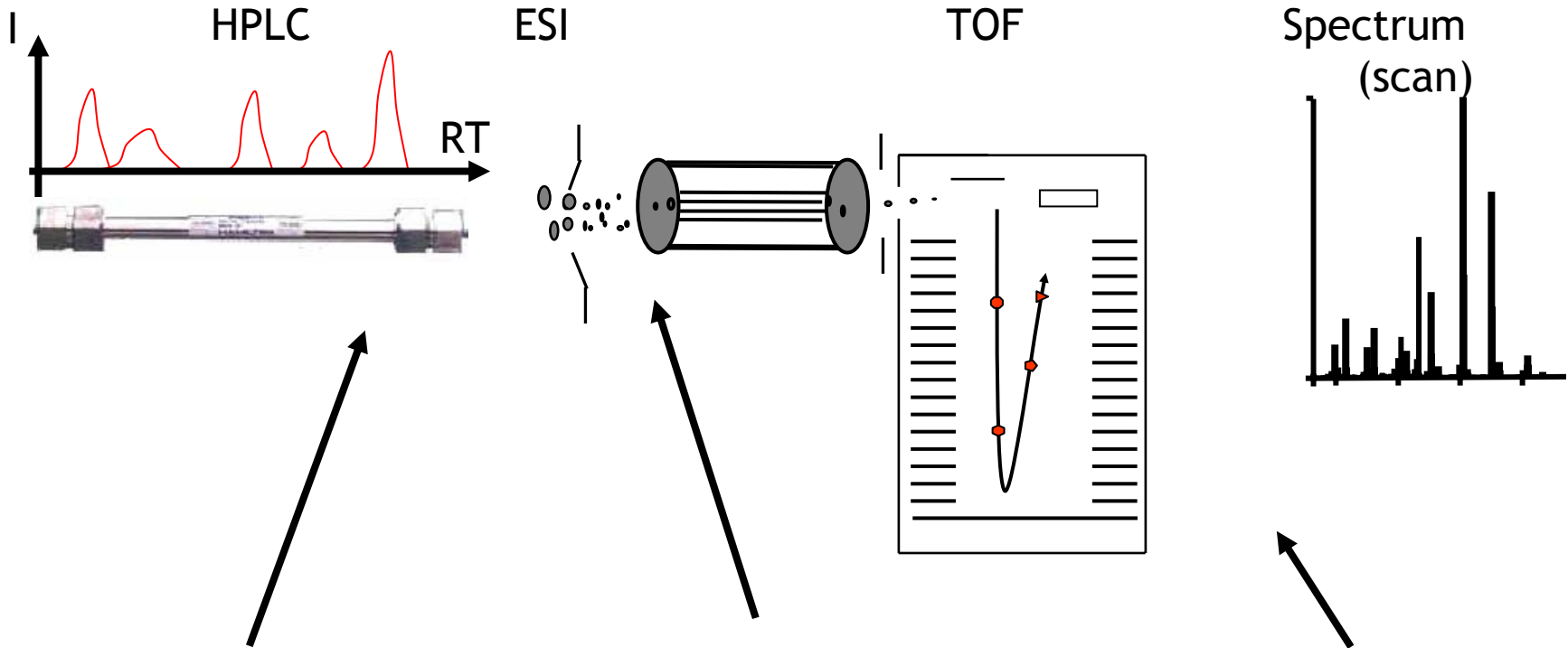
Separation (Shotgun Proteomics)



- **Key idea**

- Separation of whole proteins possible but difficult, hence digestion preferred
- Separate peptides
- Identify proteins through peptides

HPLC-MS Analysis

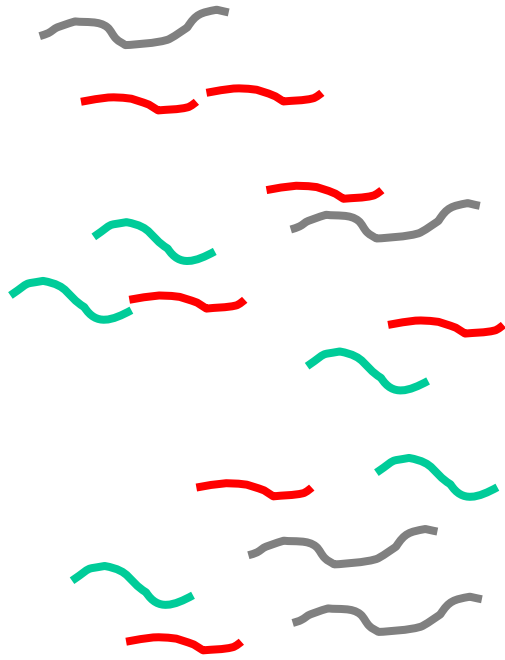


Separation 1
Different peptides
have different
retention time

Ionization
Peptide receives
z charge units

Separation 2
Detector measures
m/z

Single stage MS scan

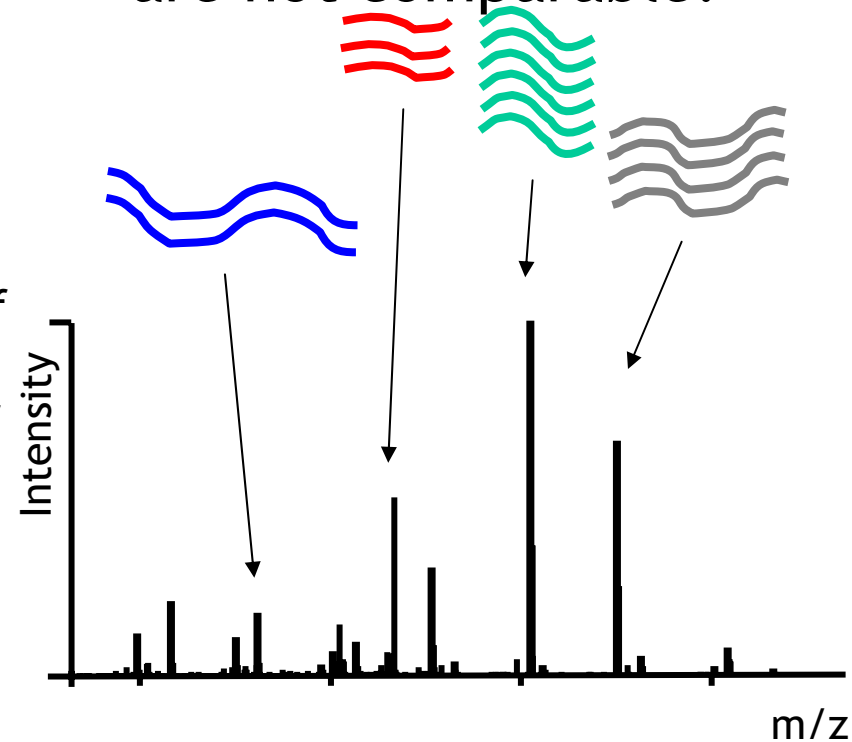


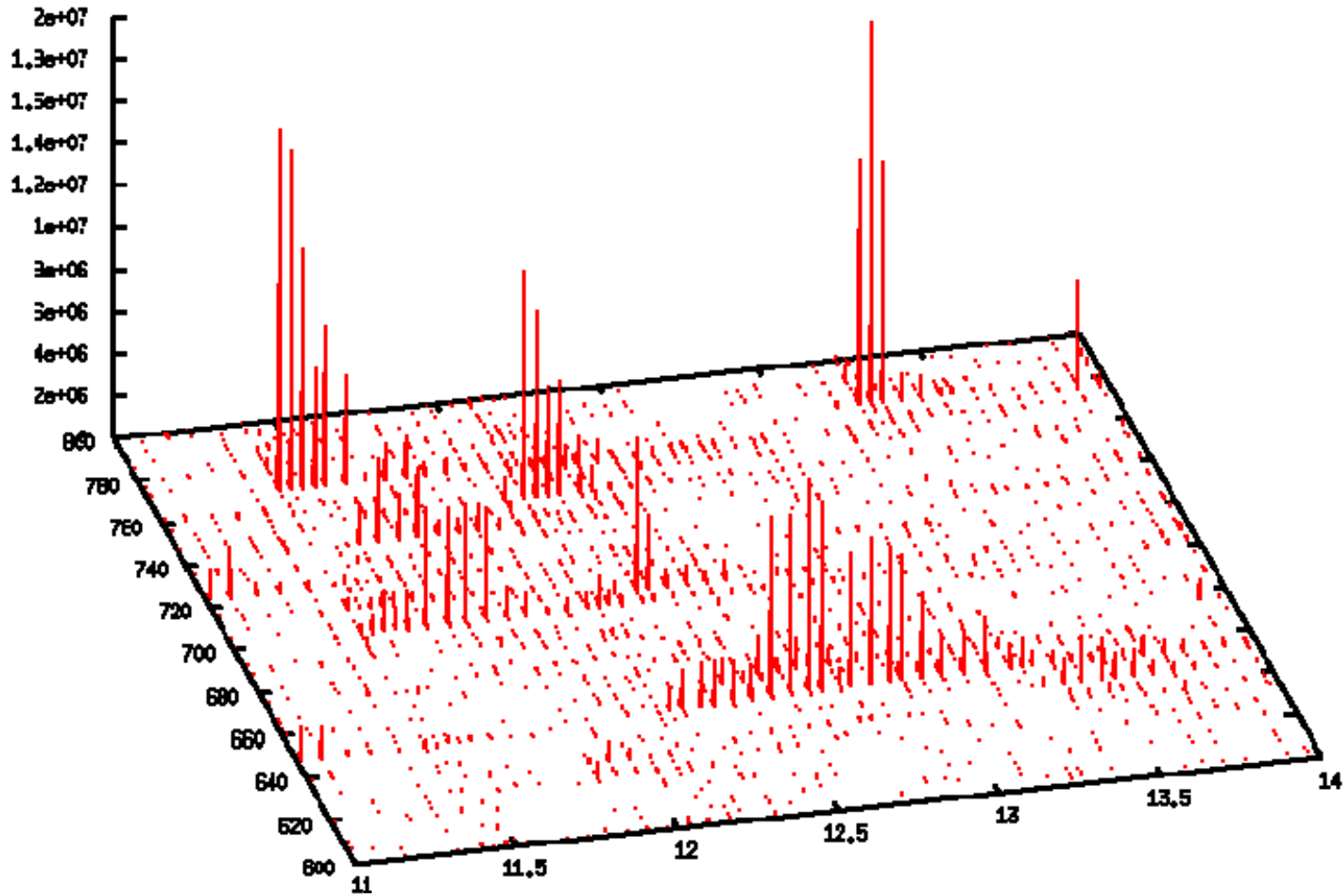
Mass spectrometry

Measures mass/charge ratio of ionized peptides for a short period of time

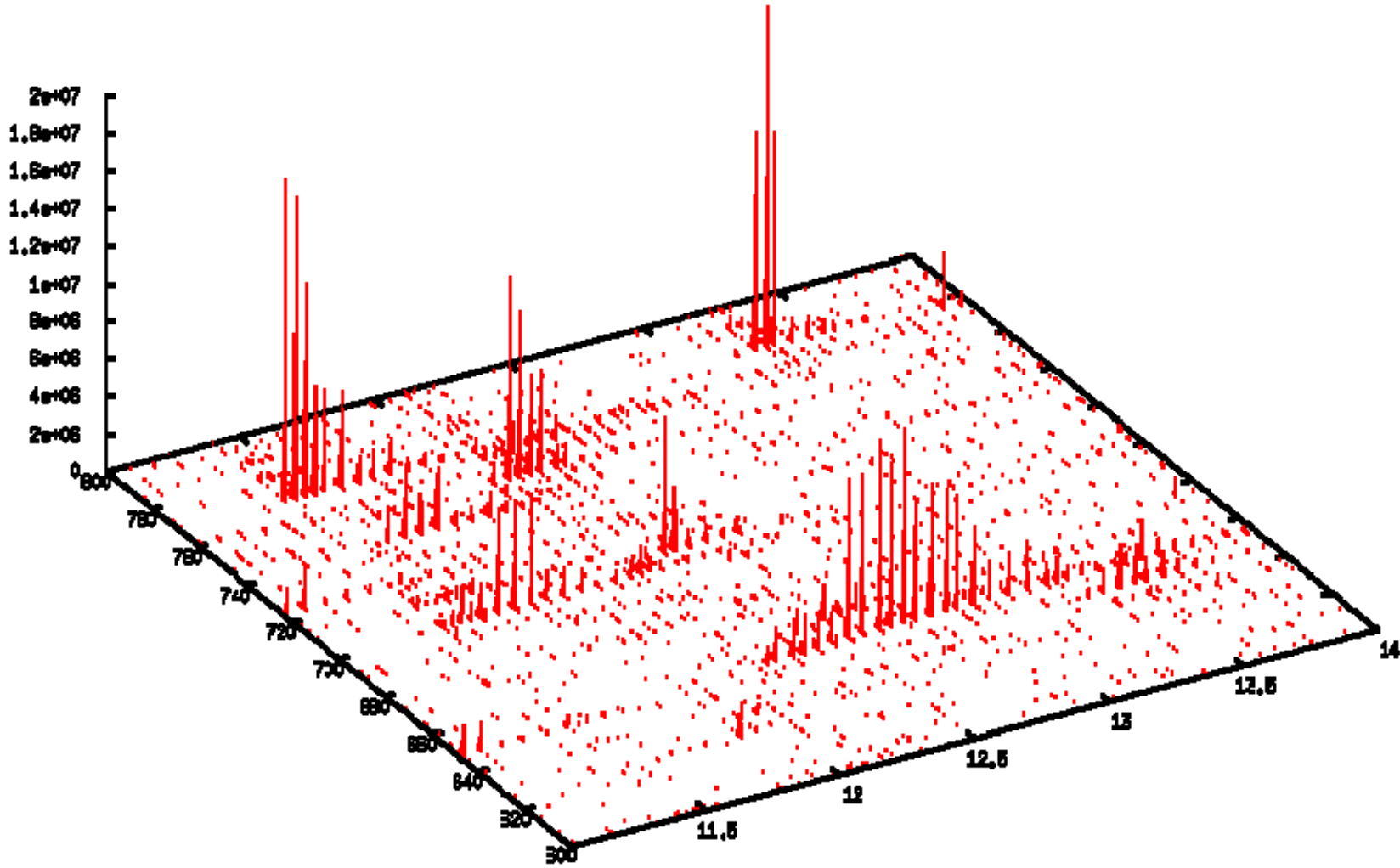


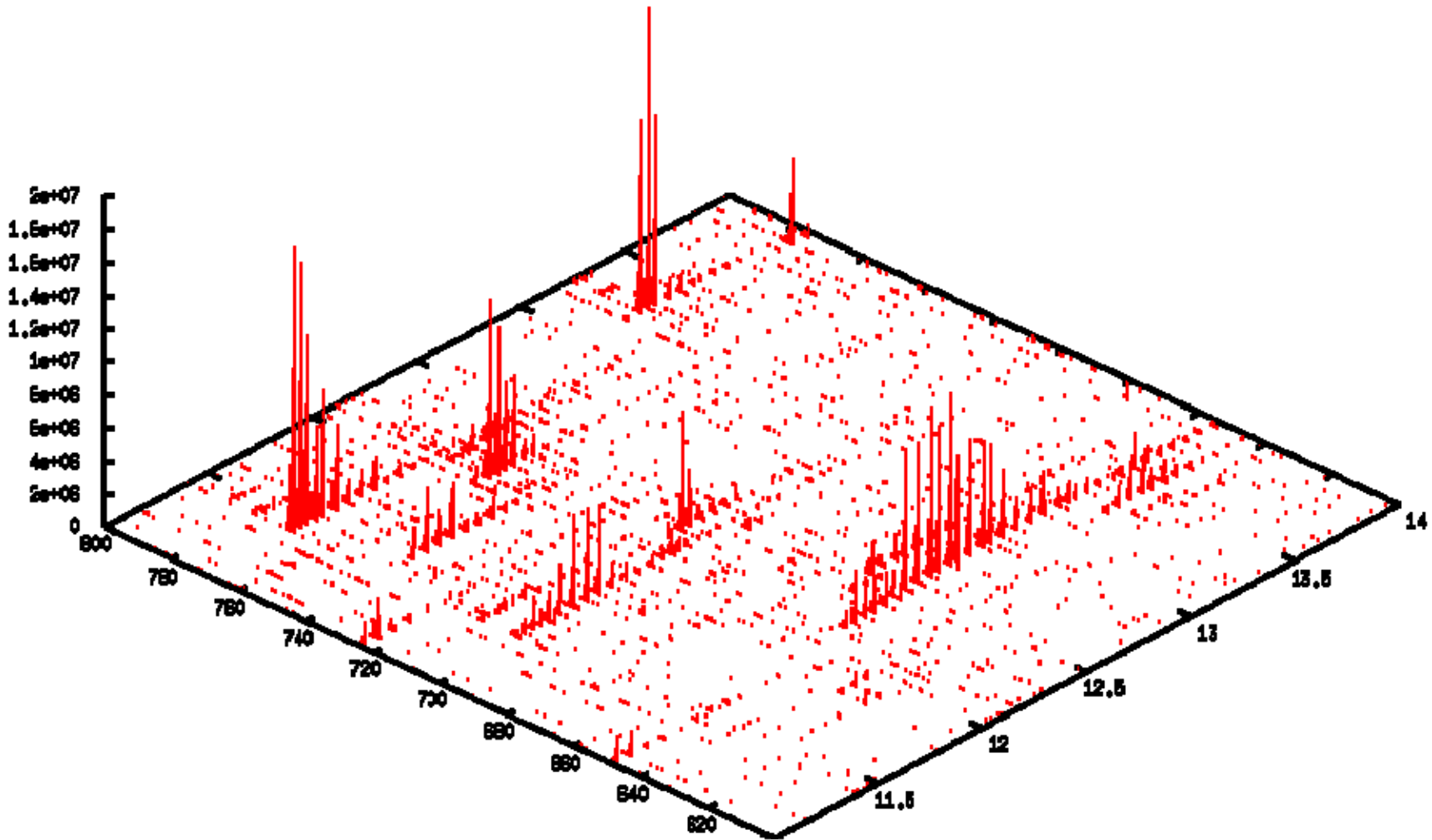
Peak intensity in scan corresponds to amount present, but intensities are not comparable!

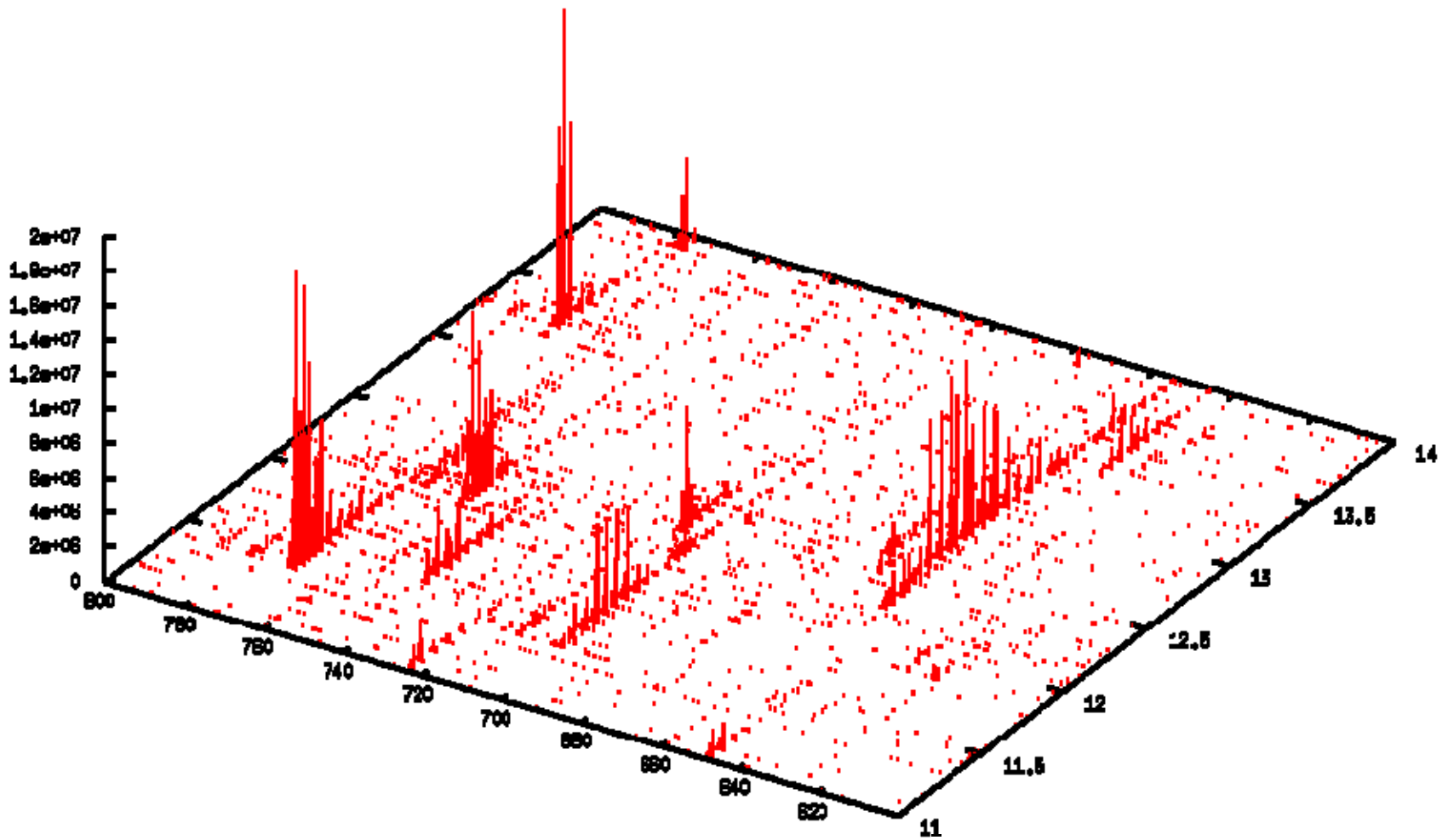


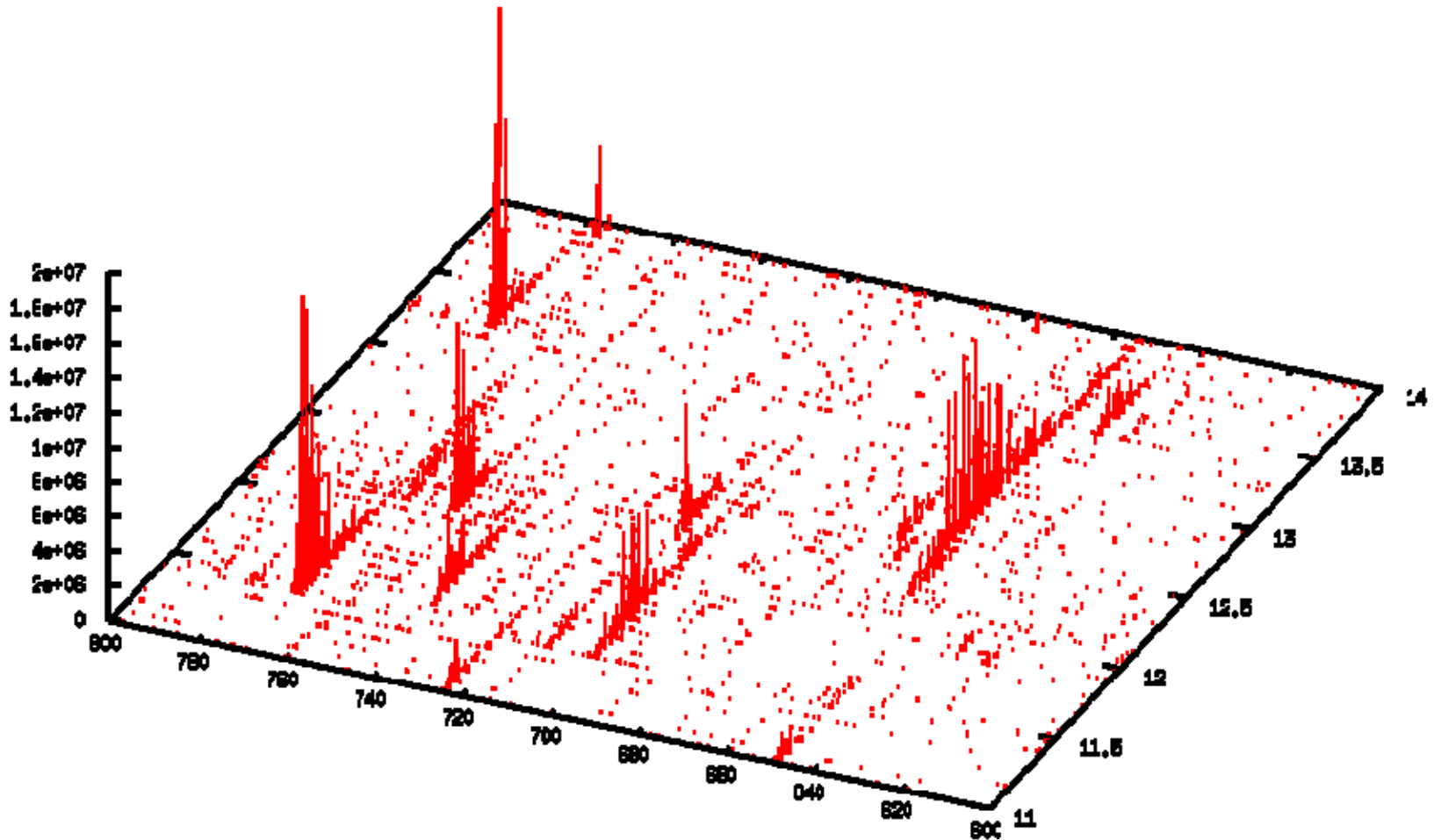


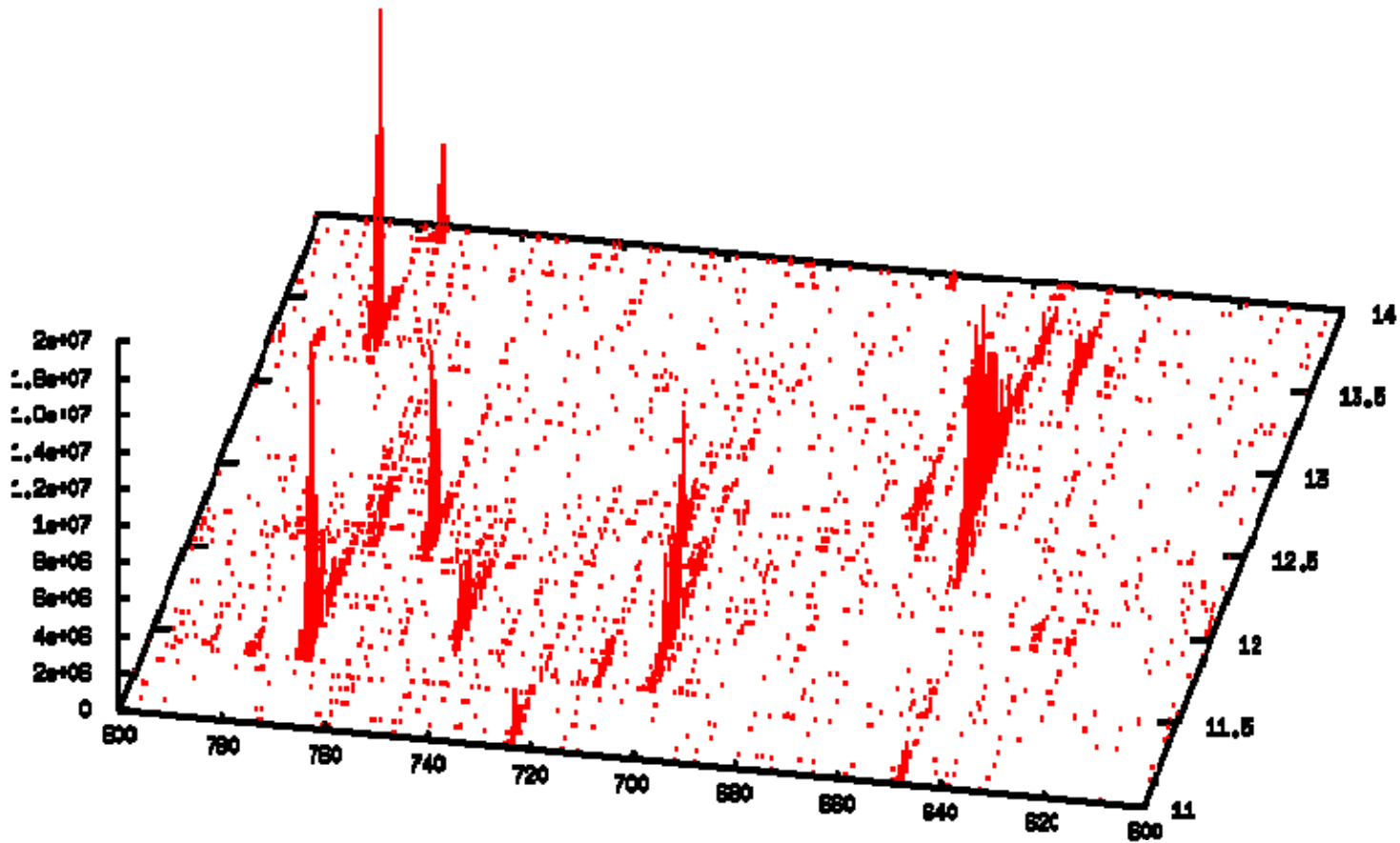
Maps

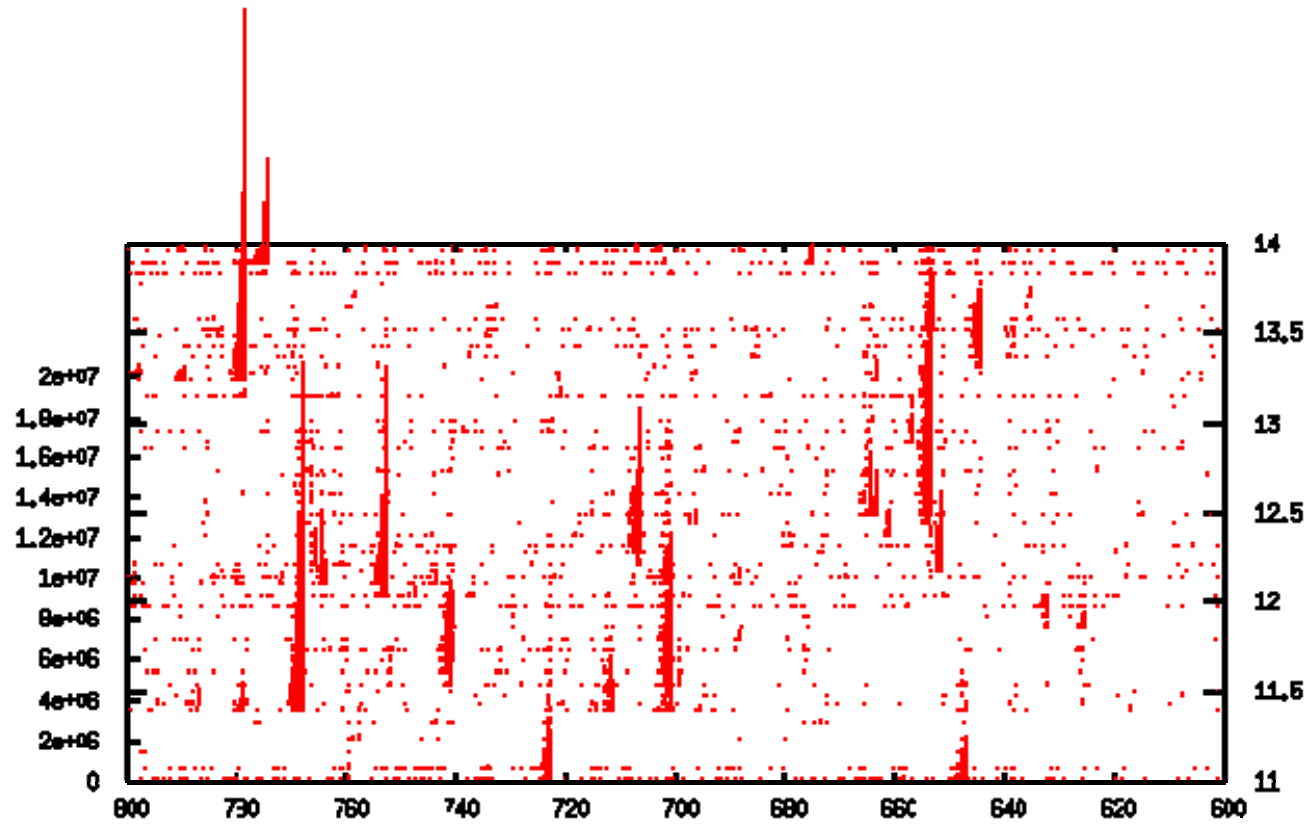


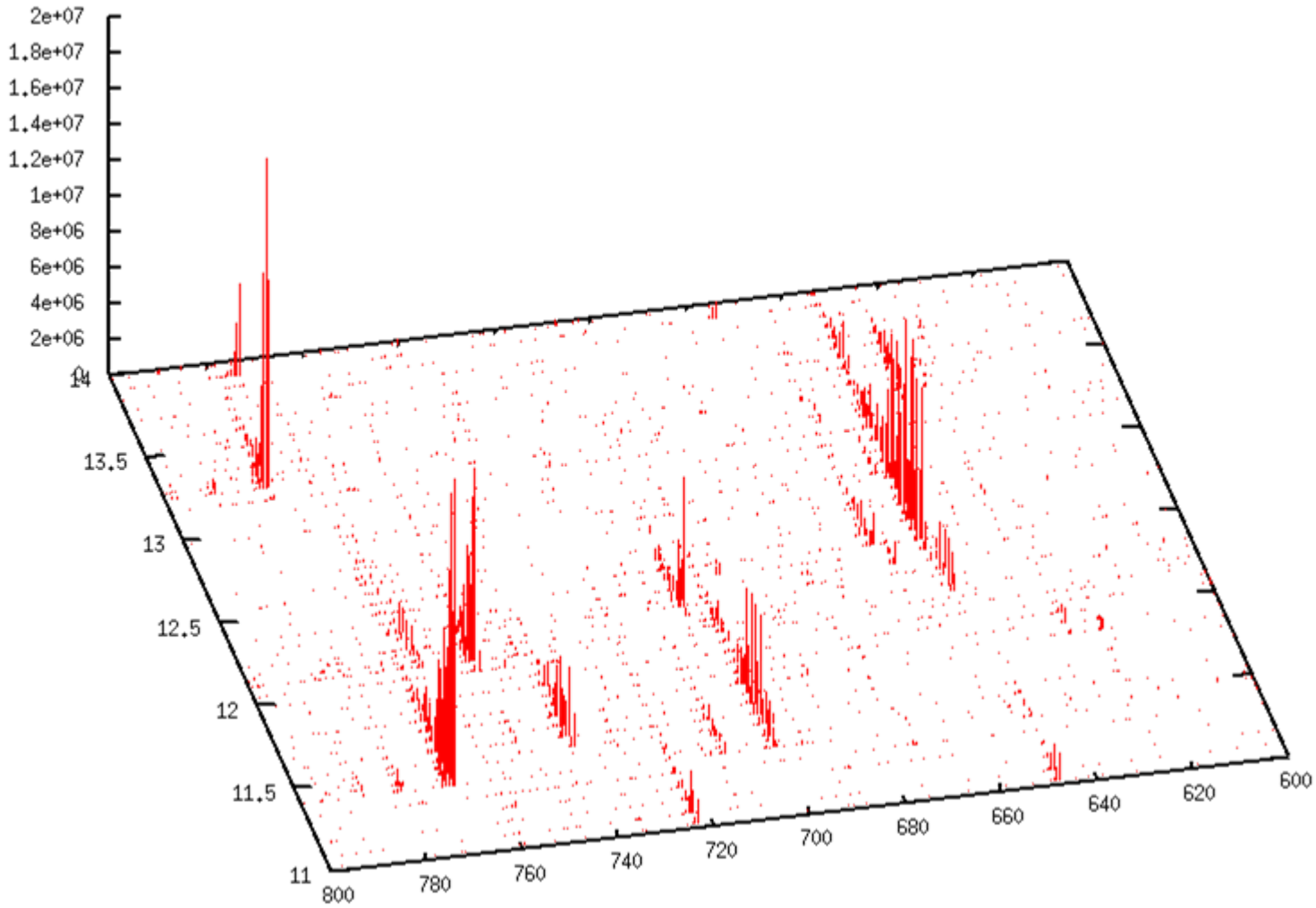


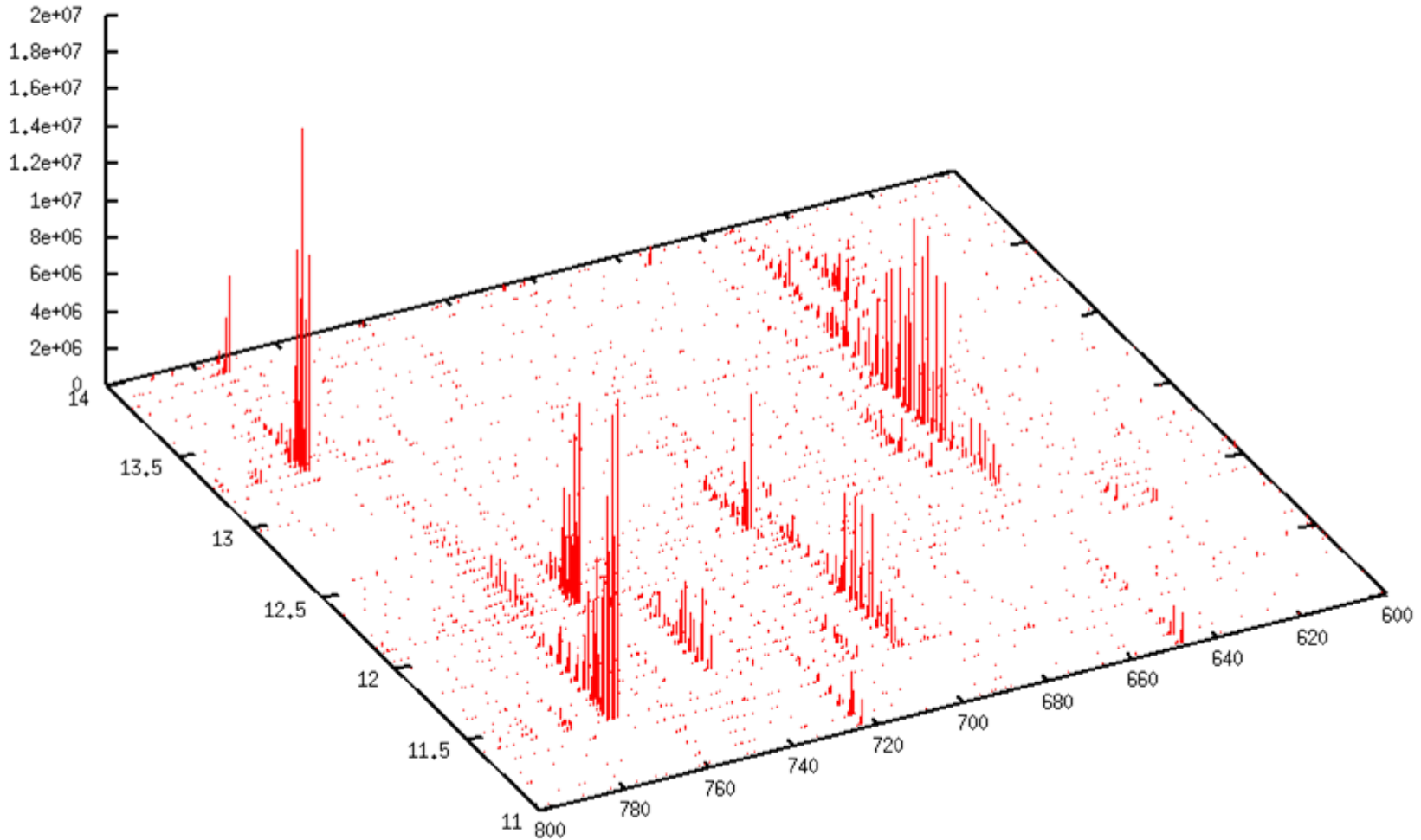




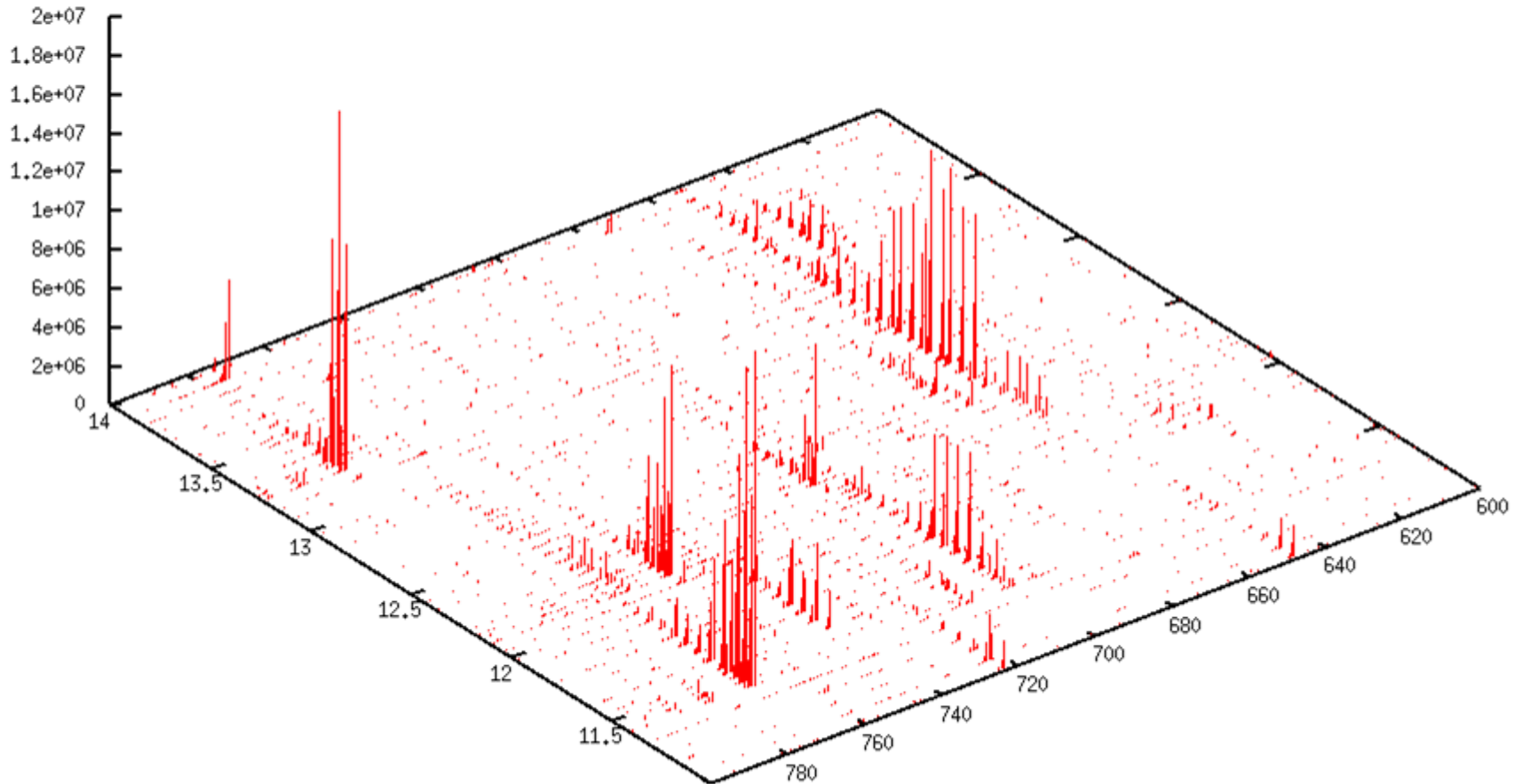


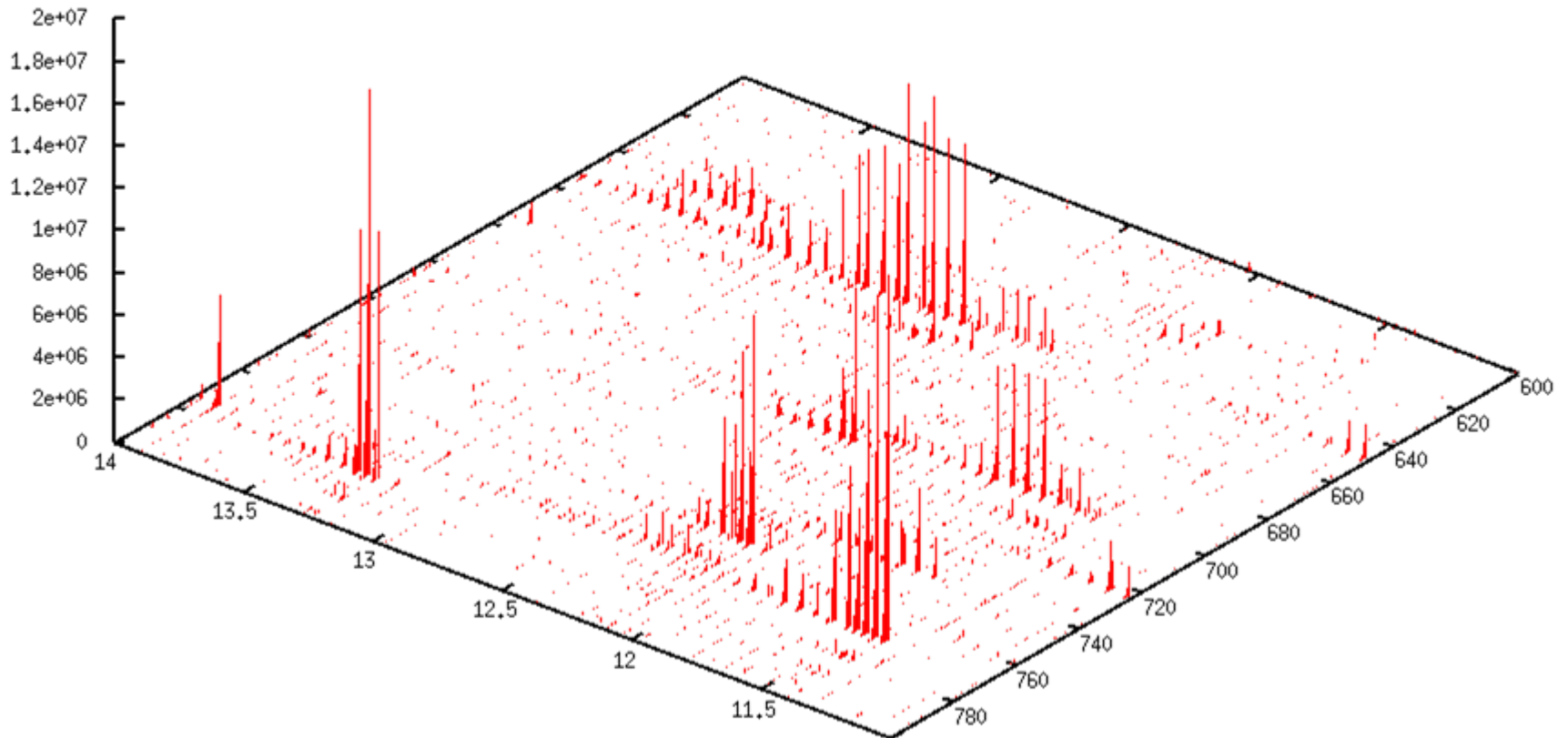


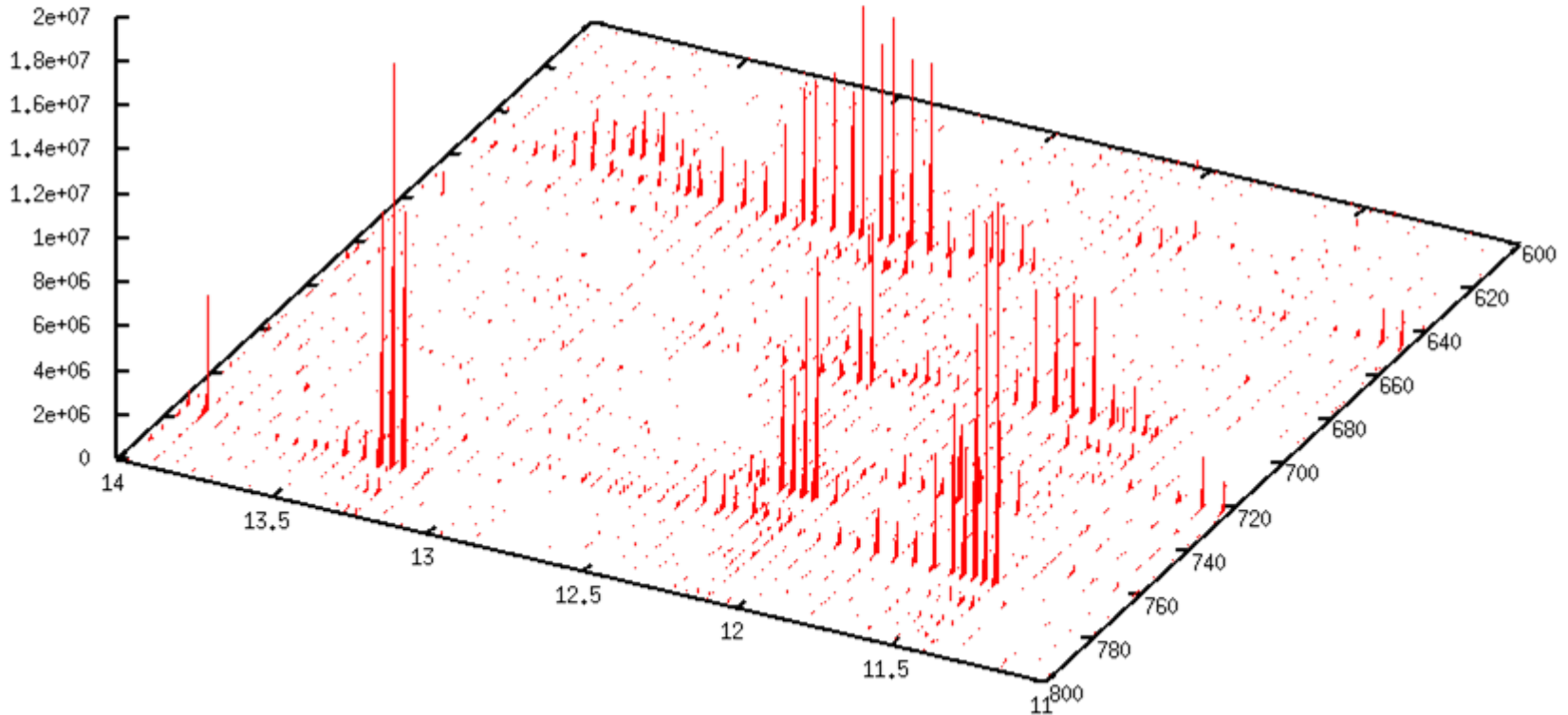


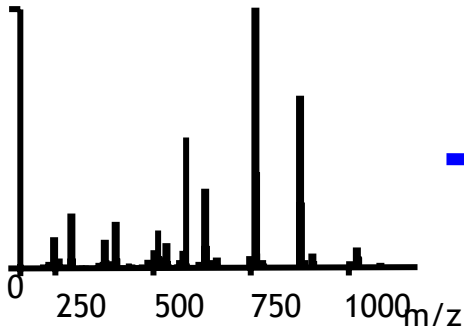


Maps



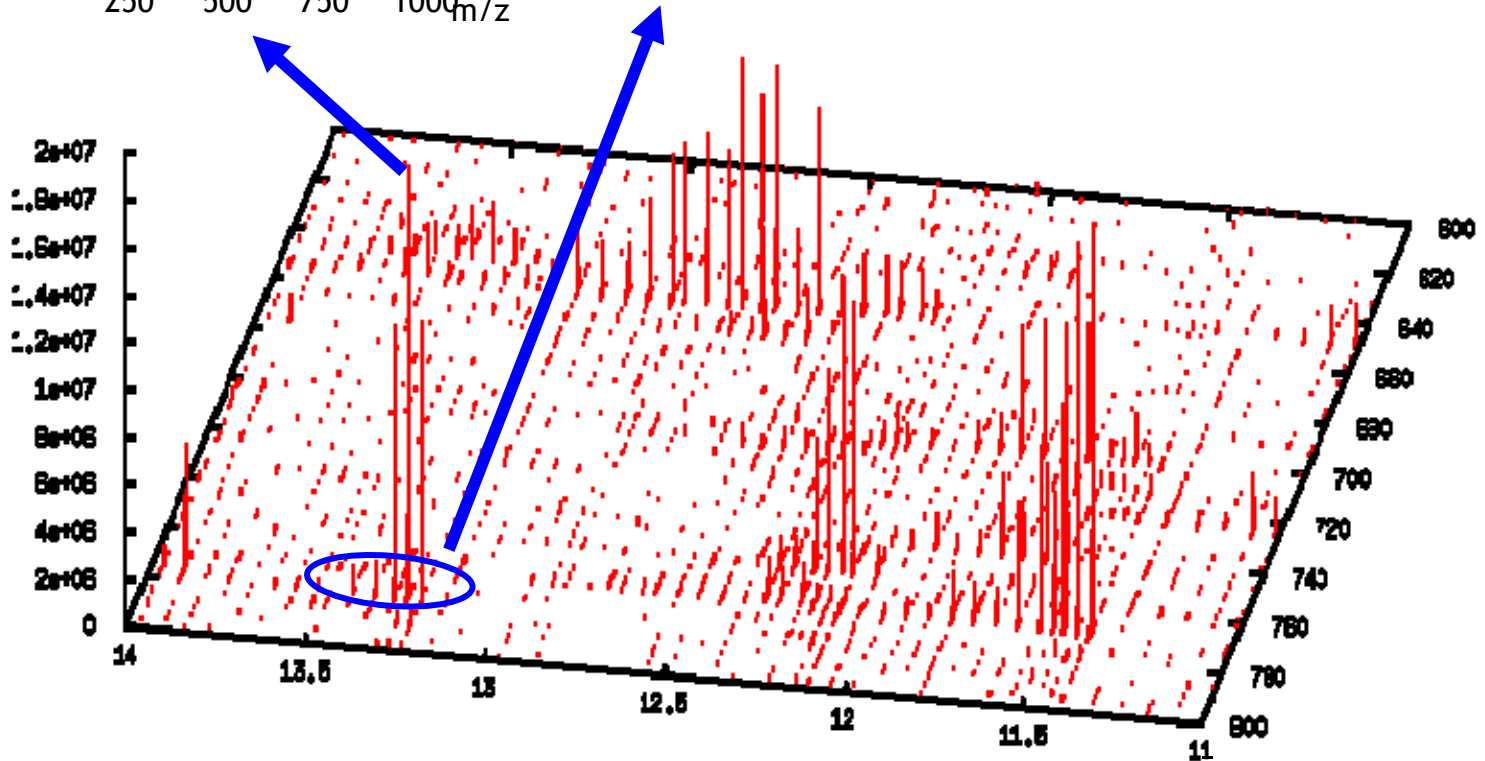






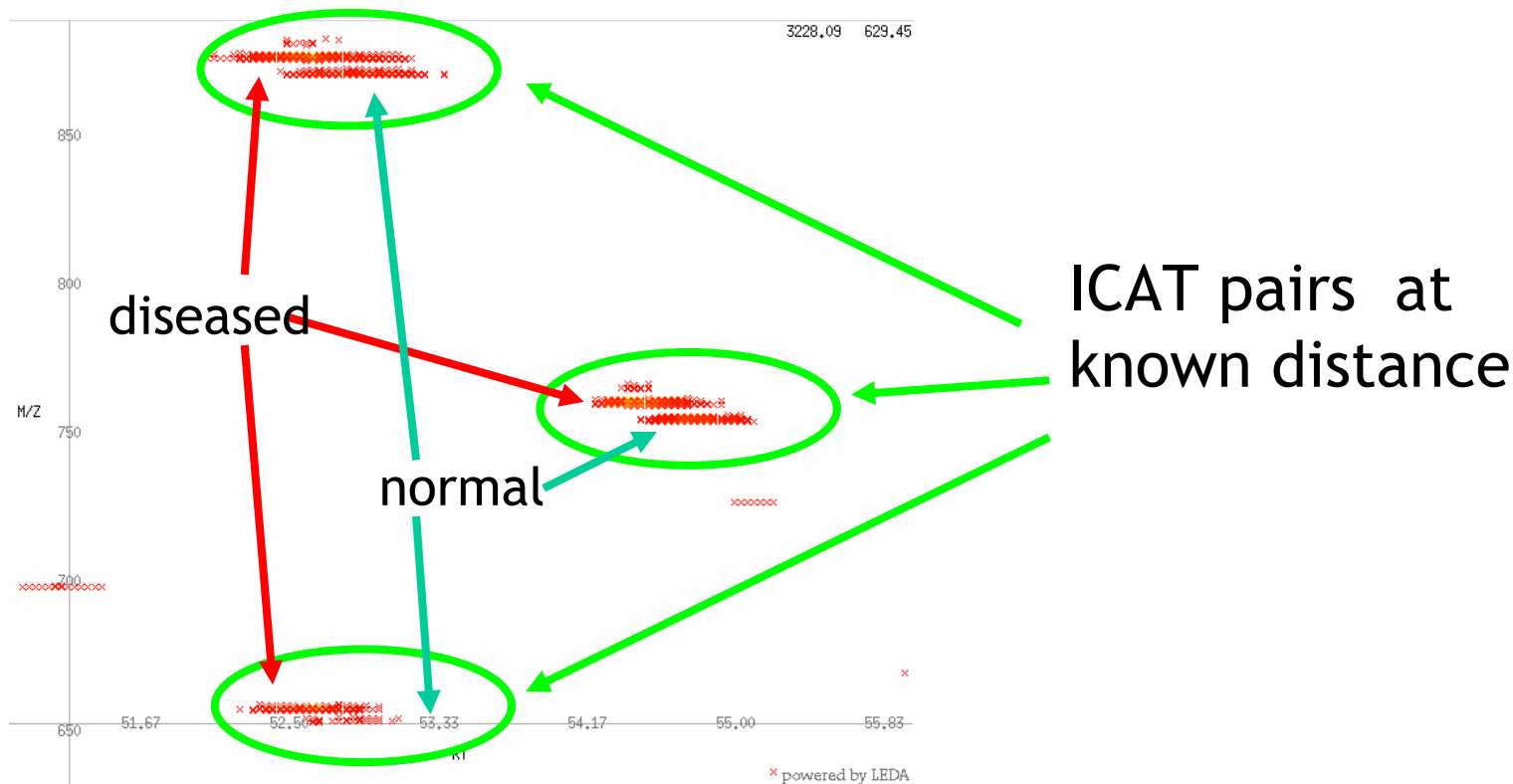
EVAFAQFGSDLASTK

15 nmol/ μ l, 3x overexpressed, ...



Two common basic approaches:

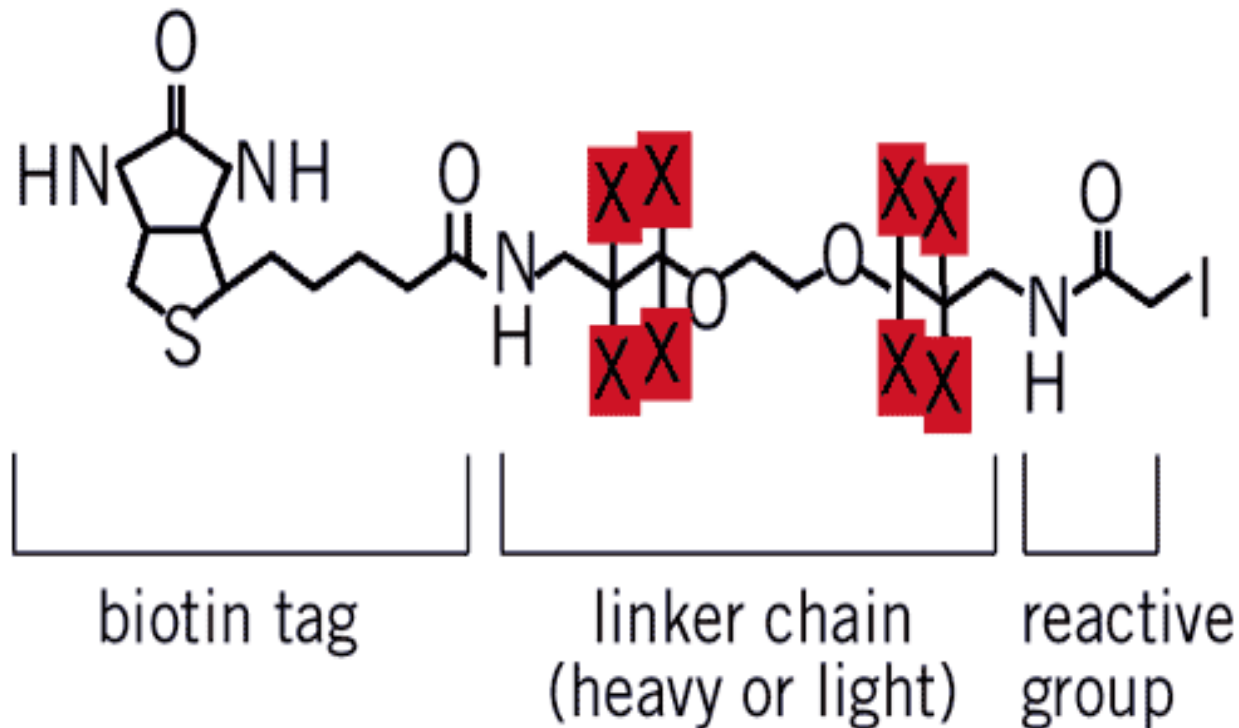
- Isotope labeling (e.g. ICAT, MeCAT, SILAC,...)
- Direct Differential Quantification (DDQ)



Isotope-Coded Affinity Tags

heavy reagent: D8-ICAT Reagent (X=deuterium)

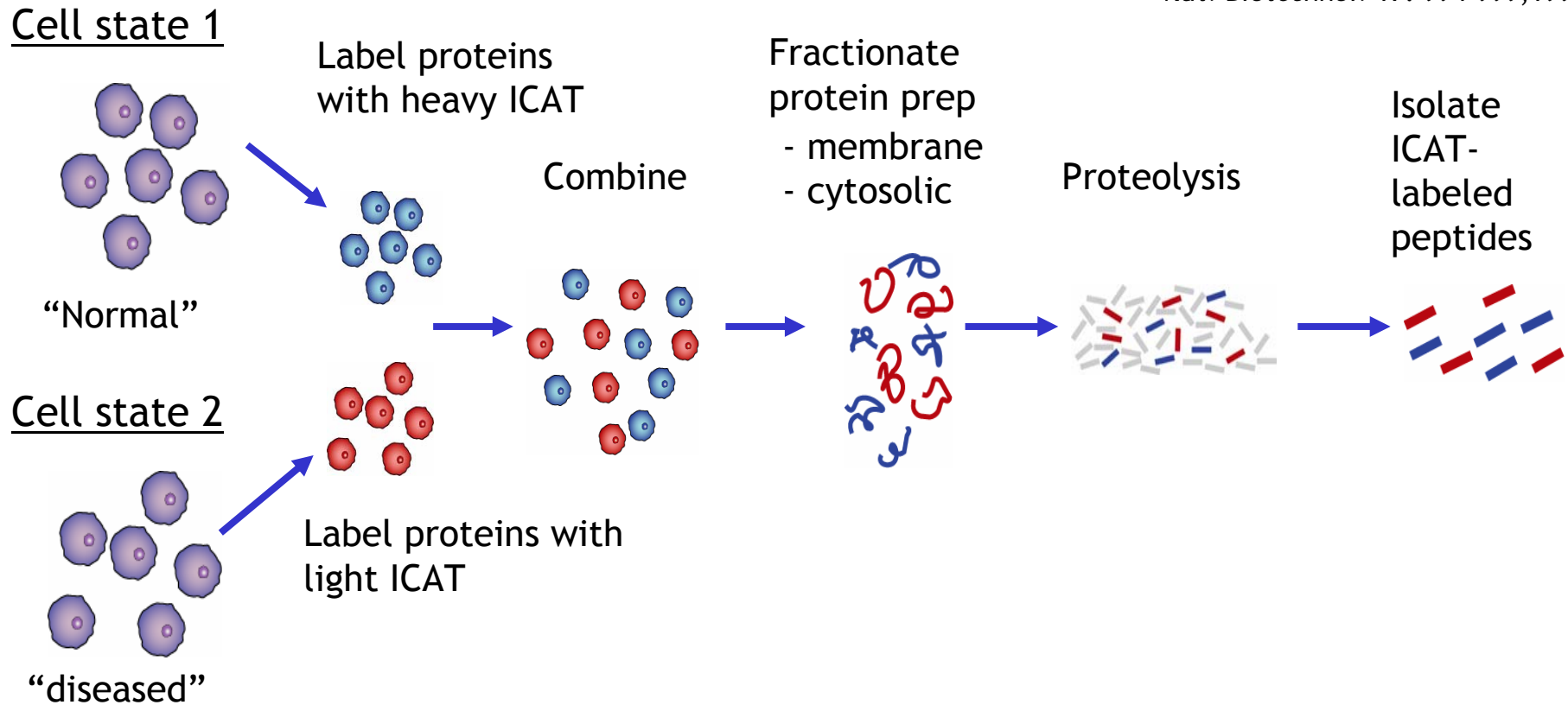
light reagent: D0-ICAT Reagent (X=hydrogen)



Isotope Labeling (ICAT)



Nat. Biotechnol. 17: 994-999,1999

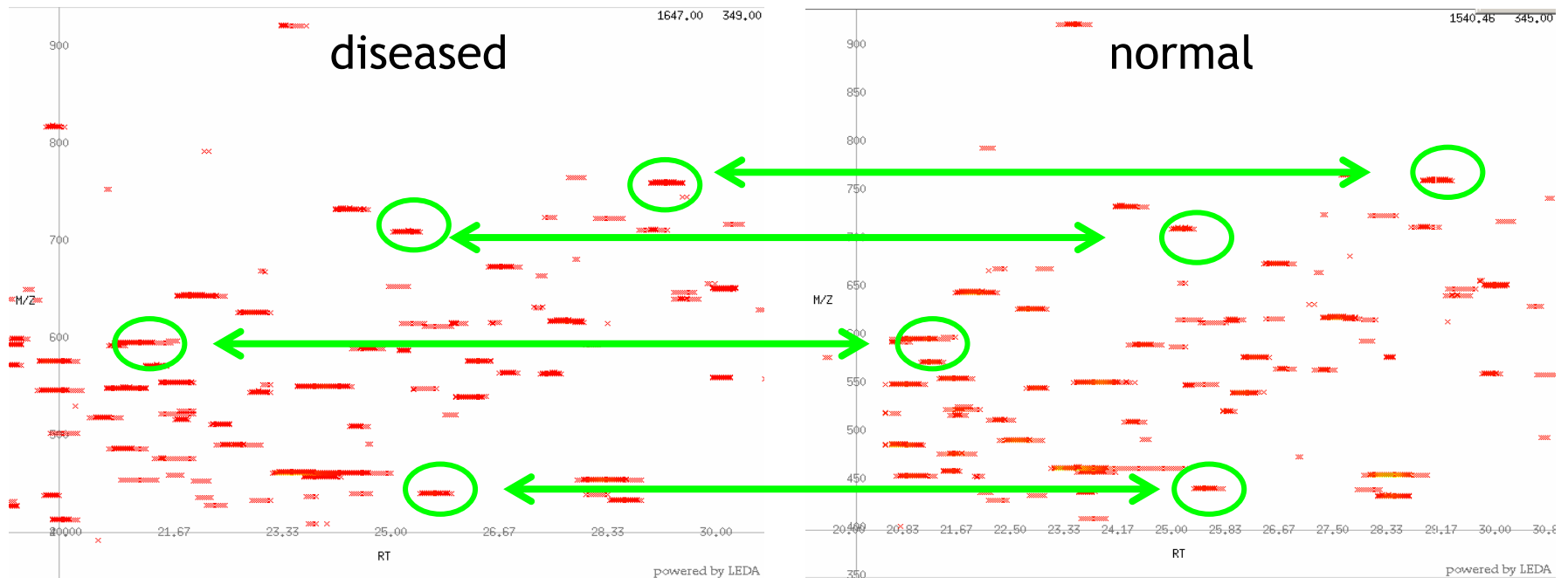


Heavy and light ICAT reagent 8 Dalton apart



Two common basic approaches:

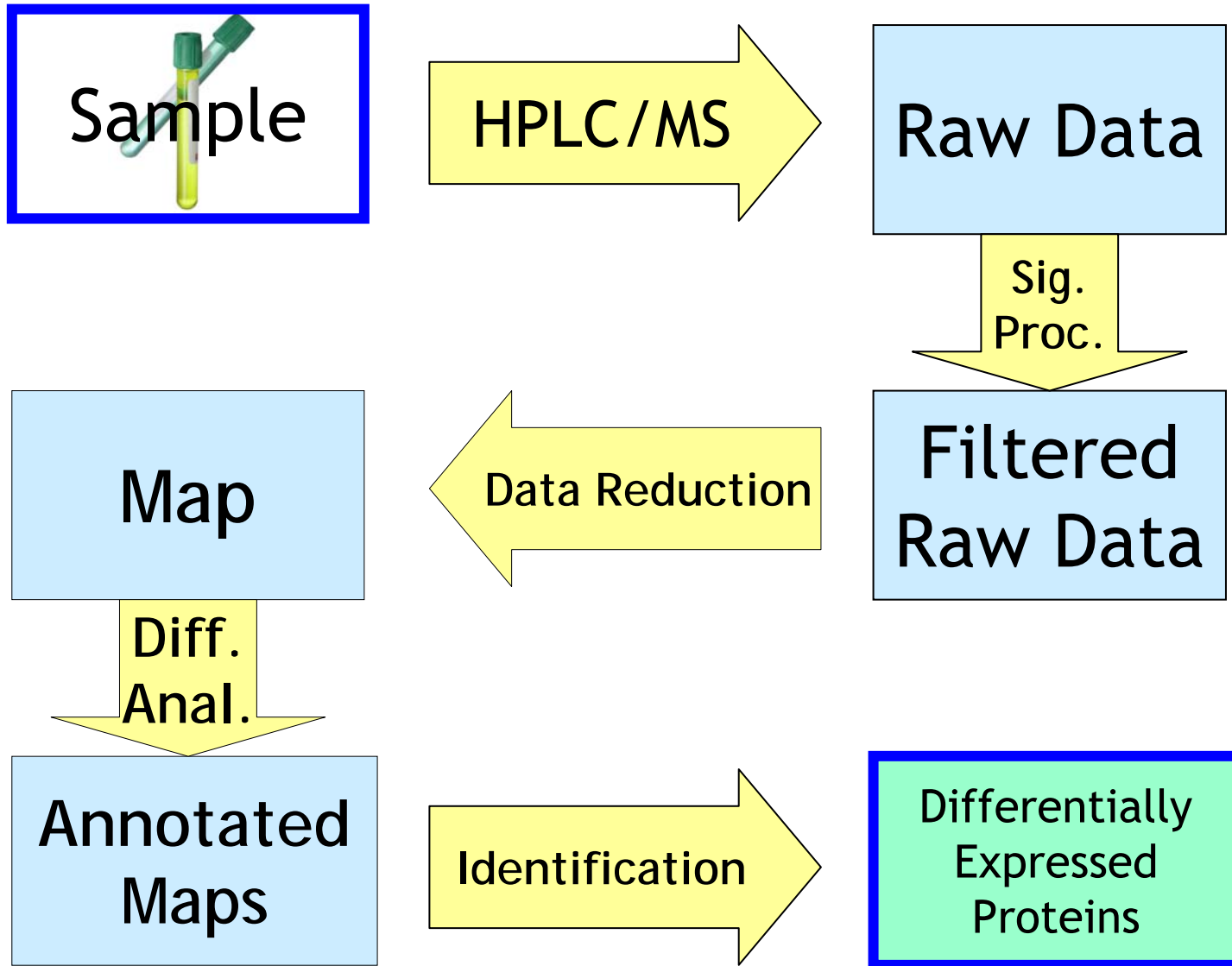
- Isotope tagging (e.g. ICAT, MeCAT)
- **Direct Differential Quantitation (DDQ)**



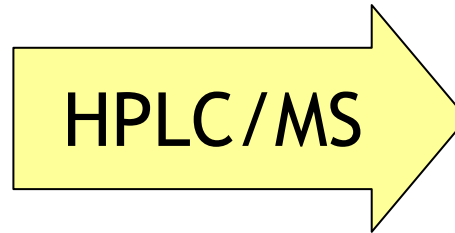
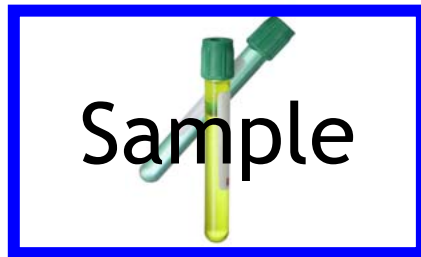
- Retention time (RT) for each scan
- Peptide mass/charge ratio (m/z)
(usually within ~ 20 ppm)
- Intensity (I)
 - \Rightarrow use m/z , RT to identify peptides
 - \Rightarrow use I to quantify peptides
(relative quantitation only!)

Maps become HUGE (10^8 Peaks)!

Proteomics Data Flow



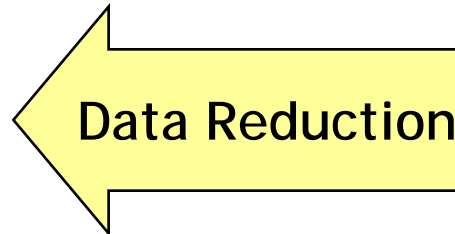
Proteomics Data Flow



10 GB



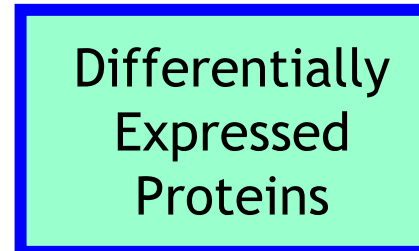
1 GB



50 MB



50 MB



1 kB

Advantages of HPLC/MS over 2D PAGE

1. Easier automation

- no robots
- no gel handling

2. Better separation power

- Monolithic columns
- Multi-dimensional chromatography

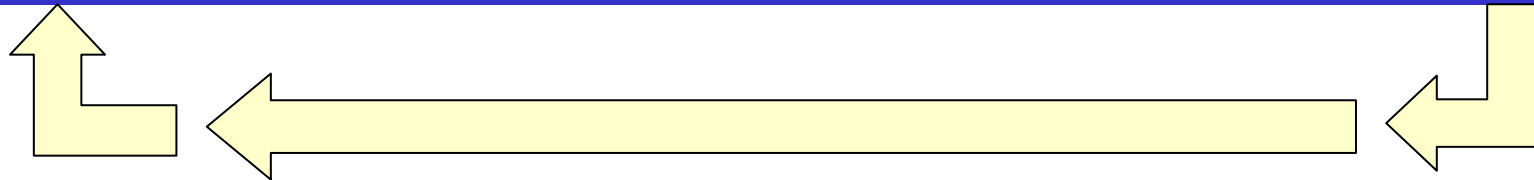
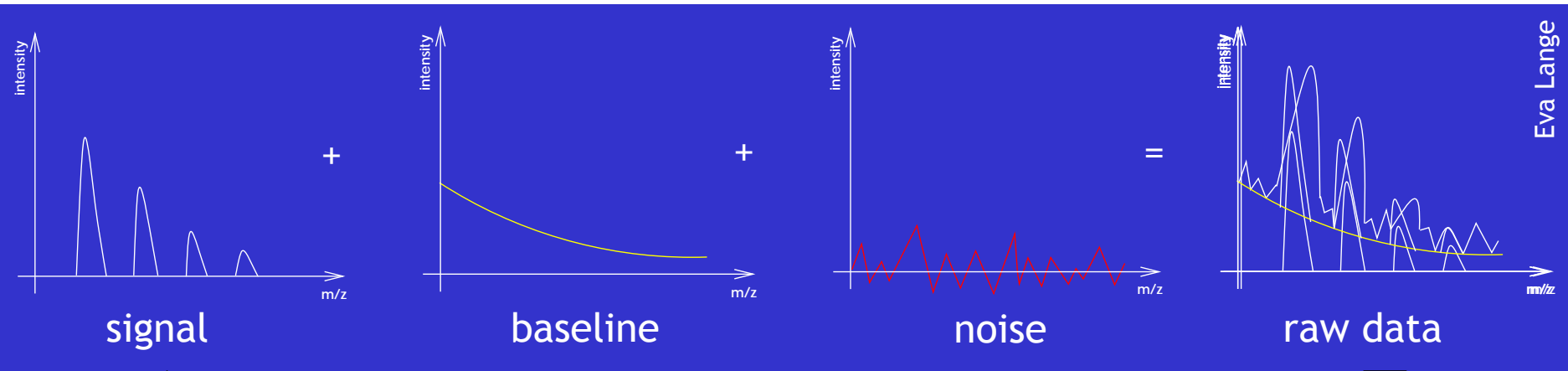
3. Simpler coupling to MS

Disadvantages of HPLC/MS

- Proteins are chopped up
- Quantitation difficult
- Huge amount of data (10 GiB/run)
- Data hard to manage/interpret

⇒ Need for Informatics!

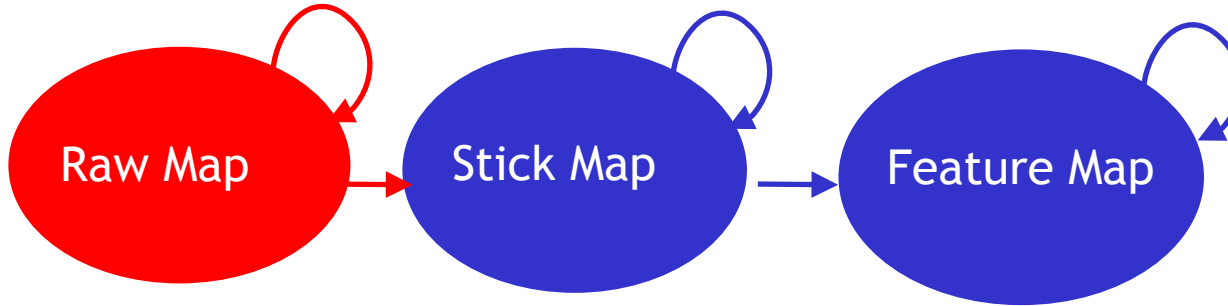
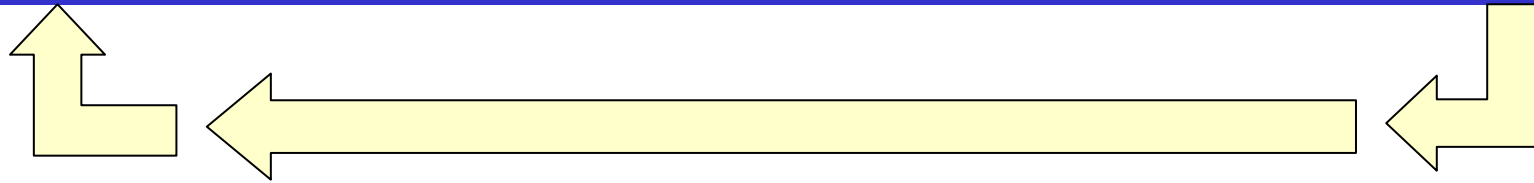
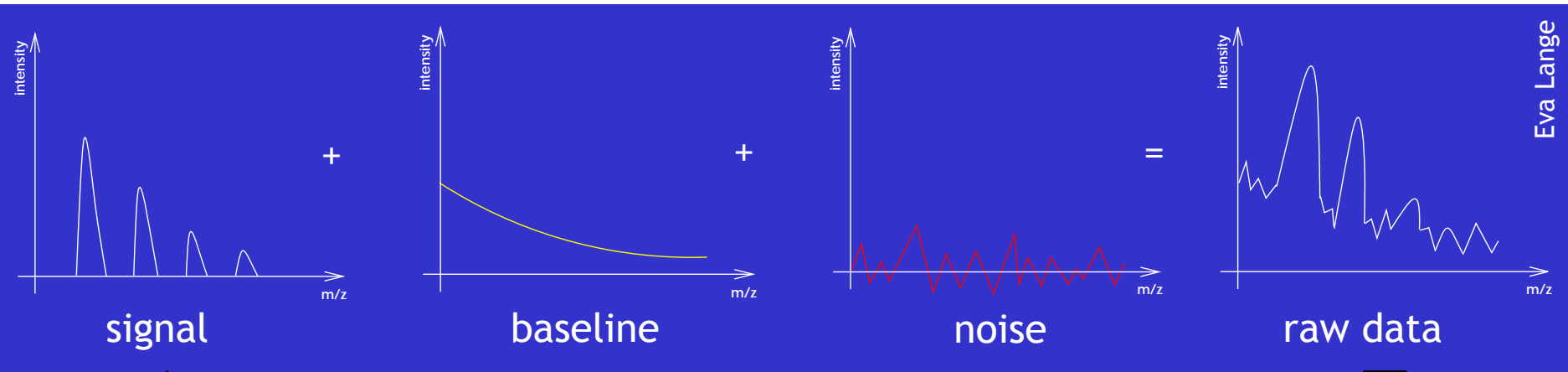
- Motivation for OpenMS
- **Bioinformatics Issues in quantitative Proteomics**
 - **Signal Processing**
 - Feature Finding
 - Map Mapping
 - Differential Quantitation
 - Identification, Clustering
 - Software Engineering, Databases

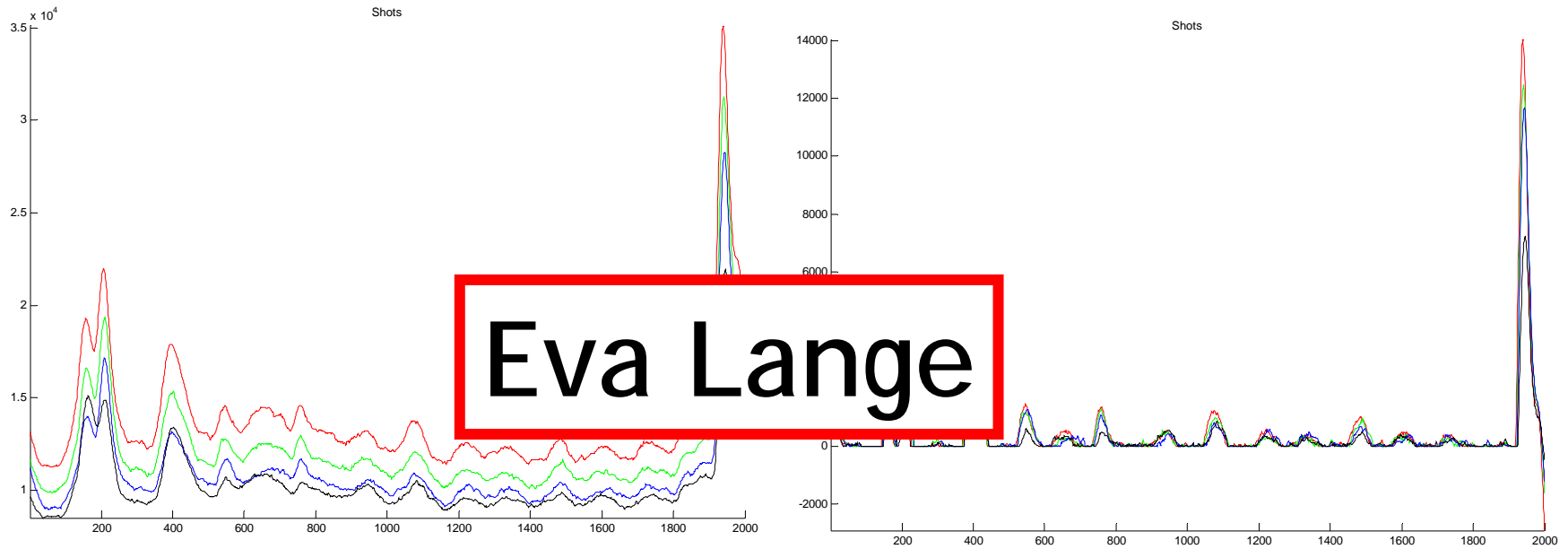


Issues

- Baseline reduction
- Noise filtering
- Calibration

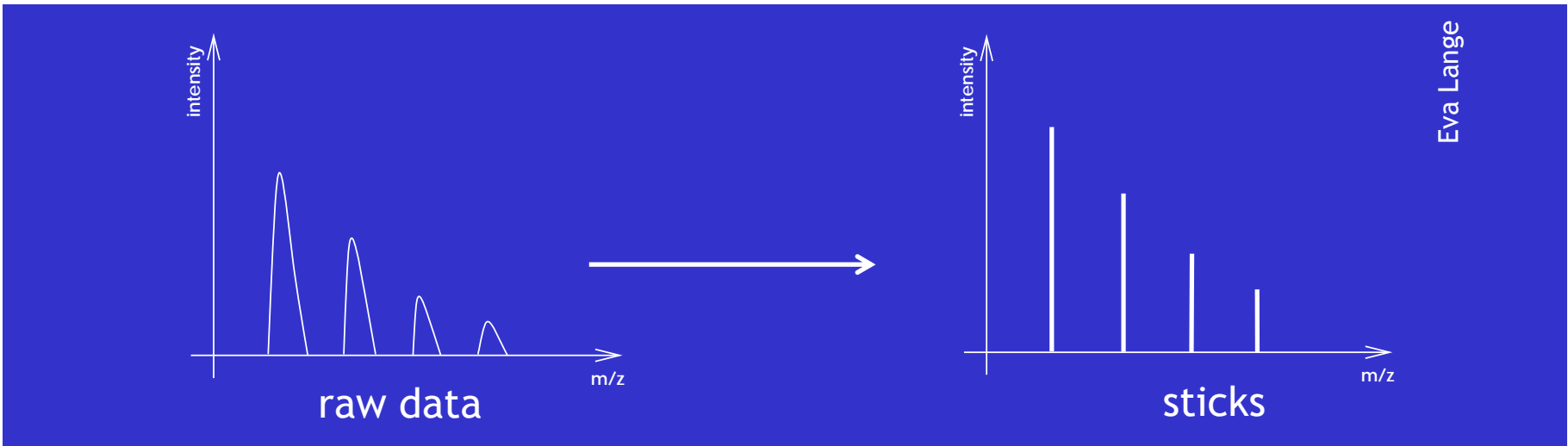
⇒ Increase reliability of data





TopHat filtering determines baseline

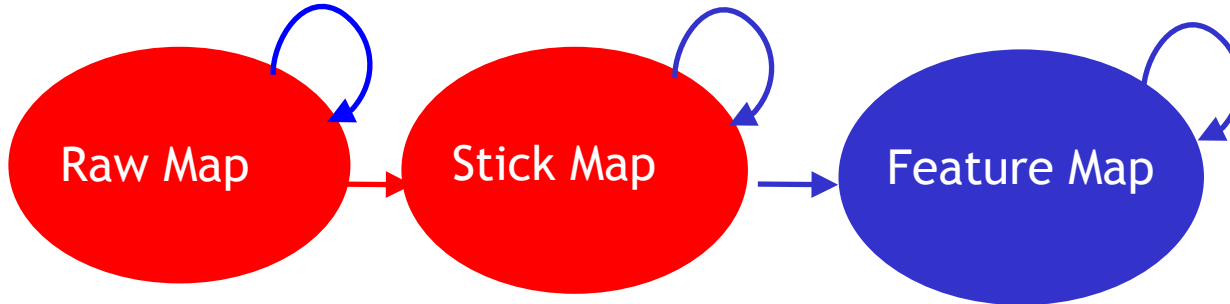
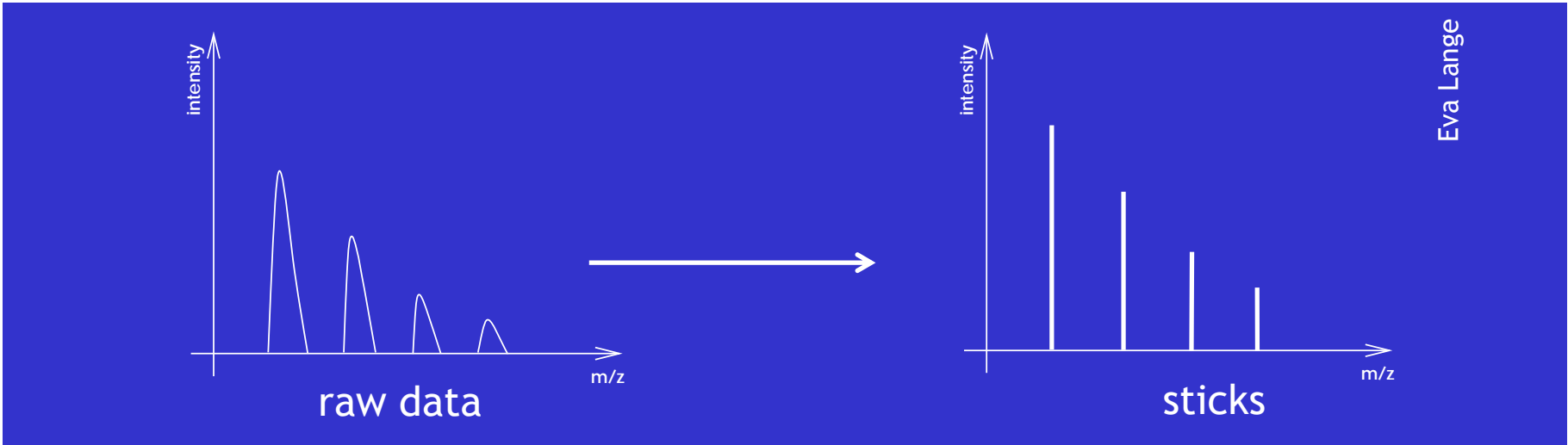
Stick conversion aka Peak picking



Issues

- „peak picking“ - Identify peaks
 - „stick conversion“ - Integrate peaks to sticks
- ⇒ **Reduce amount of data** by factor of 10 - 100

Stick conversion aka Peak picking

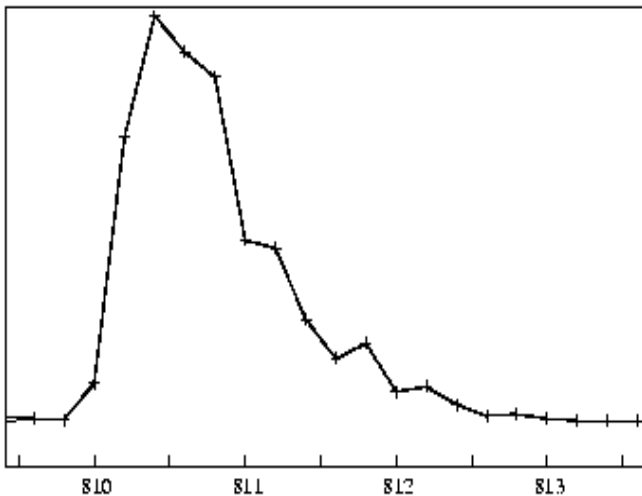
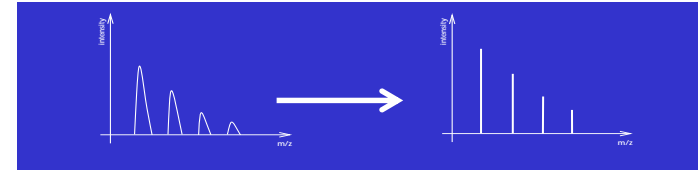


Peak picking in m/z dimension



Main objectives of a peak picking algorithm:

- precise peak positions
- run in real time

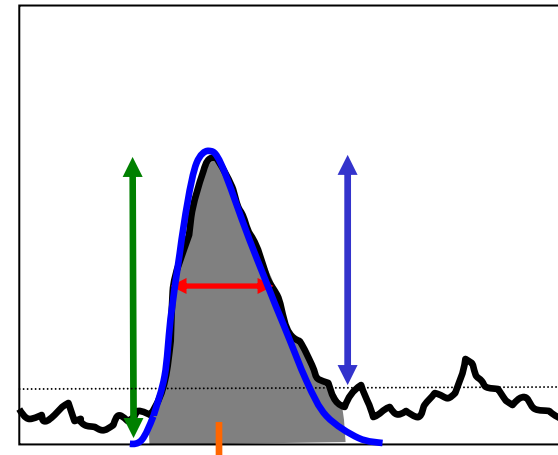


Main difficulties of a peak picking algorithm:

- considerable asymmetry of peaks
- convolution of isotopic peaks

1. Peak detection
2. Extract important peak parameter:

- centroid
- area
- height
- fwhm
- asymmetric peak shape
- signal/noise ratio

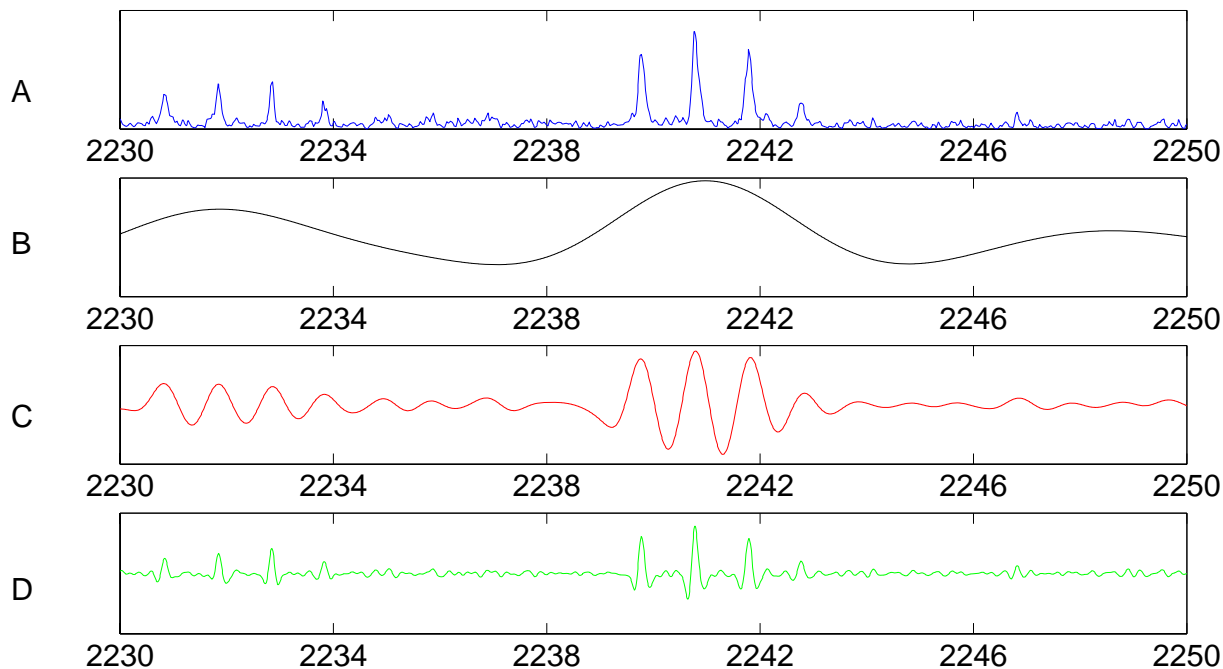


3. Optimize the peak parameter (optional)

1. Peak detection

Idea: Split the signal into different frequency ranges
(length scales)

Method: Continuous Wavelet Transformation (CWT)

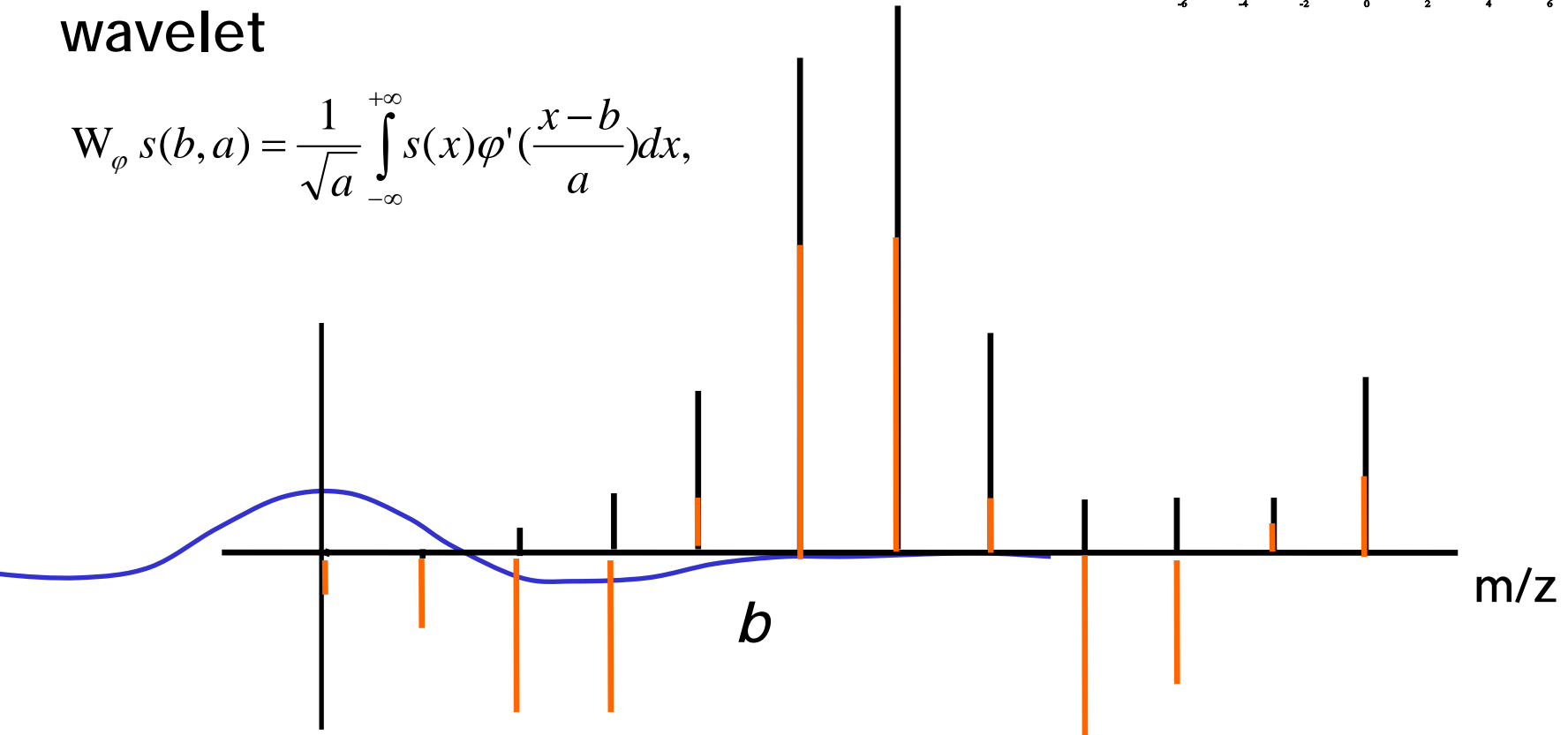
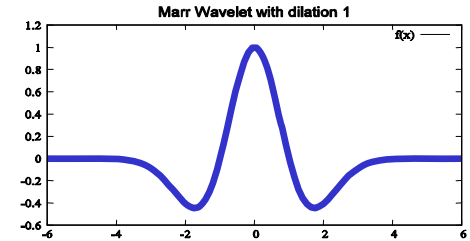


Continuous wavelet transformation



CWT is the sum over all positions b of the original (shifted) wavelet multiplied by the shifted and scaled version of the signal

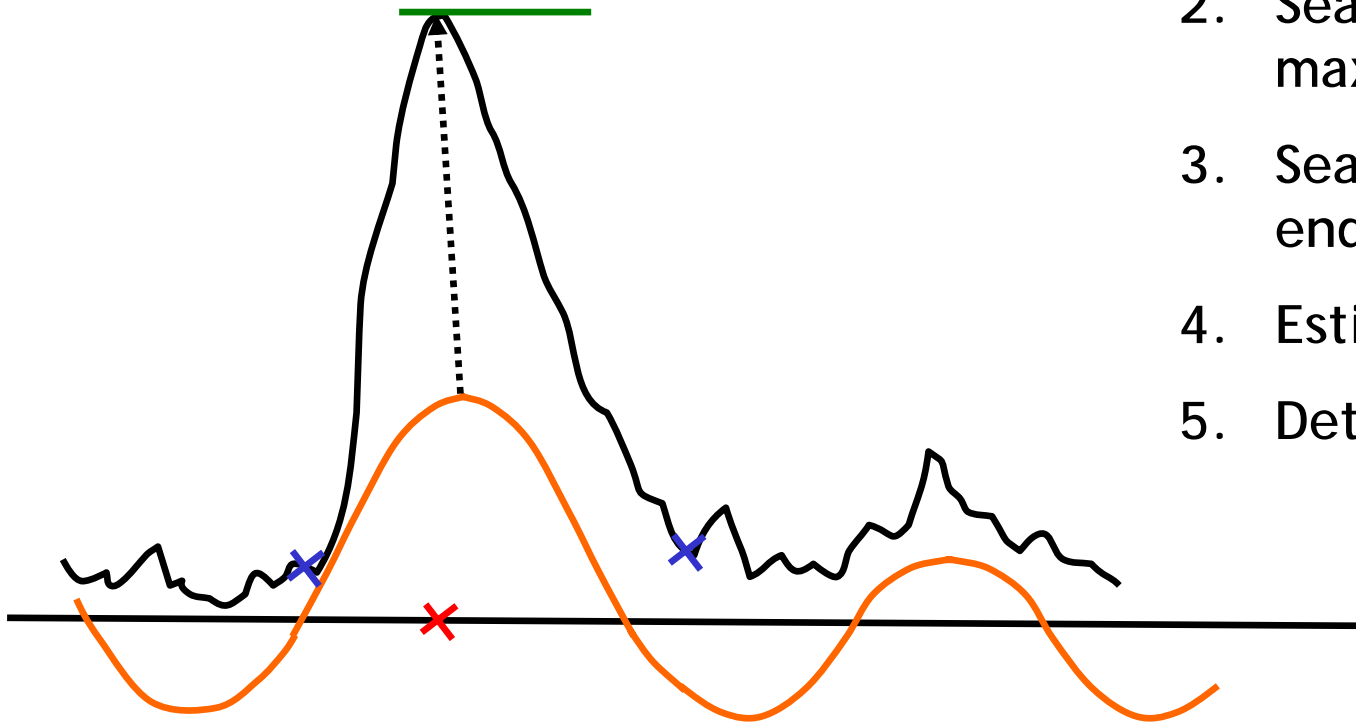
$$W_{\varphi} s(b, a) = \frac{1}{\sqrt{a}} \int_{-\infty}^{+\infty} s(x) \varphi' \left(\frac{x-b}{a} \right) dx,$$





Workflow

1. Compute the wavelet transform
2. Search for a peak's maximum
3. Search for peak-endpoints
4. Estimate the centroid
5. Determine the height



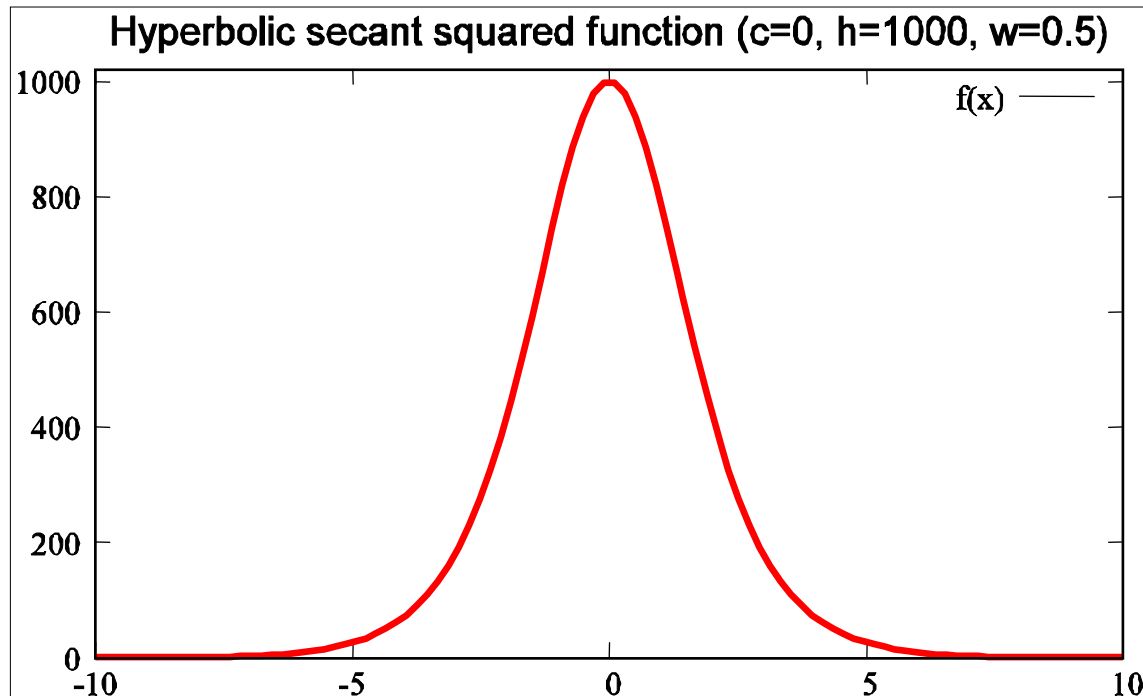
1. Compute the peak positions

Compute the CWT:

- use the marr wavelet
- predefined scale a (should correspond to the typical peak width)
- ✓ maximum position in the CWT is a good first estimate of the maximum in the data even for asymmetric peaks
- ✓ convolution can be computed very efficiently with pre-tabulated values for the wavelet

2. Extract important peak parameter

Fit functions similar to real raw mass peaks

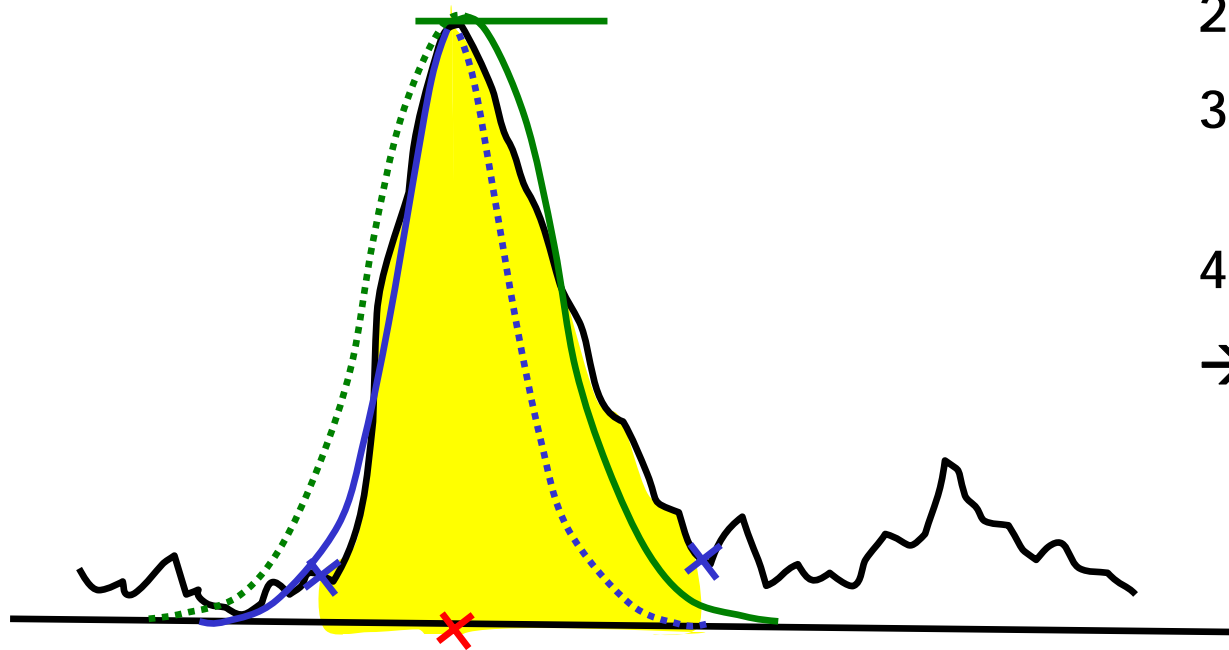


$$\text{Sech}^2_{h,w,x_0}(x) = \frac{h}{\cosh(w(x - x_0))^2}$$

2. Extract important peak parameter



Workflow



1. Estimate the peak's left area
 2. Fit a sech^2 to the left
 3. Estimate the peak's right area
 4. Fit a sech^2 to the right
- Asymmetric peak shape

2. Extract important peak parameter

Fitting asymmetric lorentzian and sech² functions to a raw peak using the peak's :

- maximum position
 - height
 - area
-
- ✓ introduces a smoothing effect
 - ✓ compute the peak's width from its area in constant time
 - ✓ yield very good approximations to the original peak shape

Peak picking algorithm

```
for all mass spectra  $s$  in experiment do  
   $box\_list = \text{splitSpectrum}(s)$   
  for all boxes  $b$  in  $box\_list$  do  
     $peak\_list = []$   
     $w_b = \text{continuousWaveletTransform}(b)$   
    while  $\text{getNextMaximumPosition}(w_b, b, x_0)$  do  
       $(x_l, x_r) = \text{searchForPeakEndPoints}(b, x_0)$   
       $c = \text{estimateCentroid}(x_l, x_r)$   
       $h = \text{intensity}(x_0)$   
       $(A_l, A_r) = \text{integrateAreas}(x_l, x_r)$   
       $f = \text{fitPeakShape}(A_l, A_r, x_0, h)$   
       $\text{push}(f, peak\_list)$   
       $\text{removeRawData}(x_l, x_r, b)$   
       $w_b = \text{continuousWaveletTransform}(b)$   
    end while  
     $\text{optimizePeakParameter}(peak\_list, b)$   
  end for  
end for
```

Peak picking algorithm

for all mass spectra s **in** *experiment* **do**

$box_list = \text{splitSpectrum}(s)$

for all boxes b **in** box_list **do**

$peak_list = []$

$w_b = \text{continuousWaveletTransform}(b)$

while $\text{getNextMaximumPosition}(w_b, b, x_0)$ **do**

$(x_l, x_r) = \text{searchForPeakEndPoints}(b, x_0)$

$c = \text{estimateCentroid}(x_l, x_r)$

$h = \text{intensity}(x_0)$

$(A_l, A_r) = \text{integrateAreas}(x_l, x_r)$

$f = \text{fitPeakShape}(A_l, A_r, x_0, h)$

$\text{push}(f, peak_list)$

$\text{removeRawData}(x_l, x_r, b)$

$w_b = \text{continuousWaveletTransform}(b)$

end while

$\text{optimizePeakParameter}(peak_list, b)$

end for

end for

Peak picking algorithm

```
for all mass spectra  $s$  in experiment do  
   $box\_list = \text{splitSpectrum}(s)$   
  for all boxes  $b$  in  $box\_list$  do  
     $peak\_list = []$   
     $w_b = \text{continuousWaveletTransform}(b)$   
    while  $\text{getNextMaximumPosition}(w_b, b, x_0)$  do  
       $(x_l, x_r) = \text{searchForPeakEndPoints}(b, x_0)$   
       $c = \text{estimateCentroid}(x_l, x_r)$   
       $h = \text{intensity}(x_0)$   
       $(A_l, A_r) = \text{integrateAreas}(x_l, x_r)$   
       $f = \text{fitPeakShape}(A_l, A_r, x_0, h)$   
       $\text{push}(f, peak\_list)$   
       $\text{removeRawData}(x_l, x_r, b)$   
       $w_b = \text{continuousWaveletTransform}(b)$   
    end while  
     $\text{optimizePeakParameter}(peak\_list, b)$   
  end for  
end for
```

Peak picking algorithm

```
for all mass spectra  $s$  in experiment do  
   $box\_list = \text{splitSpectrum}(s)$   
  for all boxes  $b$  in  $box\_list$  do  
     $peak\_list = []$   
     $w_b = \text{continuousWaveletTransform}(b)$   
    while  $\text{getNextMaximumPosition}(w_b, b, x_0)$  do  
       $(x_l, x_r) = \text{searchForPeakEndPoints}(b, x_0)$   
       $c = \text{estimateCentroid}(x_l, x_r)$   
       $h = \text{intensity}(x_0)$   
       $(A_l, A_r) = \text{integrateAreas}(x_l, x_r)$   
       $f = \text{fitPeakShape}(A_l, A_r, x_0, h)$   
       $\text{push}(f, peak\_list)$   
       $\text{removeRawData}(x_l, x_r, b)$   
       $w_b = \text{continuousWaveletTransform}(b)$   
    end while  
     $\text{optimizePeakParameter}(peak\_list, b)$   
  end for  
end for
```


Peak picking algorithm

```
for all mass spectra  $s$  in experiment do  
   $box\_list = \text{splitSpectrum}(s)$   
  for all boxes  $b$  in  $box\_list$  do  
     $peak\_list = []$   
     $w_b = \text{continuousWaveletTransform}(b)$   
    while  $\text{getNextMaximumPosition}(w_b, b, x_0)$  do  
       $(x_l, x_r) = \text{searchForPeakEndPoints}(b, x_0)$   
       $c = \text{estimateCentroid}(x_l, x_r)$   
       $h = \text{intensity}(x_0)$   
       $(A_l, A_r) = \text{integrateAreas}(x_l, x_r)$   
       $f = \text{fitPeakShape}(A_l, A_r, x_0, h)$   
       $\text{push}(f, peak\_list)$   
       $\text{removeRawData}(x_l, x_r, b)$   
       $w_b = \text{continuousWaveletTransform}(b)$   
    end while  
     $\text{optimizePeakParameter}(peak\_list, b)$   
  end for  
end for
```

Compute accurate:

- centroid
- area
- height

→ test against a set of known composition

Results are heavily affected by:

- quality of experimental data
- calibration method

Evaluation of our peak picking method

Data set 1: peptide mix (peptide standards mix #P2693 from Sigma Aldrich) of nine known peptides.

Measuring method:

HPLC-MS using a quadrupole ion trap mass spectrometer equipped with an electrospray ion source in full scan mode (m/z 500-1500)

→ low resolution ($\Delta m=0.2$) and several overlapping isotope patterns

Evaluation of our peak picking method

Evaluation:

- compare with vendor supplied Bruker Data-Analysis software (3.2) on the same spectra
- no sophisticated calibration, we only allow a constant mass offset to keep the number of fit parameters as small as possible

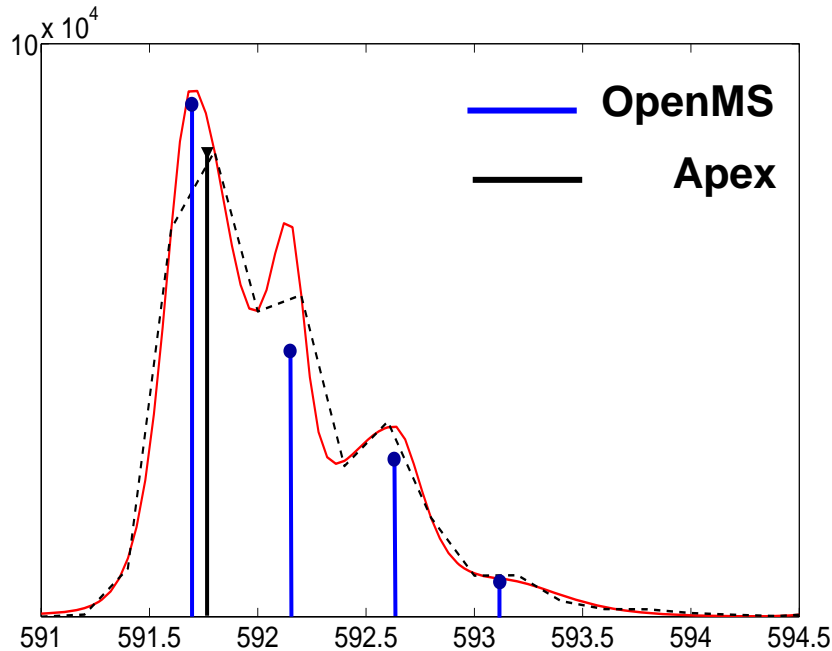
Performance criteria:

1. number of the discovered and separated isotopic patterns (given by at least three consecutive peaks) of each peptide in an expected retention time interval
2. average relative error of the monoisotopic mass

Results: LC-ESI-ion trap data

(Full scan mode (m/z 500-1500))

Data set I (Peptide mix of nine known peptides)



	mass _{theo}	Z	#occ	
			OpenMS	Apex
A	556.2693	1	22	19
B	573.3071	1	29	29
C	574.2257	1	19	15
D	1007.4365	1	8	5
E	1084.4379	1	3	2
F	1061.5614	2	3	2
G	1183.5730	2	7	0
H	1349.7360	2	8	1
I	1620.8151	2	13	0

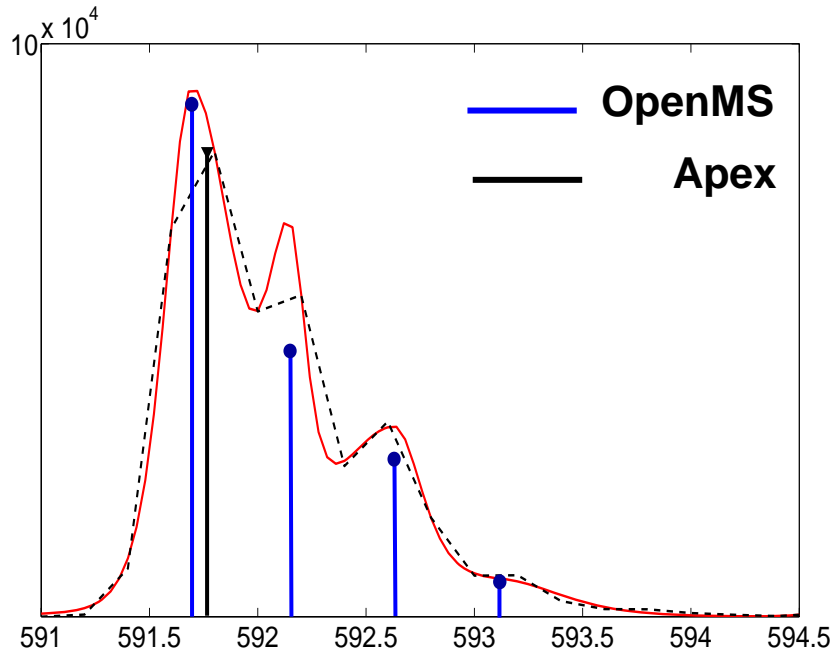
Relative error of centroids:

21ppm, 1.2ppm, 35ppm, 16ppm

Results: LC-ESI-ion trap data

(Full scan mode (m/z 500-1500))

Data set I (Peptide mix of nine known peptides)



	mass _{theo}	Z	#occ	
			OpenMS	Apex
A	556.2693	1	22	19
B	573.3071	1	29	29
C	574.2257	1	19	15
D	1007.4365	1	8	5
E	1084.4379	1	3	2
F	1061.5614	2	3	2
G	1183.5730	2	7	0
H	1349.7360	2	8	1
I	1620.8151	2	13	0

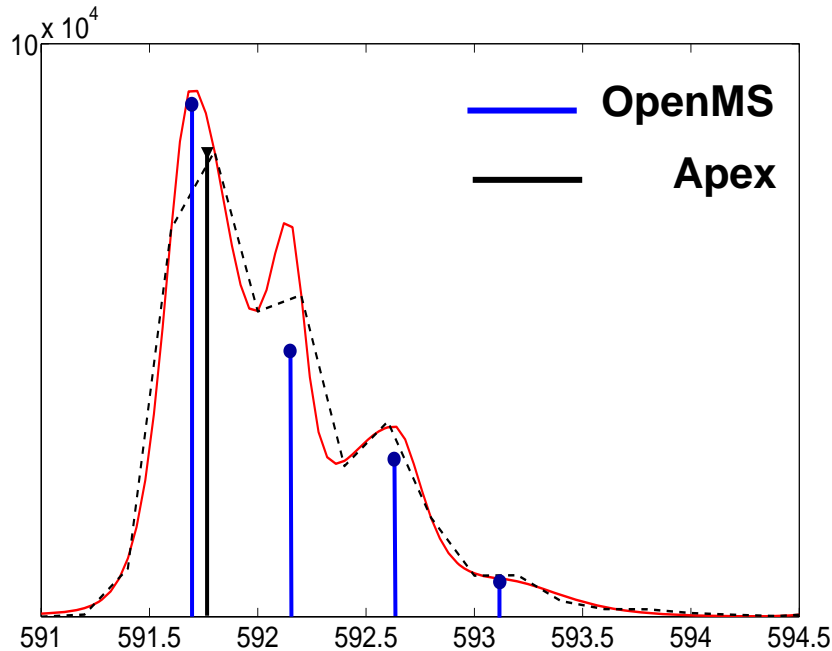
Relative error of centroids:

21ppm, 1.2ppm, 35ppm, 16ppm

Results: LC-ESI-ion trap data

(Full scan mode (m/z 500-1500))

Data set I (Peptide mix of nine known peptides)



	mass _{theo}	Z	#occ	
			OpenMS	Apex
A	556.2693	1	22	19
B	573.3071	1	29	29
C	574.2257	1	19	15
D	1007.4365	1	8	5
E	1084.4379	1	3	2
F	1061.5614	2	3	2
G	1183.5730	2	7	0
H	1349.7360	2	8	1
I	1620.8151	2	13	0

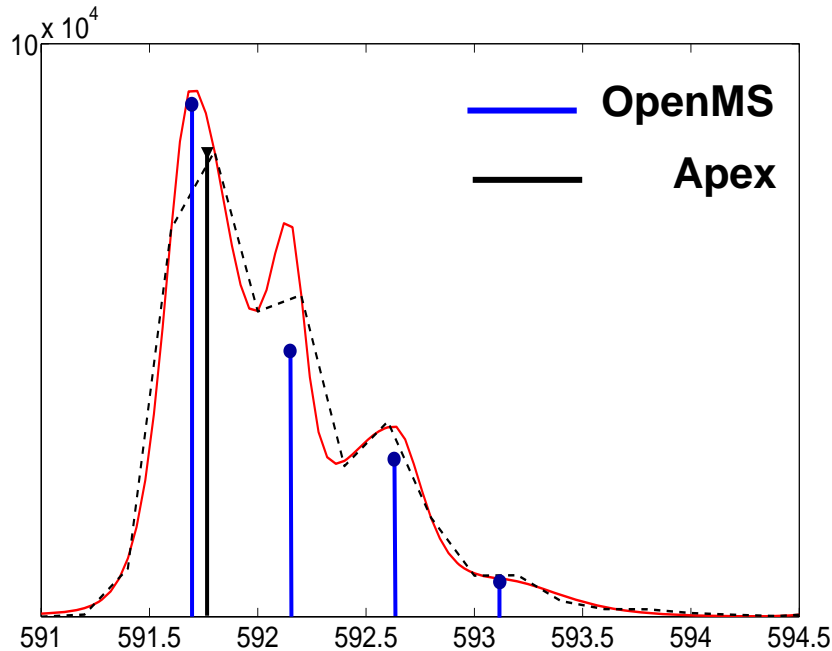
Relative error of centroids:

21ppm, 1.2ppm, 35ppm, 16ppm

Results: LC-ESI-ion trap data

(Full scan mode (m/z 500-1500))

Data set I (Peptide mix of nine known peptides)



	mass _{theo}	Z	rel. err. [ppm]	
			OpenMS	Apex
A	556.2693	1	31	39
B	573.3071	1	16	16
C	574.2257	1	44	21
D	1007.4365	1	25	94
E	1084.4379	1	18	12
F	1061.5614	2	56	64
G	1183.5730	2	15	-
H	1349.7360	2	28	13
I	1620.8151	2	37	-

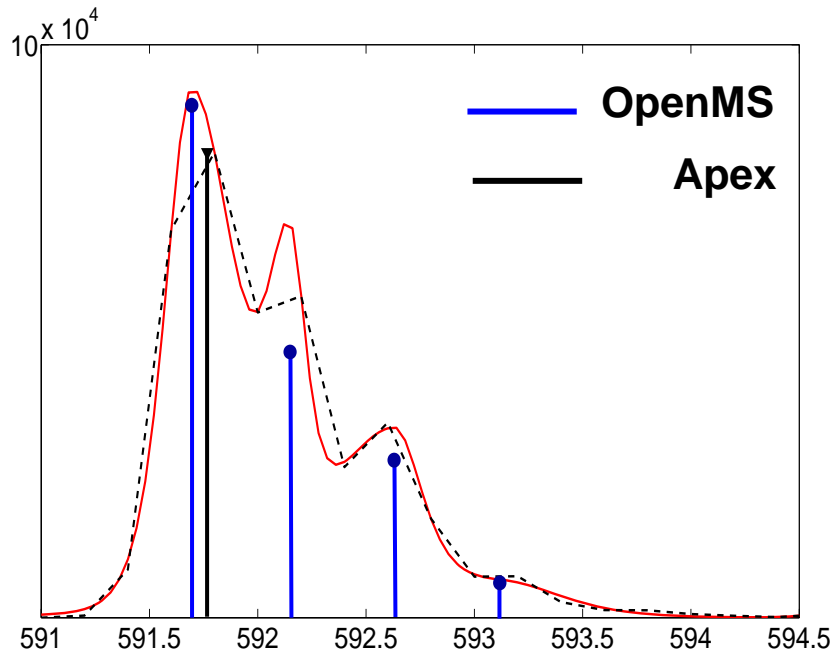
Relative error of centroids:

21ppm, 1.2ppm, 35ppm, 16ppm

Results: LC-ESI-ion trap data

(Full scan mode (m/z 500-1500))

Data set I (Peptide mix of nine known peptides)



	mass _{theo}	Z	rel. err. [ppm]	
			OpenMS	Apex
A	556.2693	1	31	39
B	573.3071	1	16	16
C	574.2257	1	44	21
D	1007.4365	1	25	94
E	1084.4379	1	18	12
F	1061.5614	2	56	64
G	1183.5730	2	15	-
H	1349.7360	2	28	13
I	1620.8151	2	37	-

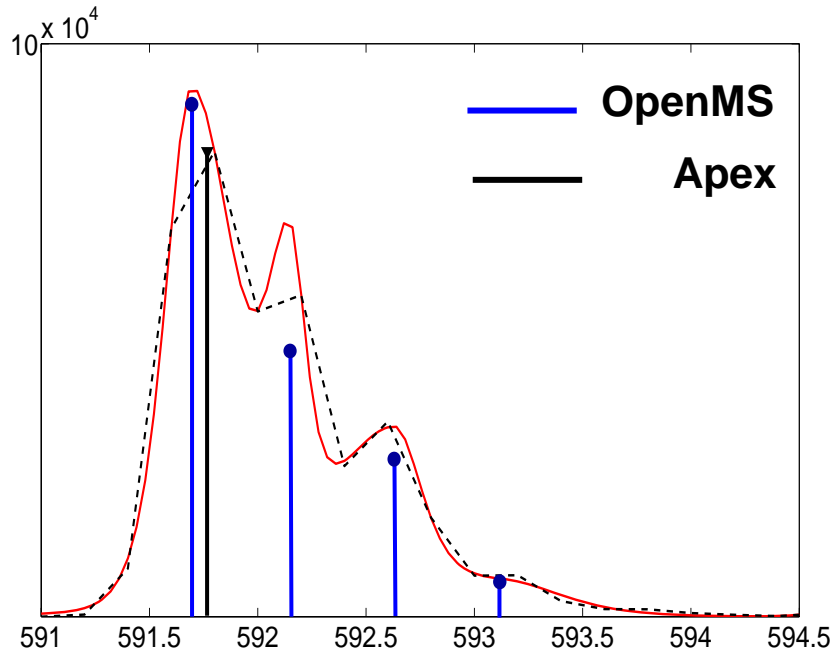
Relative error of centroids:

21ppm, 1.2ppm, 35ppm, 16ppm

Results: LC-ESI-ion trap data

(Full scan mode (m/z 500-1500))

Data set I (Peptide mix of nine known peptides)



	mass _{theo}	Z	rel. err. [ppm]	
			OpenMS	Apex
A	556.2693	1	31	39
B	573.3071	1	16	16
C	574.2257	1	44	21
D	1007.4365	1	25	94
E	1084.4379	1	18	12
F	1061.5614	2	56	64
G	1183.5730	2	15	-
H	1349.7360	2	28	13
I	1620.8151	2	37	-

Relative error of centroids:

21ppm, 1.2ppm, 35ppm, 16ppm

Sanity check of our peak picking method

Data set II: tryptic digest of bovine serum albumin (BSA from Sigma Aldrich).

Measuring method:

MALDI-TOF using a Ultraflex II Lift mass spectrometer operated in the reflectron mode and using Panorama delayed ion extraction (m/z 500-5000)

Sanity check of our peak picking method

Evaluation:

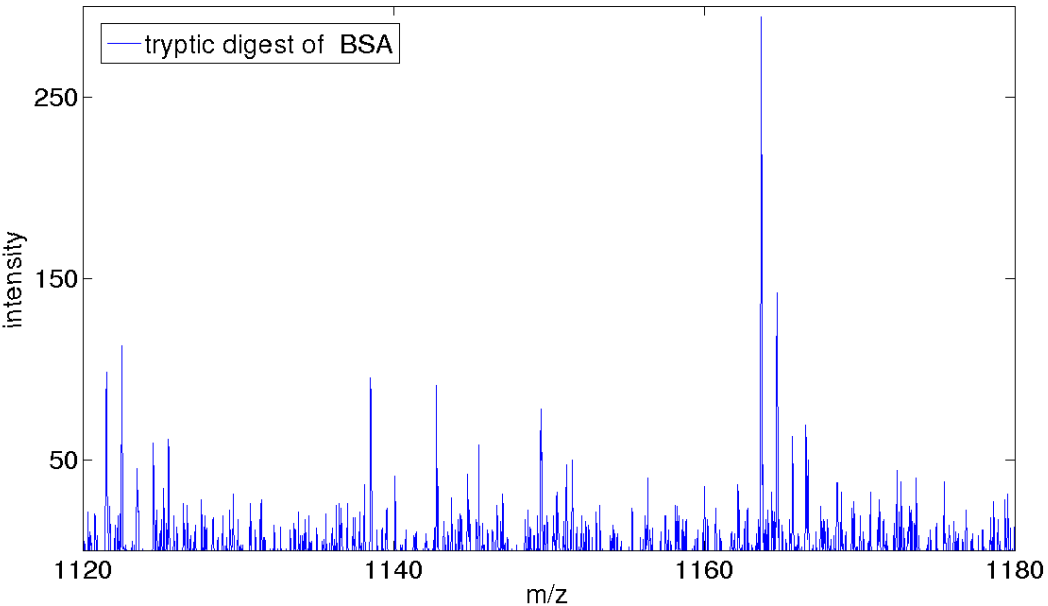
- compare with vendor supplied Bruker FlexAnalysis software on the same spectra
- no sophisticated calibration

Performance criterion:

1. compute the sequence coverage using the determined peak list as a MASCOT peptide mass fingerprinting query

Results: MALDI-TOF data

Data set II (tryptic digest of bovine serum)

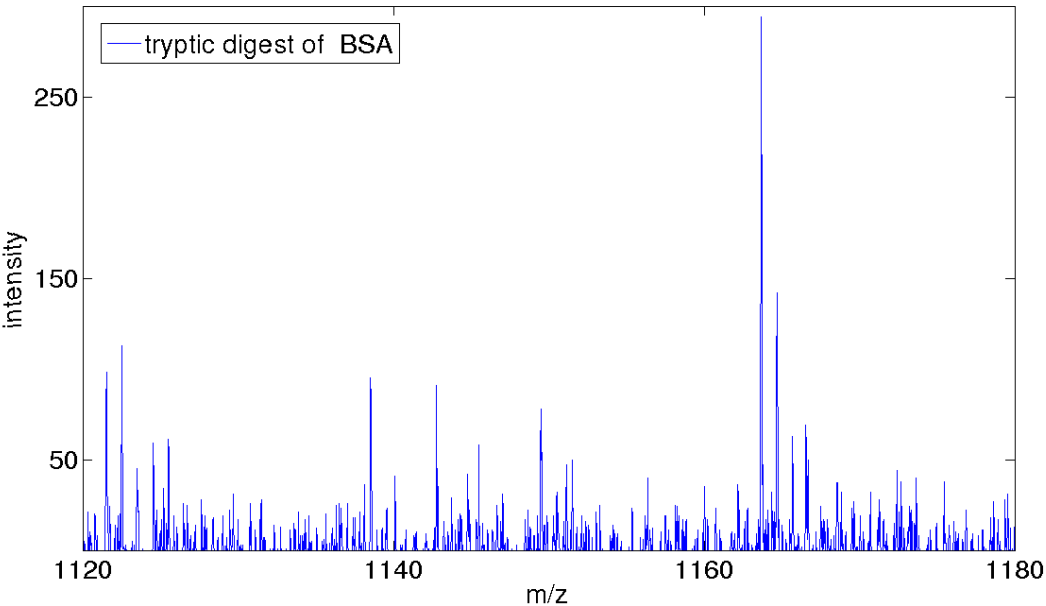


Sequence coverage	
OpenMS	Bruker
52% - 67%	44%

rel. err. [ppm]	
OpenMS	Bruker
80 - 93	95

Results: MALDI-TOF data

Data set II (tryptic digest of bovine serum)

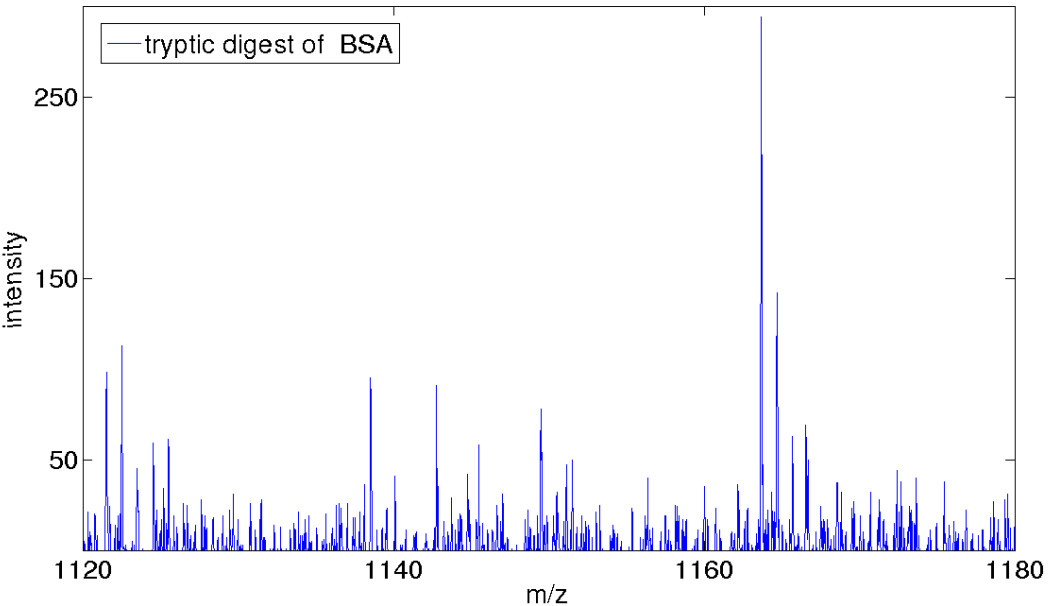


Sequence coverage	
OpenMS	Bruker
52% - 67%	44%

rel. err. [ppm]	
OpenMS	Bruker
80 - 93	95

Results: MALDI-TOF data

Data set II (tryptic digest of bovine serum)

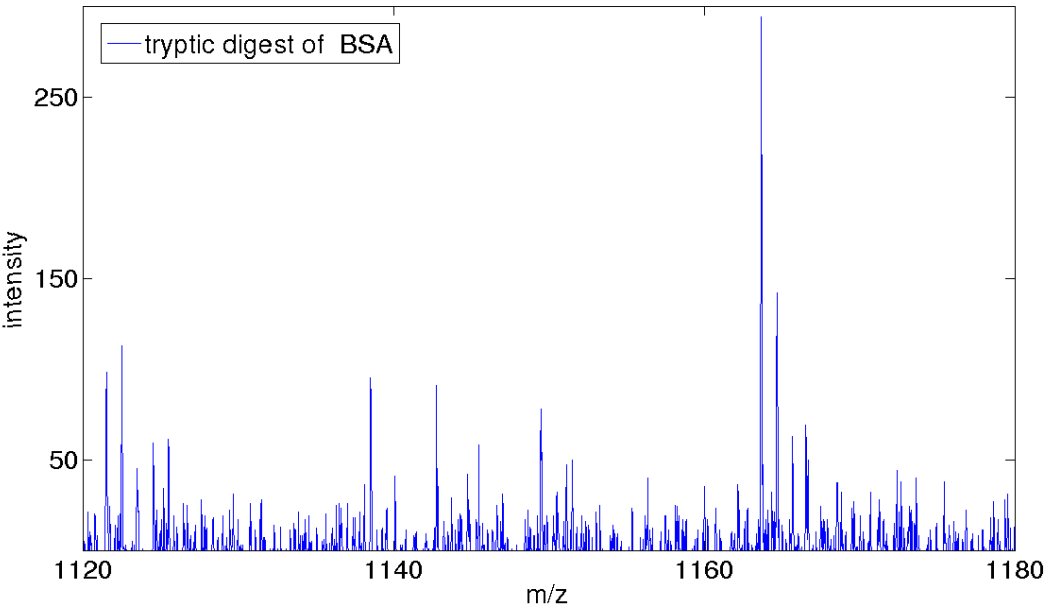


Sequence coverage	
OpenMS	Bruker
52% - 67%	44%

rel. err. [ppm]	
OpenMS	Bruker
80 - 93	95

Results: MALDI-TOF data

Data set II (tryptic digest of bovine serum)

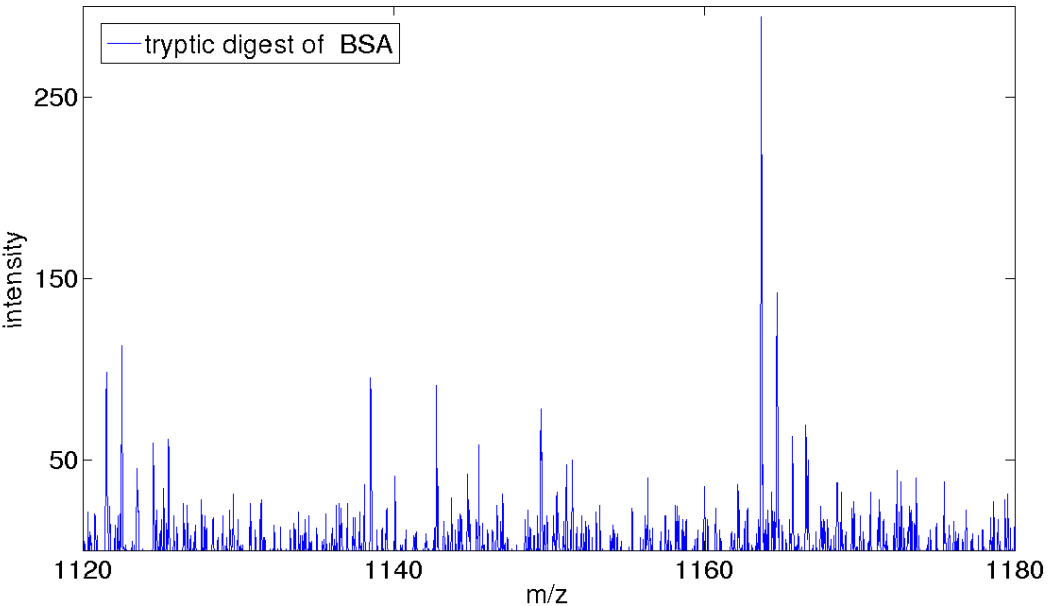


Sequence coverage	
OpenMS	Bruker
52% - 67%	44%

rel. err. [ppm]	
OpenMS	Bruker
80 - 93	95

Results: MALDI-TOF data

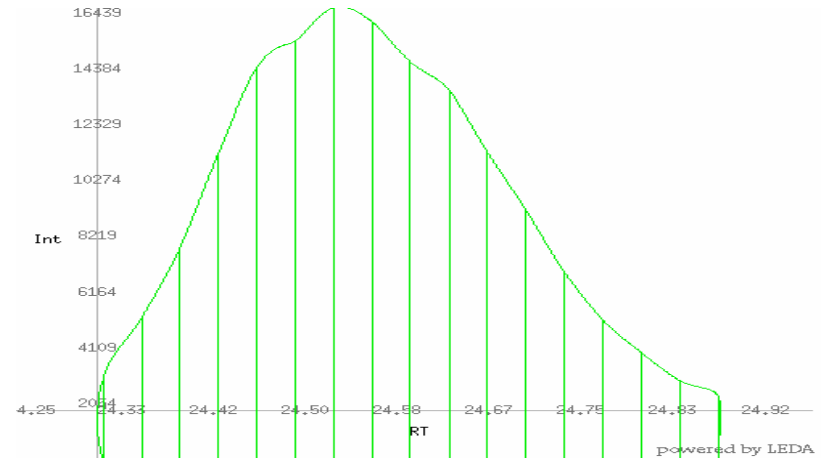
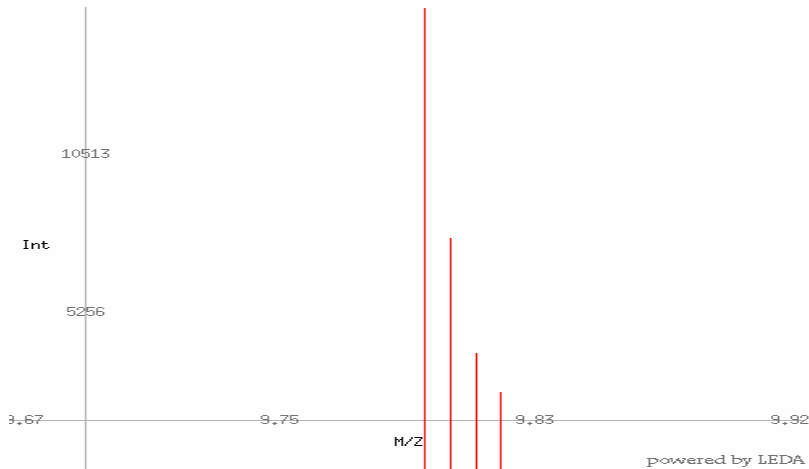
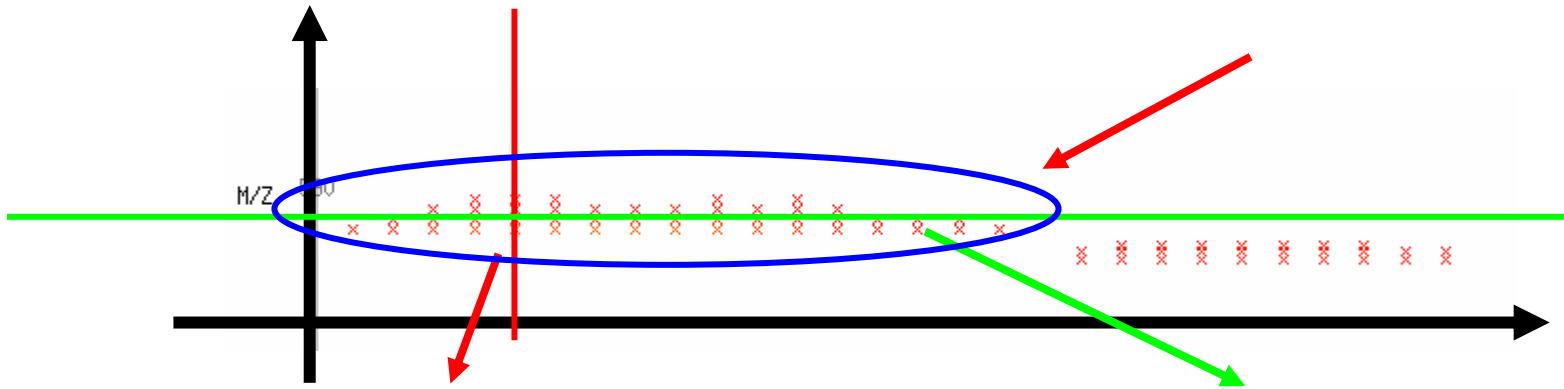
Data set II (tryptic digest of bovine serum)



Sequence coverage	
OpenMS	Bruker
52% - 67%	44%

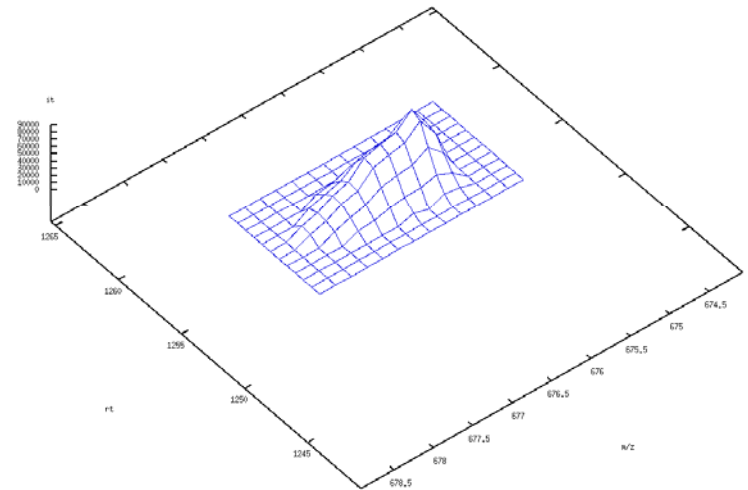
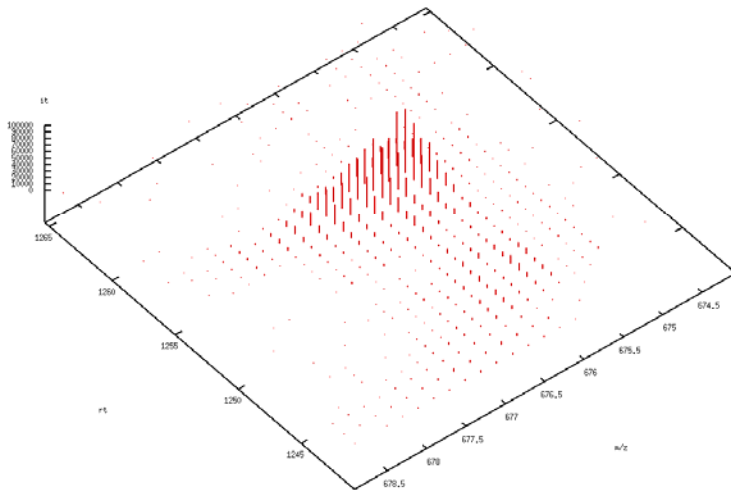
rel. err. [ppm]	
OpenMS	Bruker
80 - 93	95

- Motivation for OpenMS
- Bioinformatics Issues in quantitative Proteomics
 - Signal Processing
 - **Feature Finding**
 - Map Mapping
 - Differential Quantitation
 - Identification, Clustering
 - Software Engineering, Databases



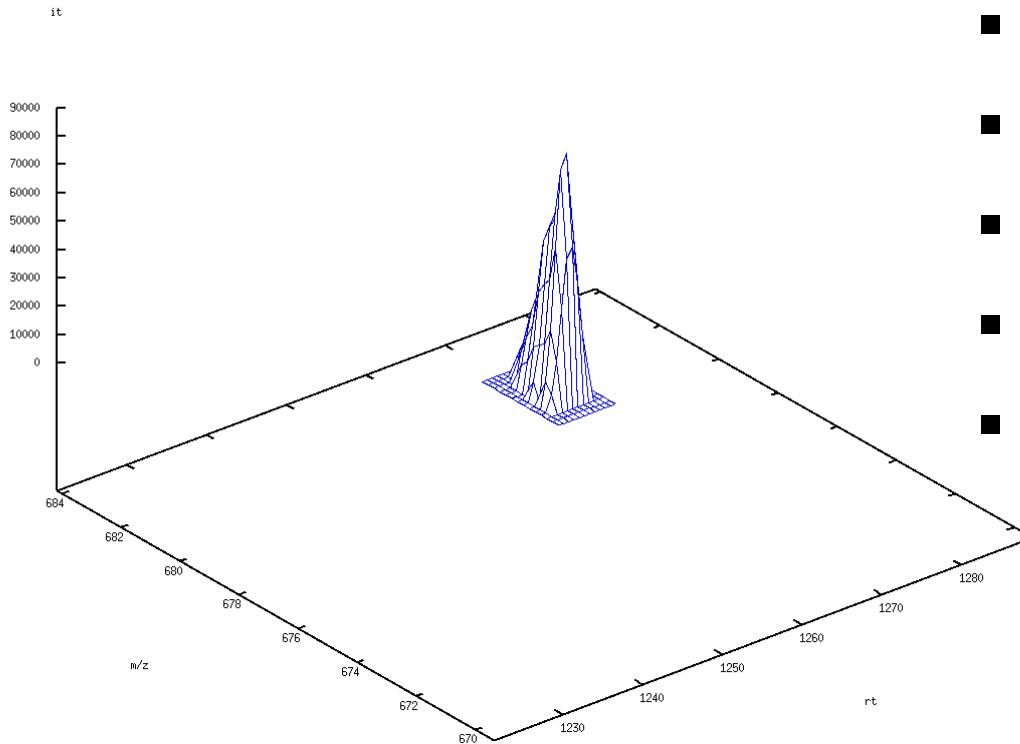
Capture ALL peaks belonging to a peptide for quantitation!

- A two-dimensional *model* has to be adjusted to the raw data

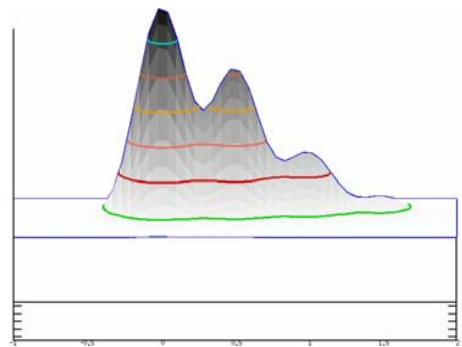
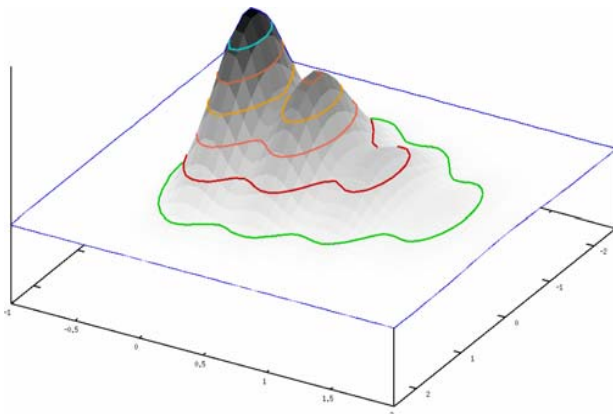


Attributes of a feature:

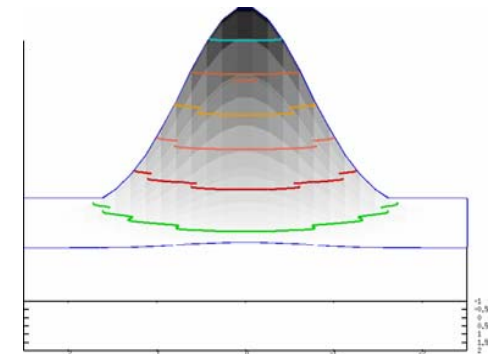
- mass-to-charge ratio
- retention time
- intensity, **volume**
- quality
-



feature model = isotope pattern + elution profile



m/z



rt

- Natural isotopes occur with well-known abundances
- Peak intensities determined by binomial convolution
- Depend on molecular formula of peptide

^{12}C 98.90%

^{13}C 1.10%

^{14}N 99.63%

^{15}N 0.37%

^{16}O 99.76%

^{17}O 0.04%

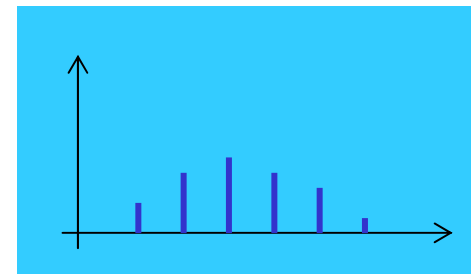
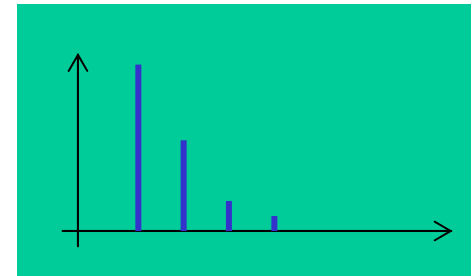
^{18}O 0.20%

^1H 99.98%

^2H 0.02%

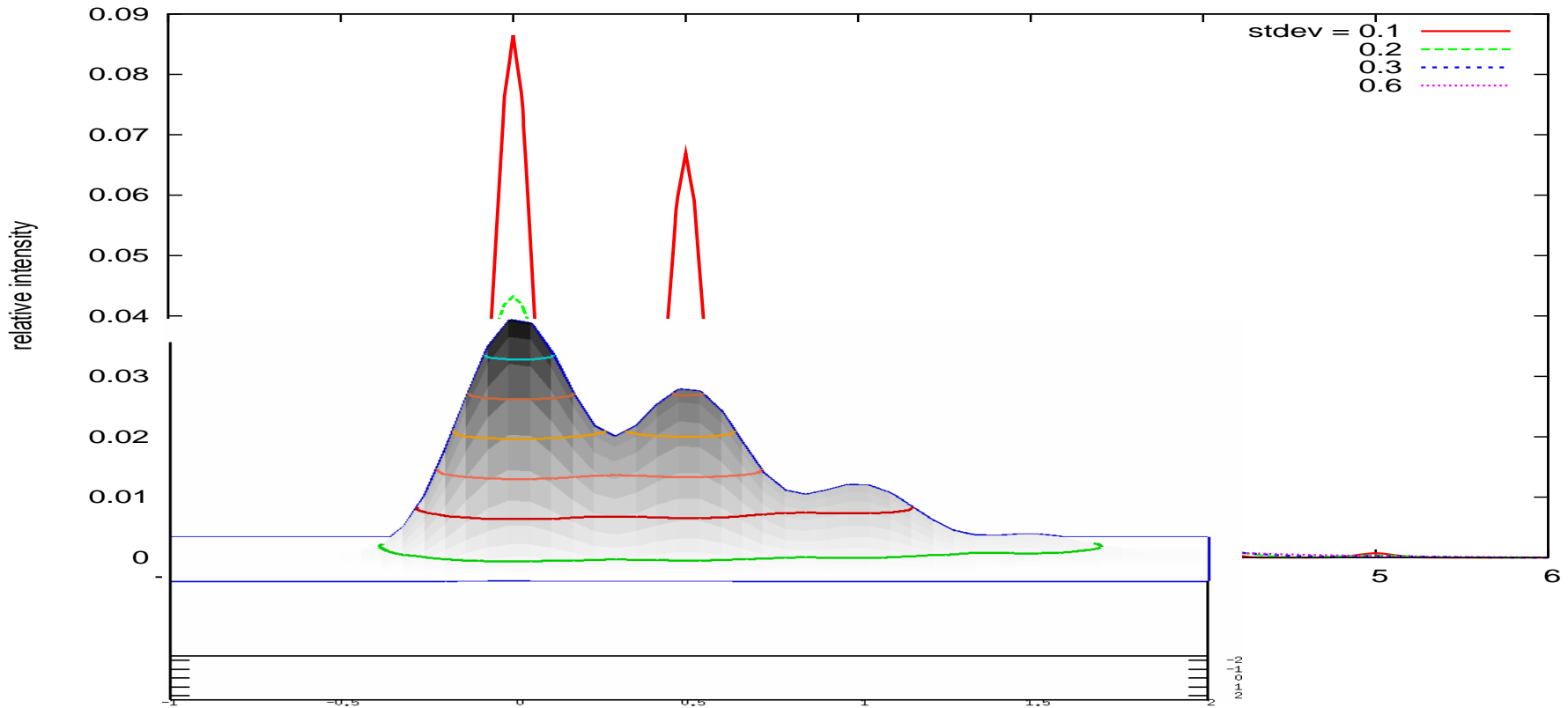
Randomly sample peptides from protein database and compute average composition.
Table with average isotope pattern for specified peptide mass

m [Da]	P (k=0)	P (k=1)	P (k=2)	P (k=3)	P (k=4)
1000	0.55	0.30	0.10	0.02	0.00
2000	0.30	0.33	0.21	0.09	0.03
3000	0.17	0.28	0.25	0.15	0.08
4000	0.09	0.20	0.24	0.19	0.12



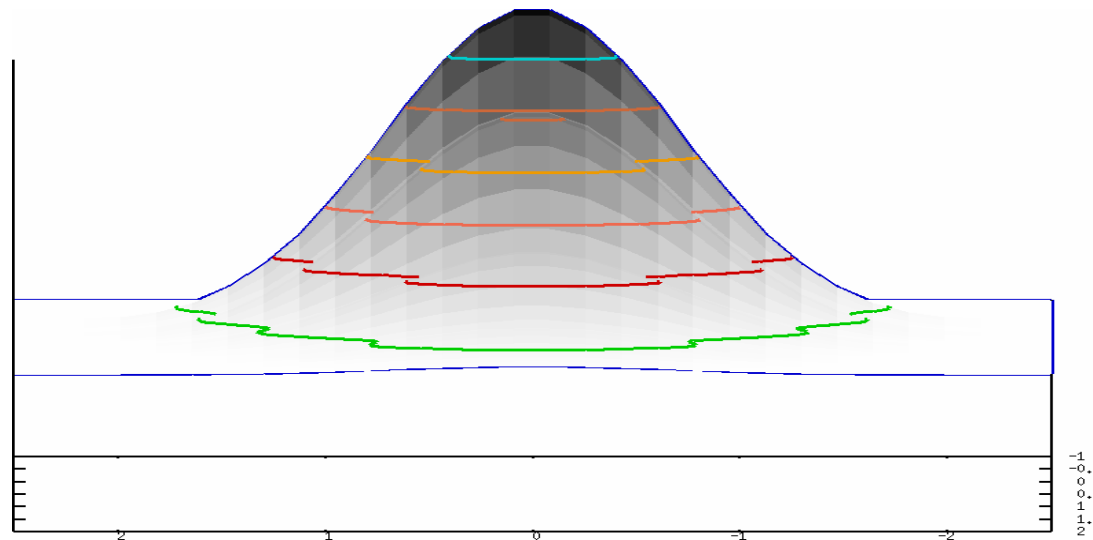
- Isotope patterns

Effect of the smoothing width on an average isotope pattern at mass 1350

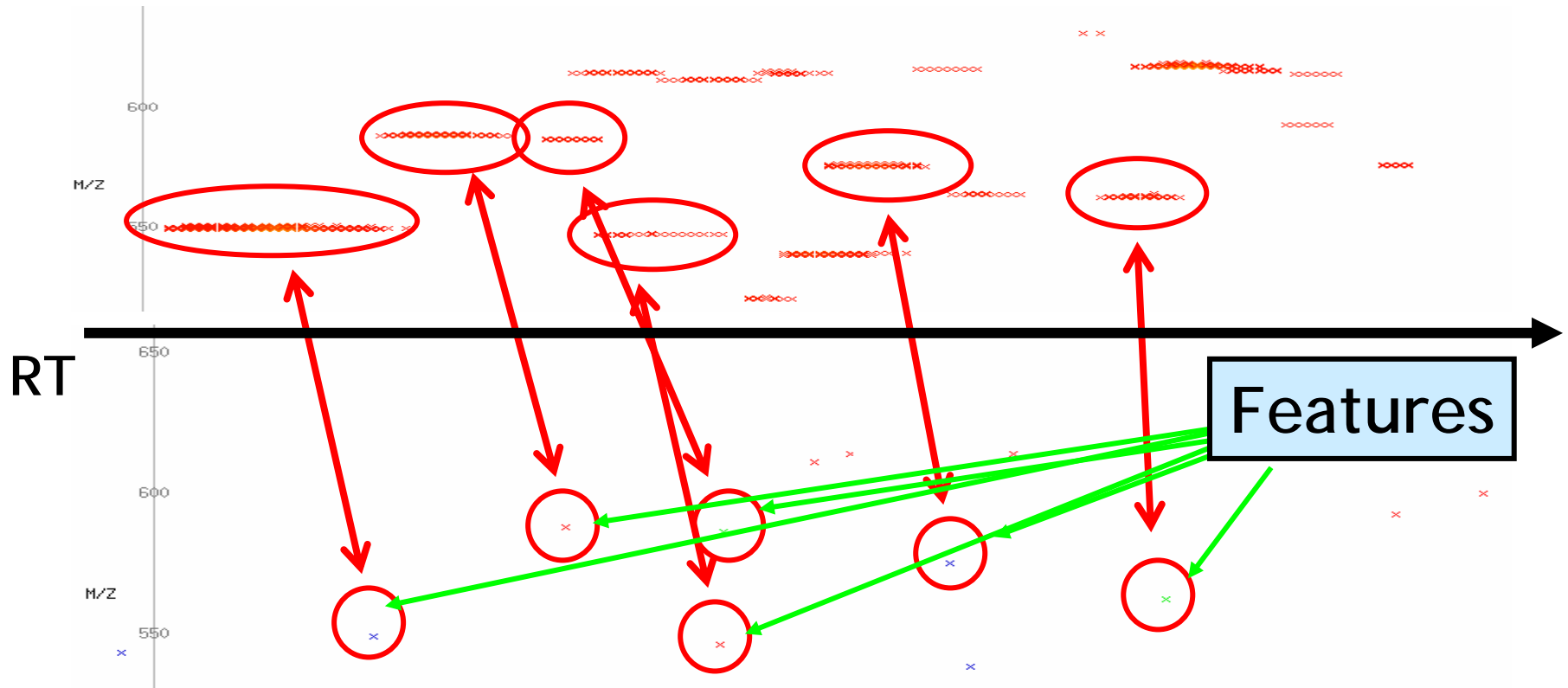




- Elution profiles
 - Currently modeled by a normal distribution, other shapes possible



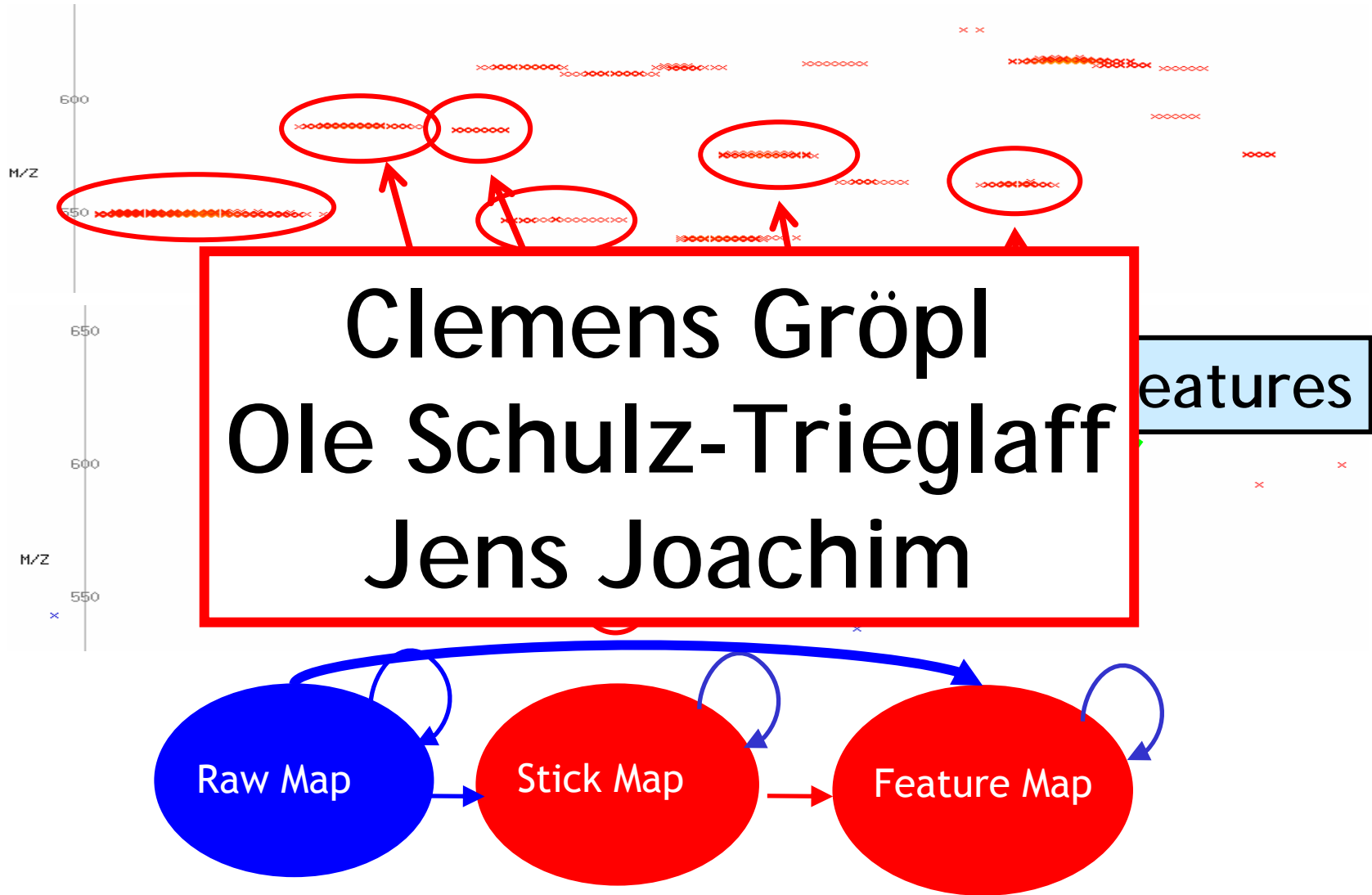
Peak Aggregation



Peak aggregation to features yields

- 1,000-fold compression in data volume
- reduction to meaningful entity: a peptide
- accurate quantitation: integration of all related peaks

Peak Aggregation



- Noise/baseline filtering
- Peak picking, integration
- Peak aggregation
 - De-isotoping: joining peaks from same isotope pattern
 - De-charging: joining features from same peptide, different charge states

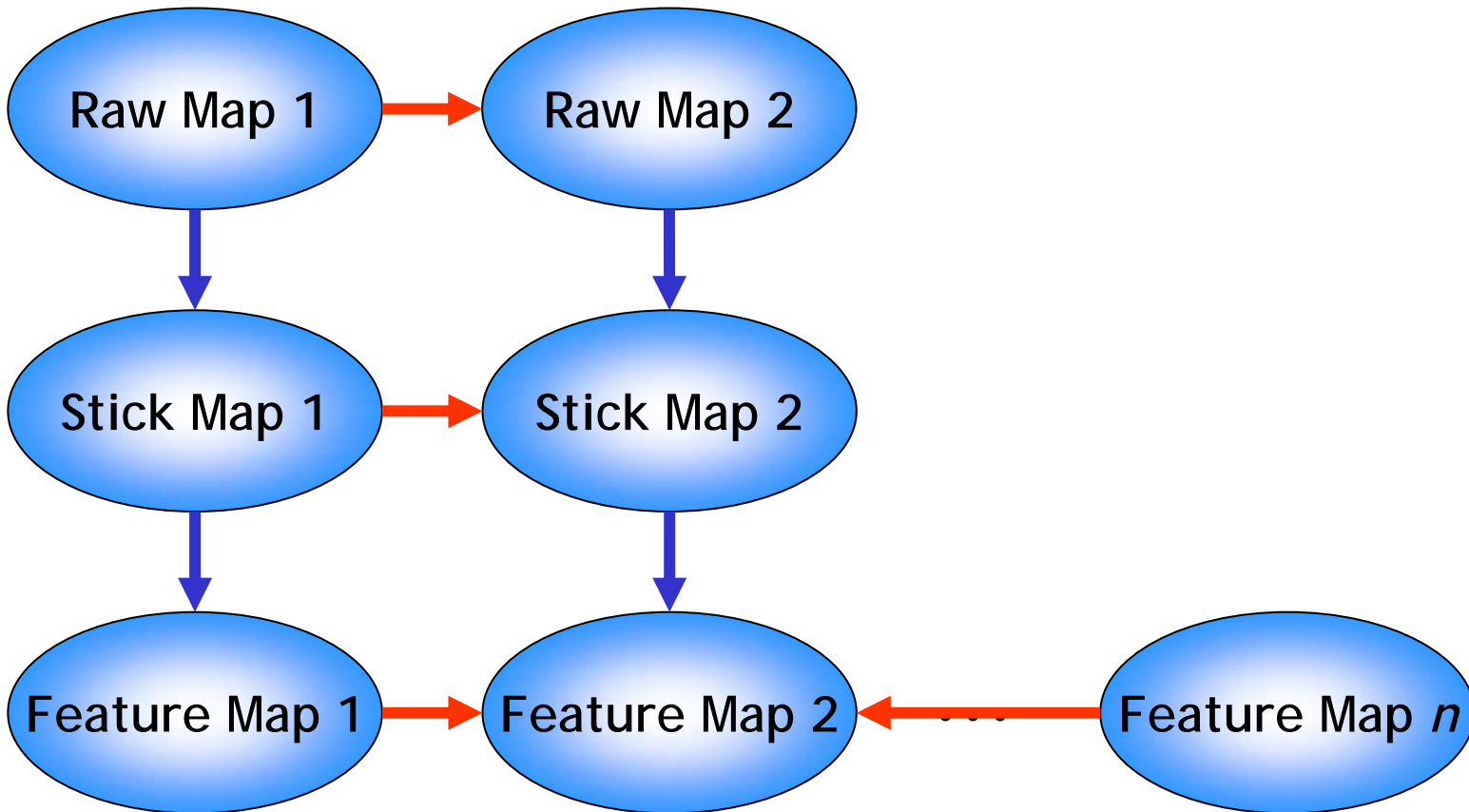
⇒ Aggregated Map

Reduction in Data Volume: up to 10^4

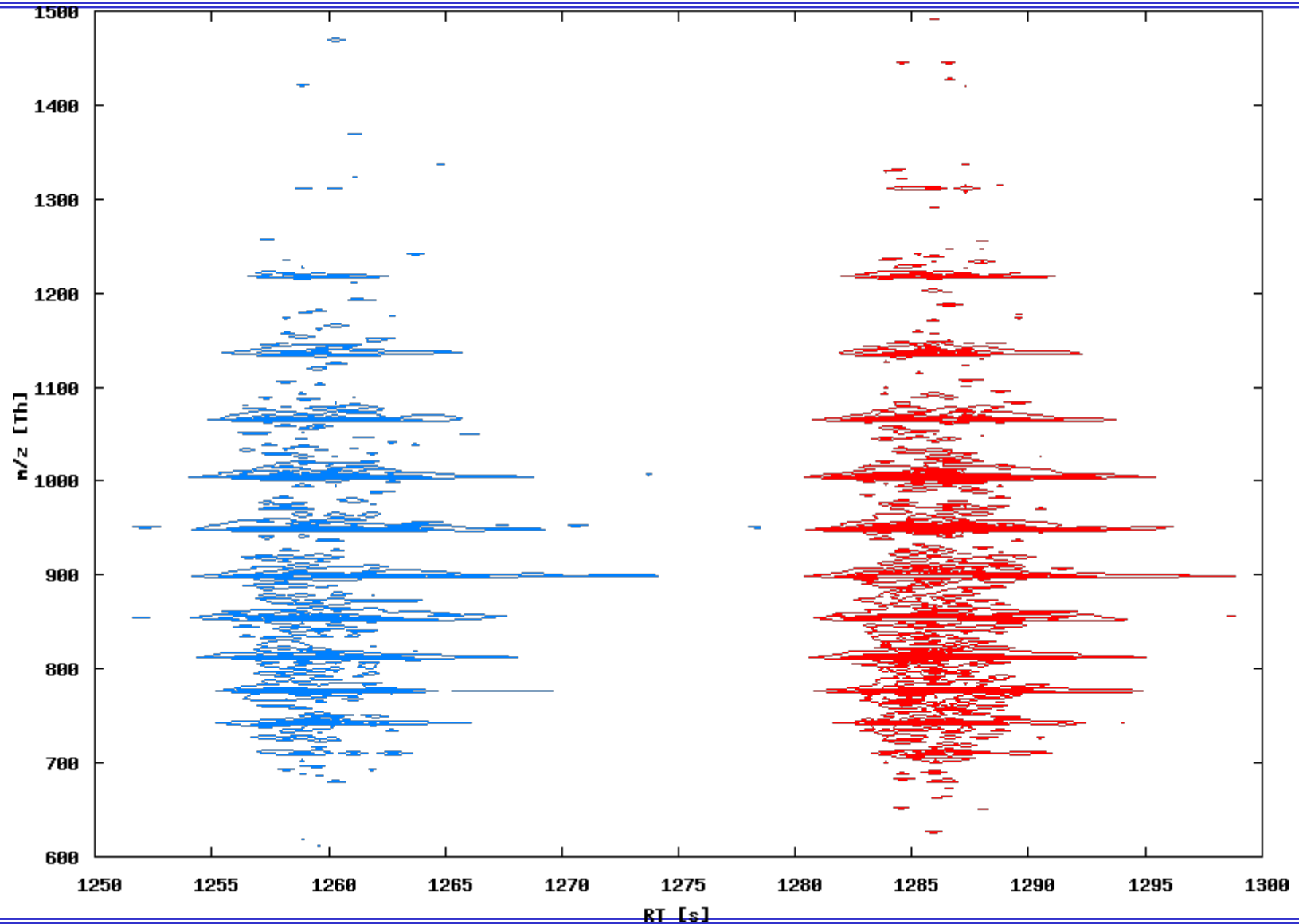
- Motivation for OpenMS
- Bioinformatics Issues in quantitative Proteomics
 - Signal Processing
 - Feature Finding
 - **Map Mapping**
 - Differential Quantitation
 - Identification, Clustering
 - Software Engineering, Databases
- Subsequent talks give more details

Map Mapping - At All Levels

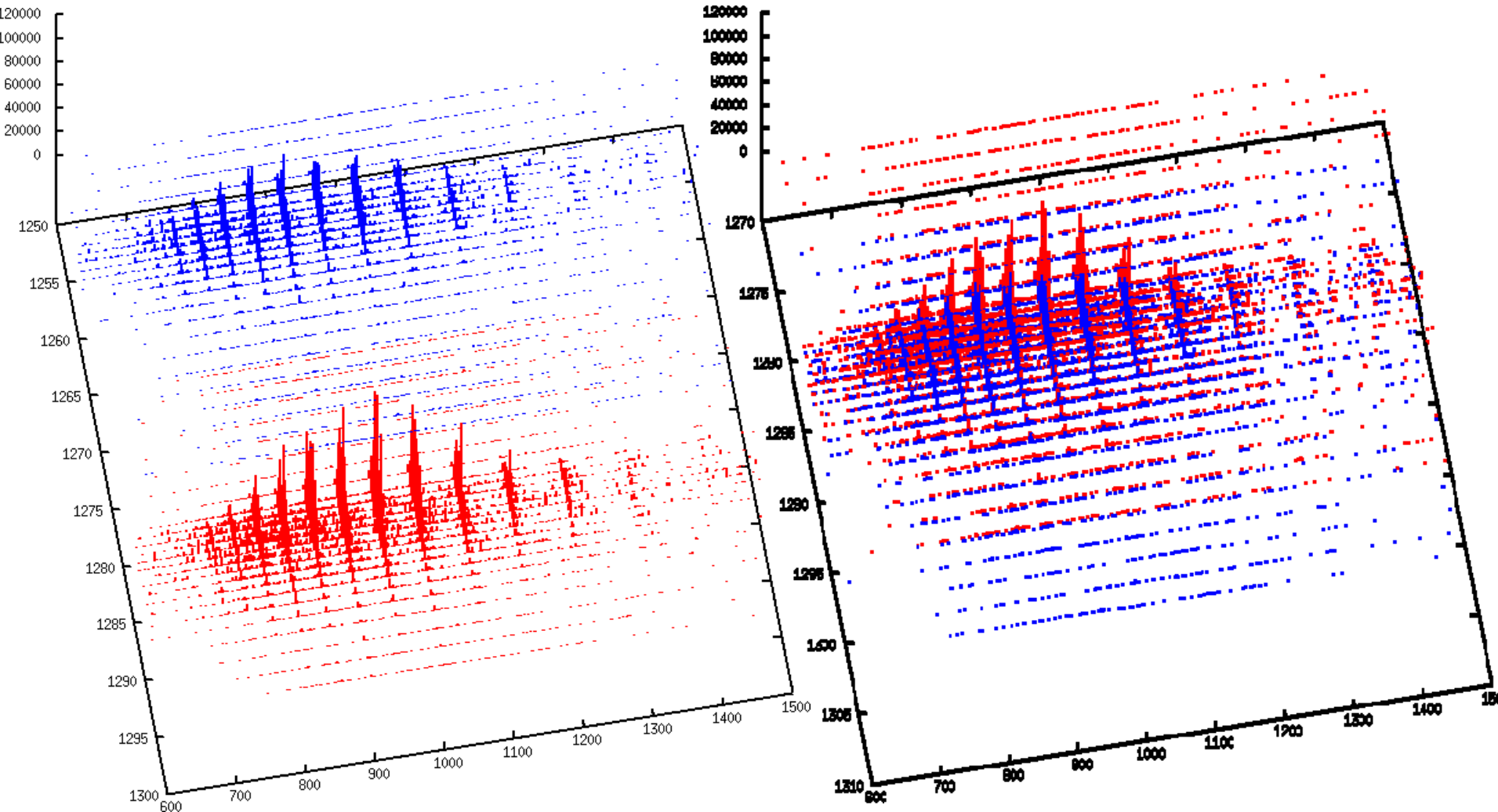
→ Data reduction → Mapping



Mapping of Raw Maps

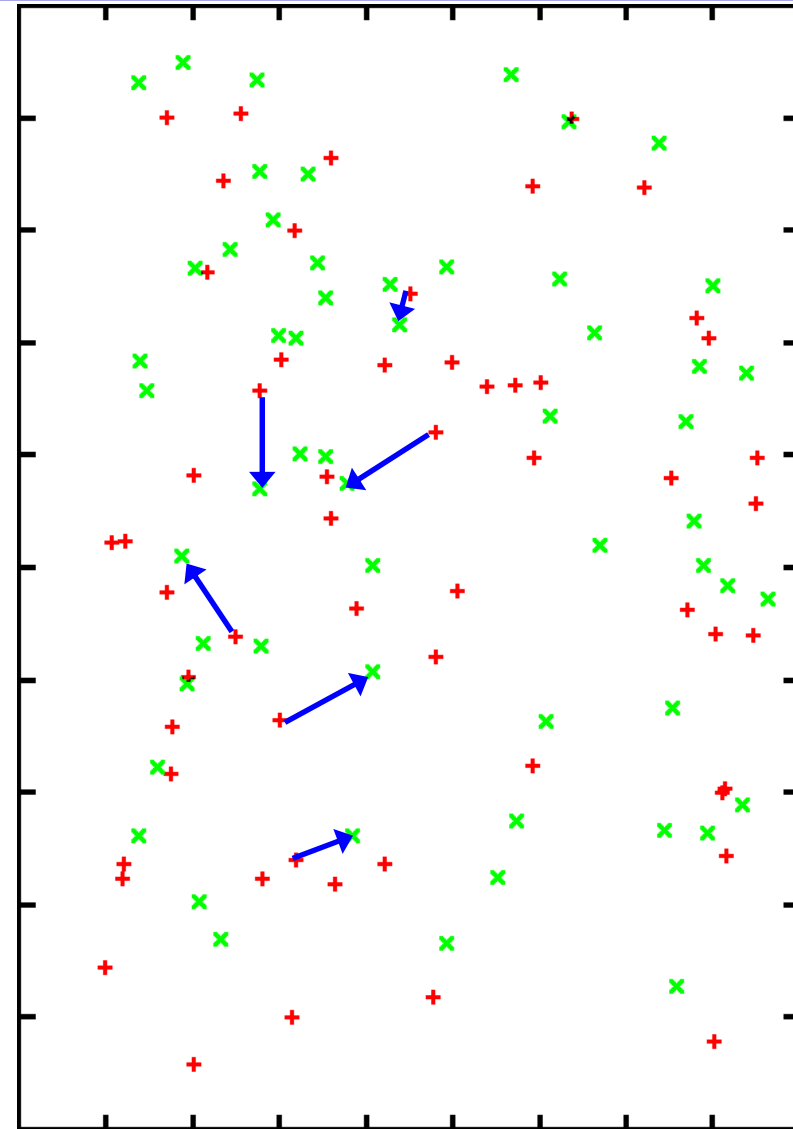


Mapping Discrete Peaks

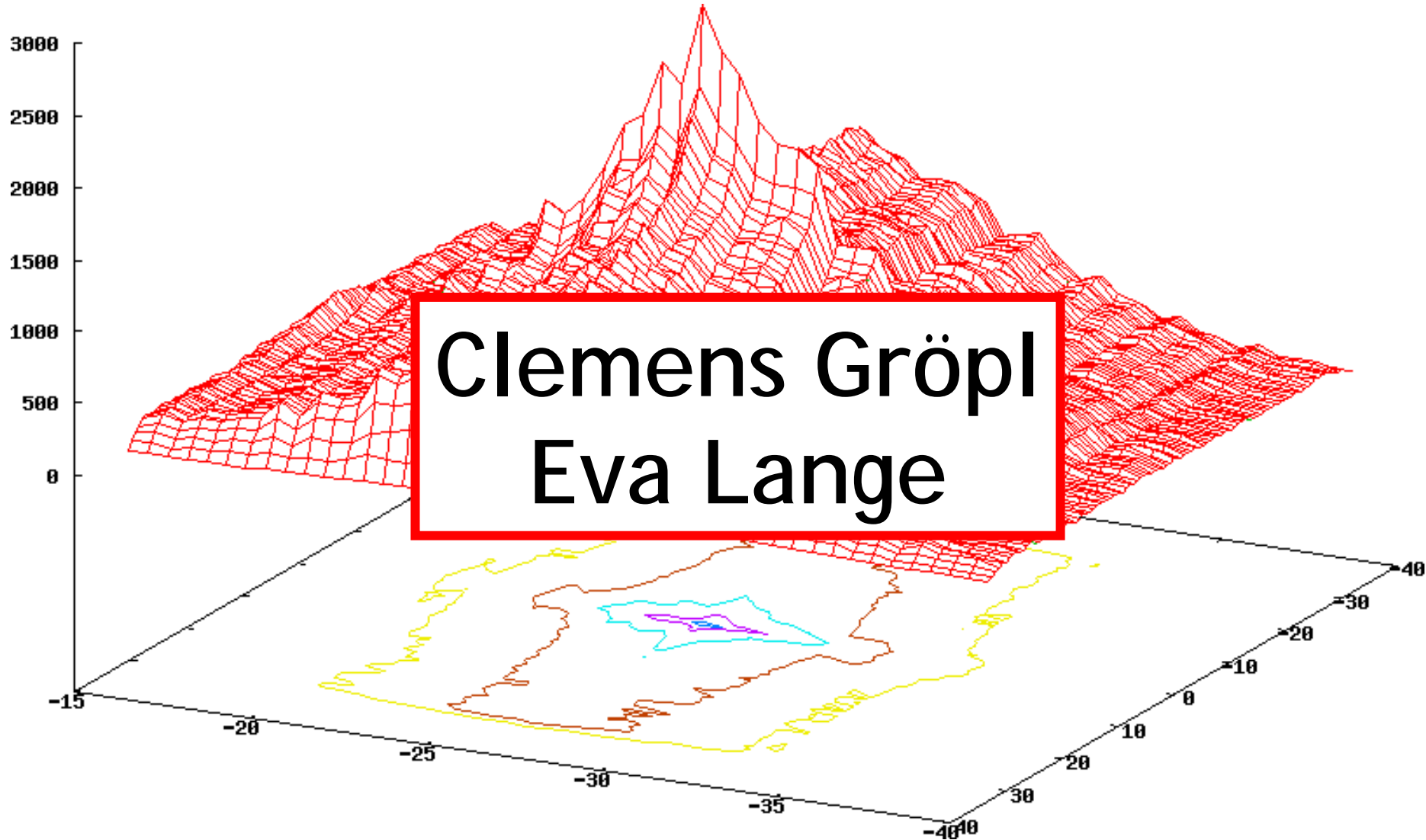


Geometric Hashing

- Consider difference vectors \rightarrow from features in map 1 (“+”) to features in map 2 (“x”)
- Restrict difference vector to reasonable differences along RT and m/z dimensions
- Hash these vectors into bins (“2D histogram”)
- In sufficiently similar maps, the optimal translation will be the most frequent one

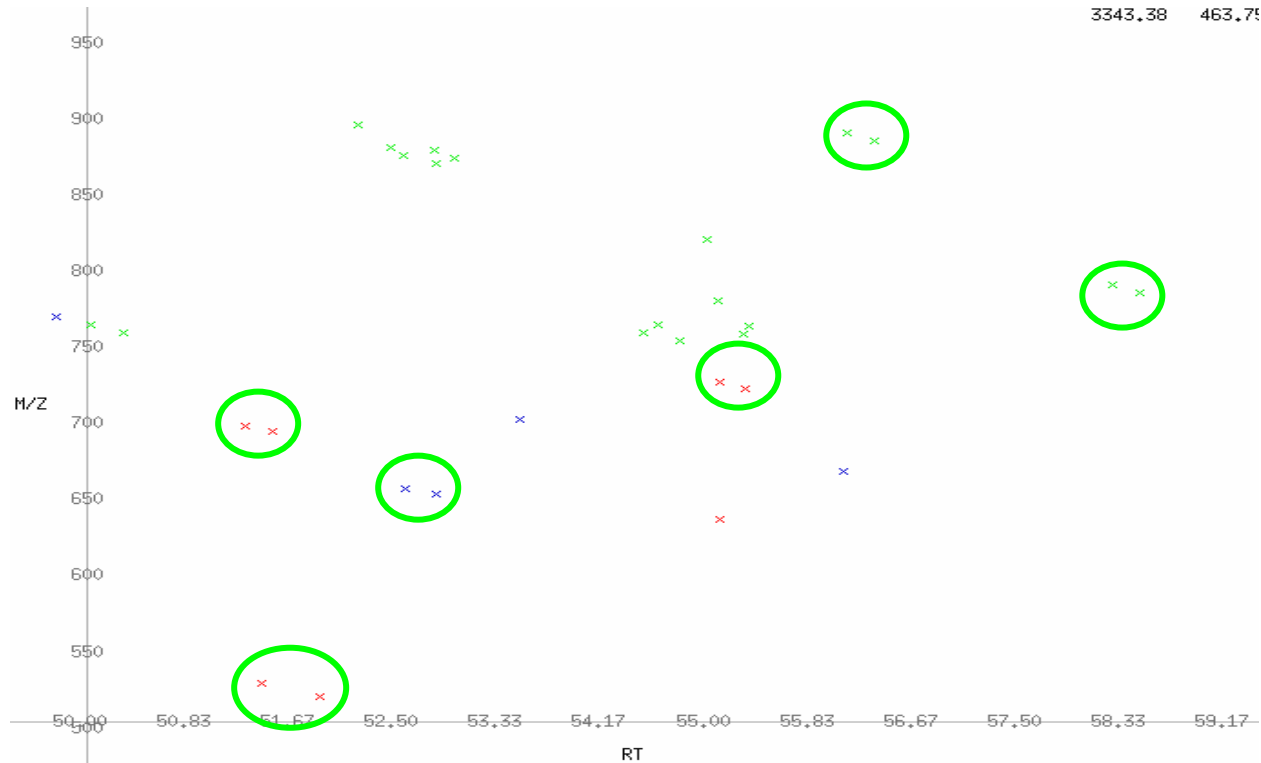


Finding the Offset Vector



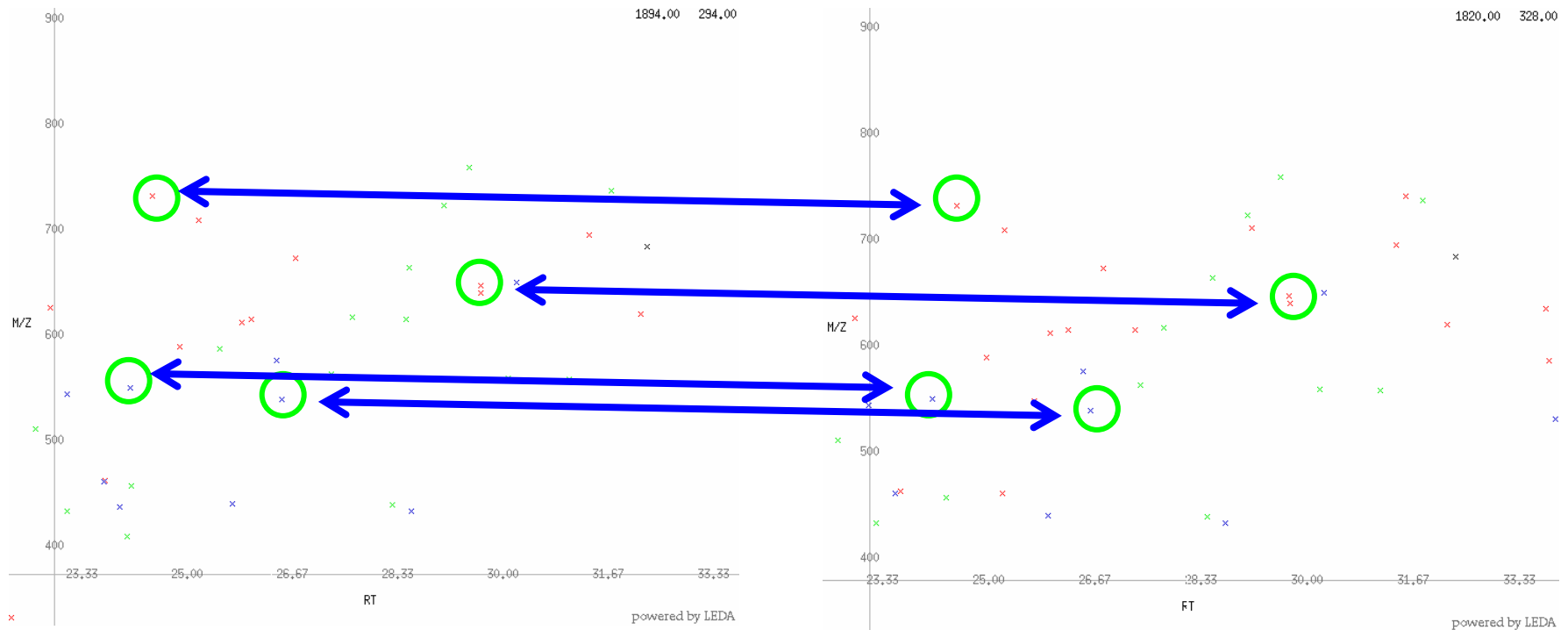
- Motivation for OpenMS
- Bioinformatics Issues in quantitative Proteomics
 - Signal Processing
 - Feature Finding
 - Map Mapping
 - **Differential Quantitation**
 - Identification, Clustering
 - Software Engineering, Databases
- Subsequent talks give more details

- Identify „differential peptides“
- Depends on quantitation technique
 - Labeling techniques:
Find pairs **within same map**
 - Unlabeled (DDQ):
Assign matching pairs **across maps**
- Usually done manually
- In HT settings this *has* to be done automatically!
- Reduction to features helps enormously



Assuming charge z and n Cys in a peptide,
check for pairs $(8n)/z$ Thomson apart in ONE MAP

DDQ - Peptide Assignment



Assign pairs across two or more maps
RT of peptide may vary between maps
⇒ compute suitable mapping
(linear function usually suffices)



Computational Geometry Algorithms Library

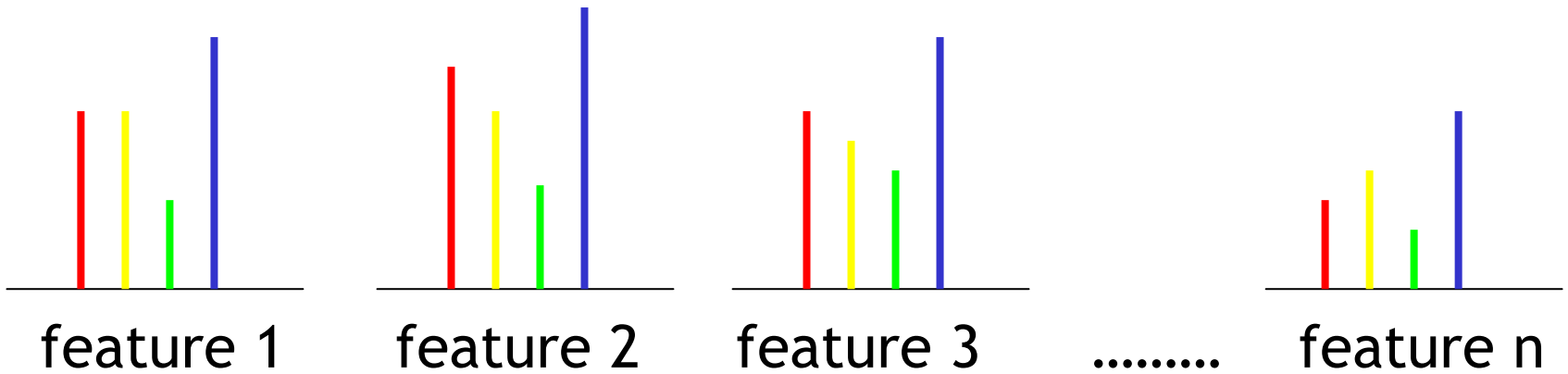
- Use efficient data structures for range searches (e.g. Kd-trees or range trees)
- Employ memory management for cached range queries to allow for large data sets (e.g. large stick maps)
- Making use of robust and efficient implementations of these data structures in CGAL (Computational Geometry Algorithms Library)
- CGAL provides standard data structures for d -dimensional computational geometry

- Intensities will vary between maps, even for the same sample (e.g. volume error during injection)
- Usually this results in a constant ratio of feature intensities
- Map normalization is required to determine this ratio for accurate relative quantitation
- Assuming a sufficient similarity between two maps, the majority of features will be identical, only few will differ in concentration between the two samples

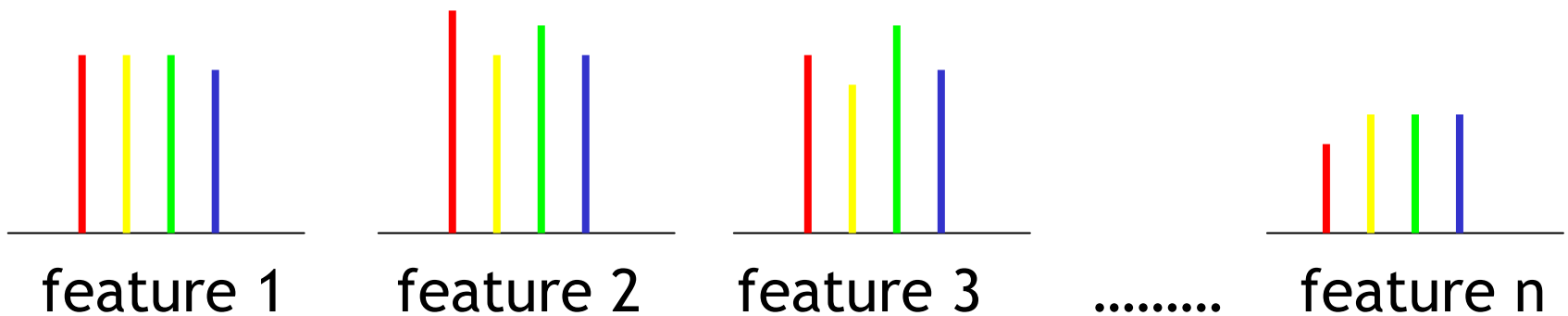
Run several experiments on the same sample

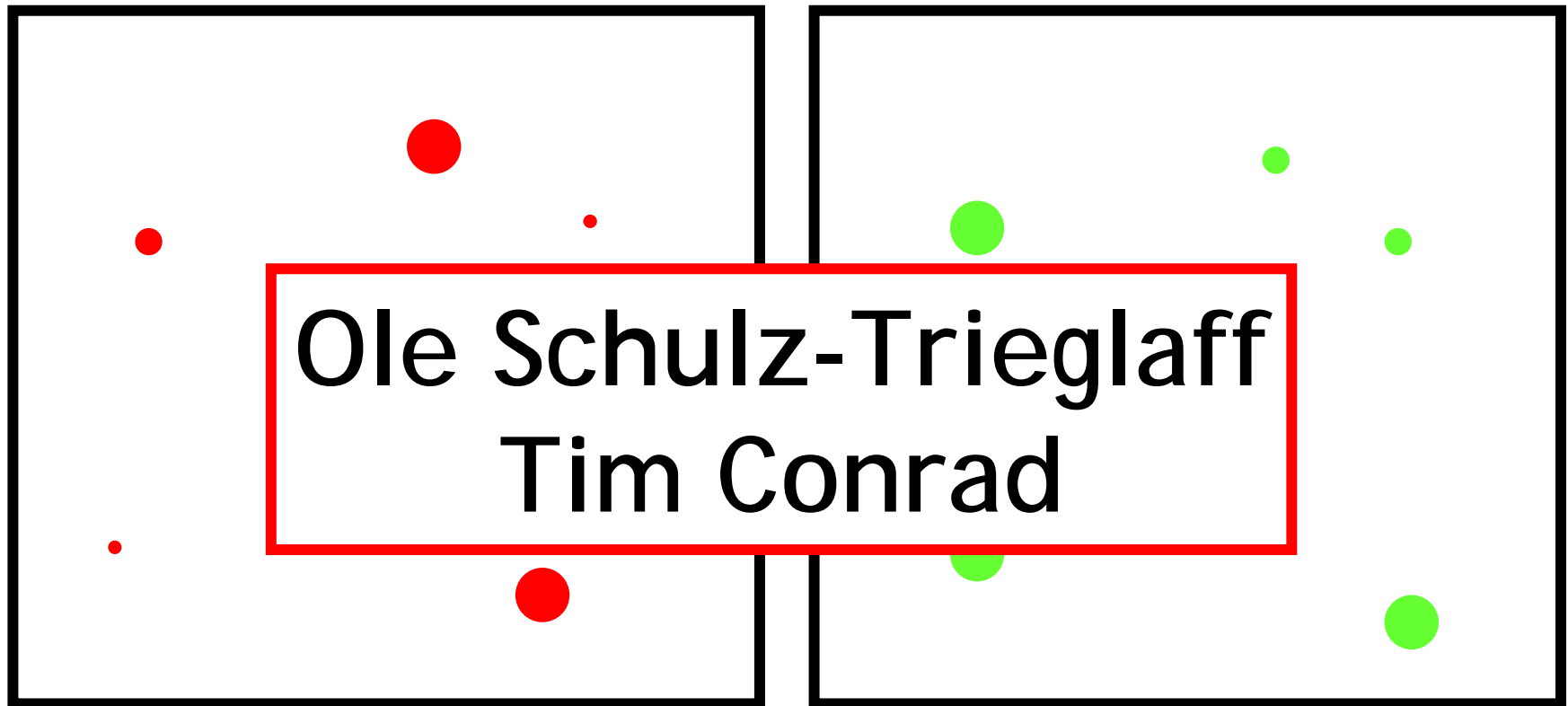
- Too time-consuming for manual assignment
- Automatic assignment increases reliability
 1. Features that show up most of the time are real
 2. Features that do not are probably noise
 3. Allows statements about statistical significance

Normalize and group measurements



Compute scaling between pairs of maps, group them



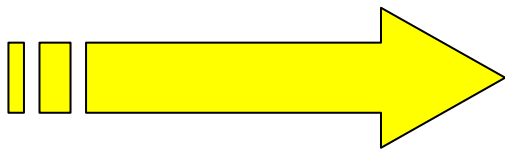
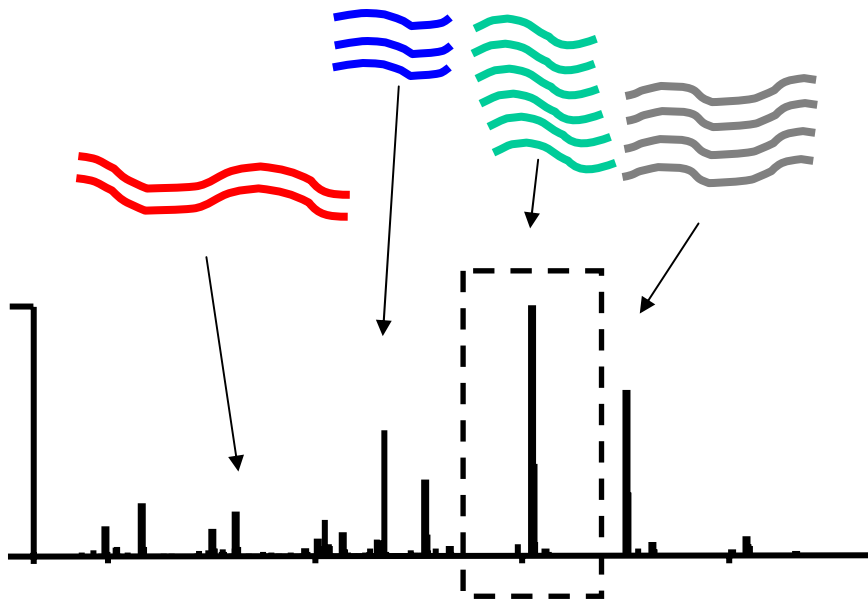


Marker in patient

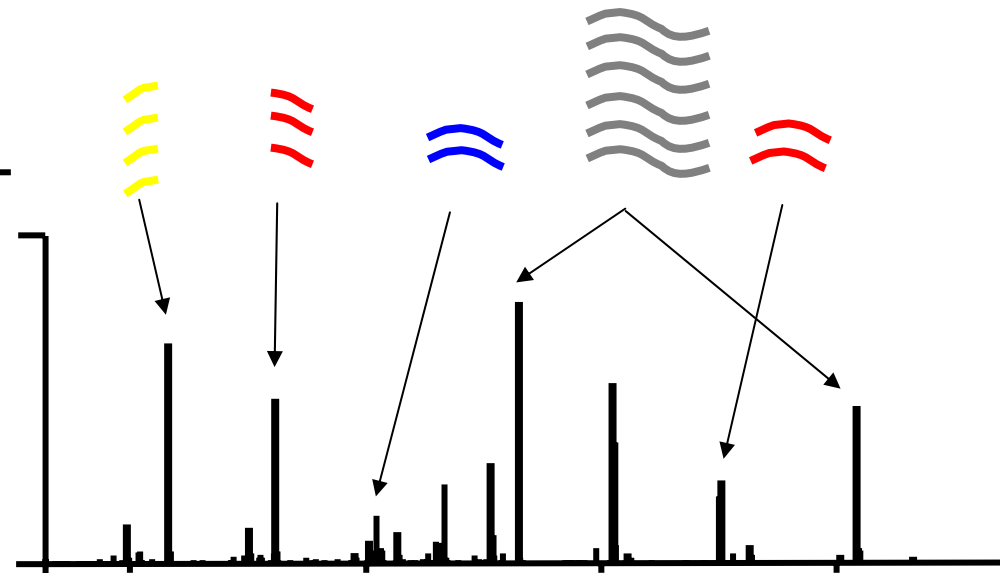
Marker in control

- Motivation for OpenMS
- Bioinformatics Issues in quantitative Proteomics
 - Signal Processing
 - Feature Finding
 - Map Mapping
 - Differential Quantitation
 - **Identification, Clustering**
 - Software Engineering, Databases

Peptide Identification (MS-MS)

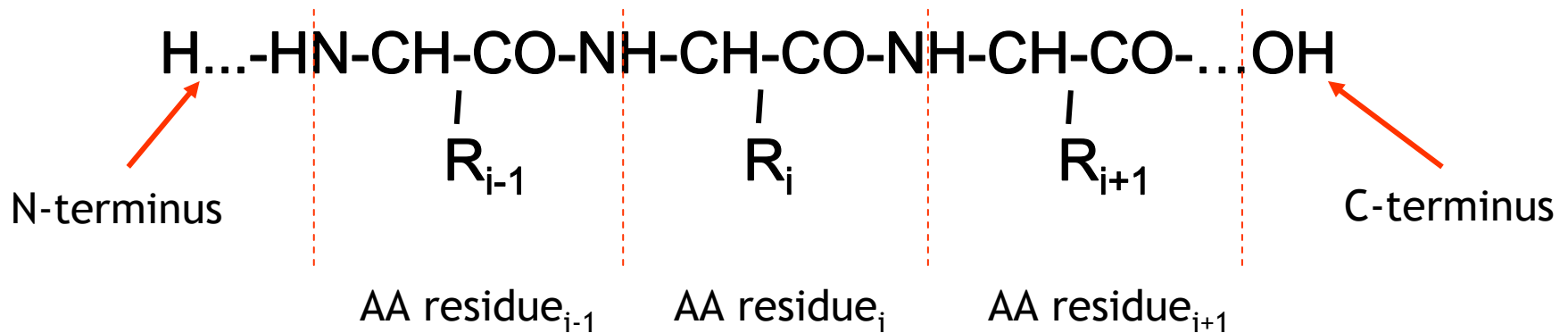


Secondary fragmentation

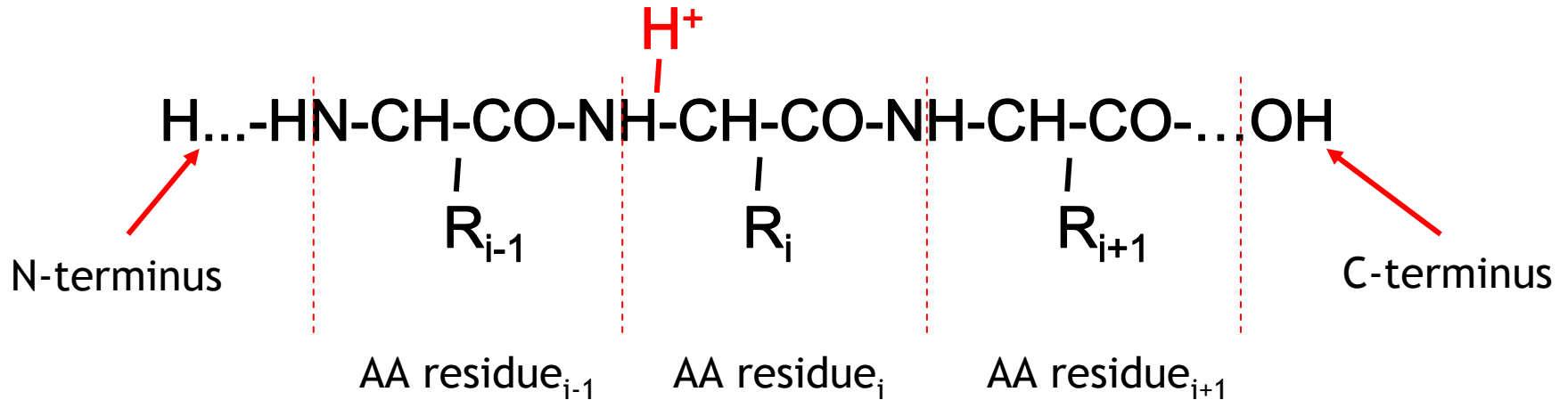


Fragment Spectrum

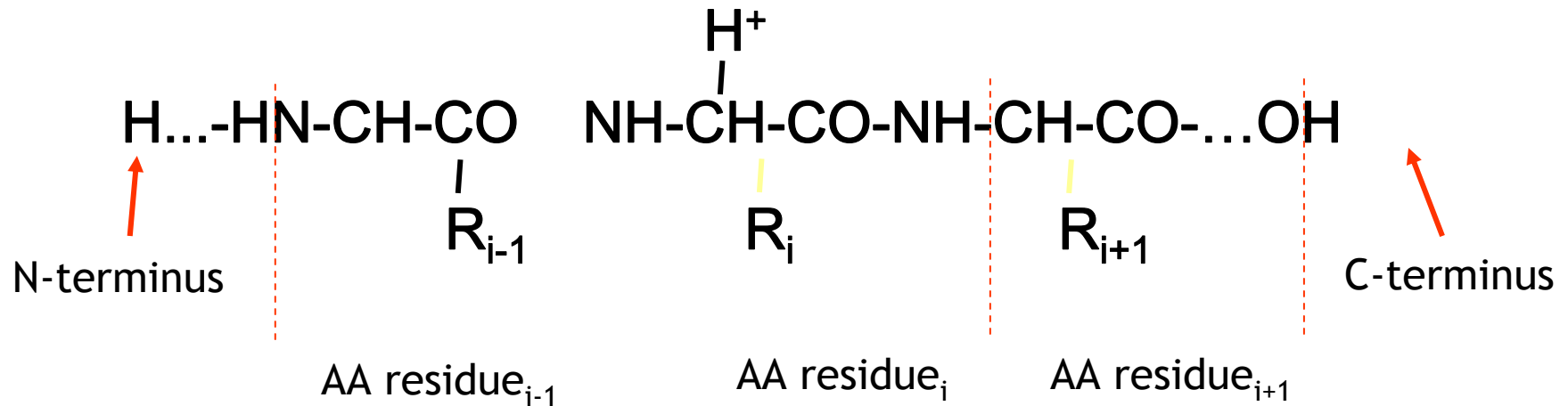
The peptide backbone breaks to form fragments with characteristic masses.



The peptide backbone breaks to form fragments with characteristic masses.



The peptide backbone breaks to form fragments with characteristic masses.

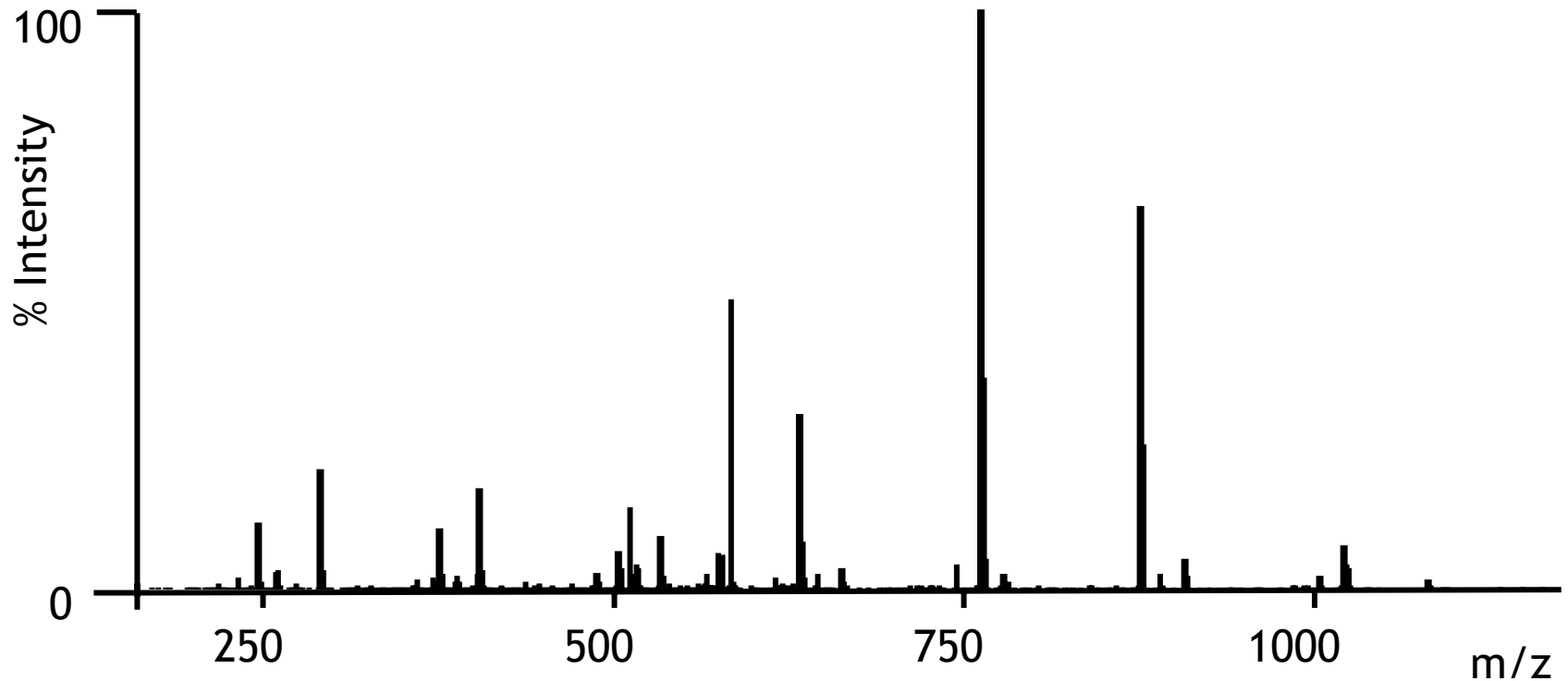


Ionized peptide fragment

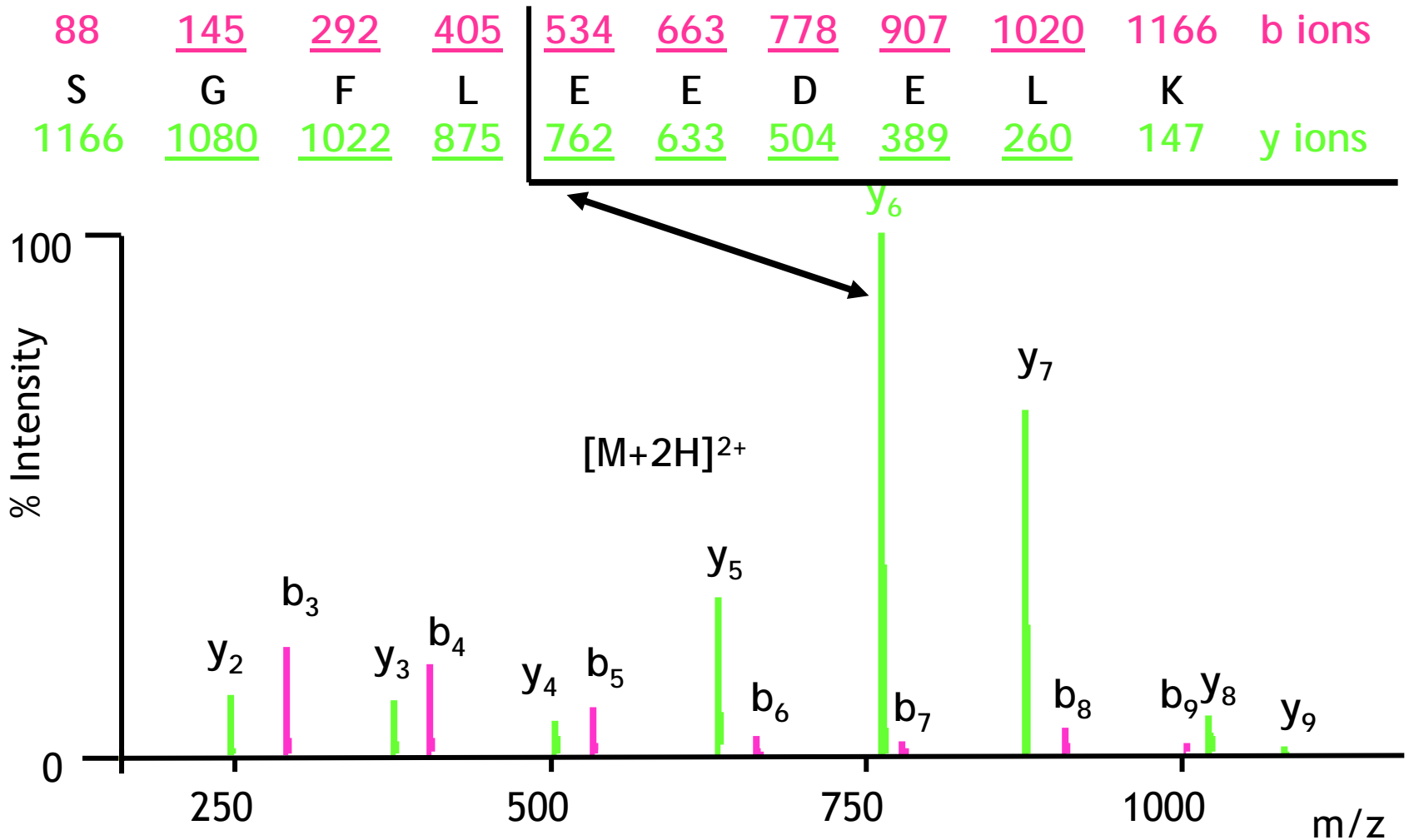
Peptide Ions in Spectrum



<u>88</u>	<u>145</u>	<u>292</u>	<u>405</u>	<u>534</u>	<u>663</u>	<u>778</u>	<u>907</u>	<u>1020</u>	1166	b ions
S	G	F	L	E	E	D	E	L	K	
1166	<u>1080</u>	<u>1022</u>	<u>875</u>	<u>762</u>	<u>633</u>	<u>504</u>	<u>389</u>	<u>260</u>	147	y ions

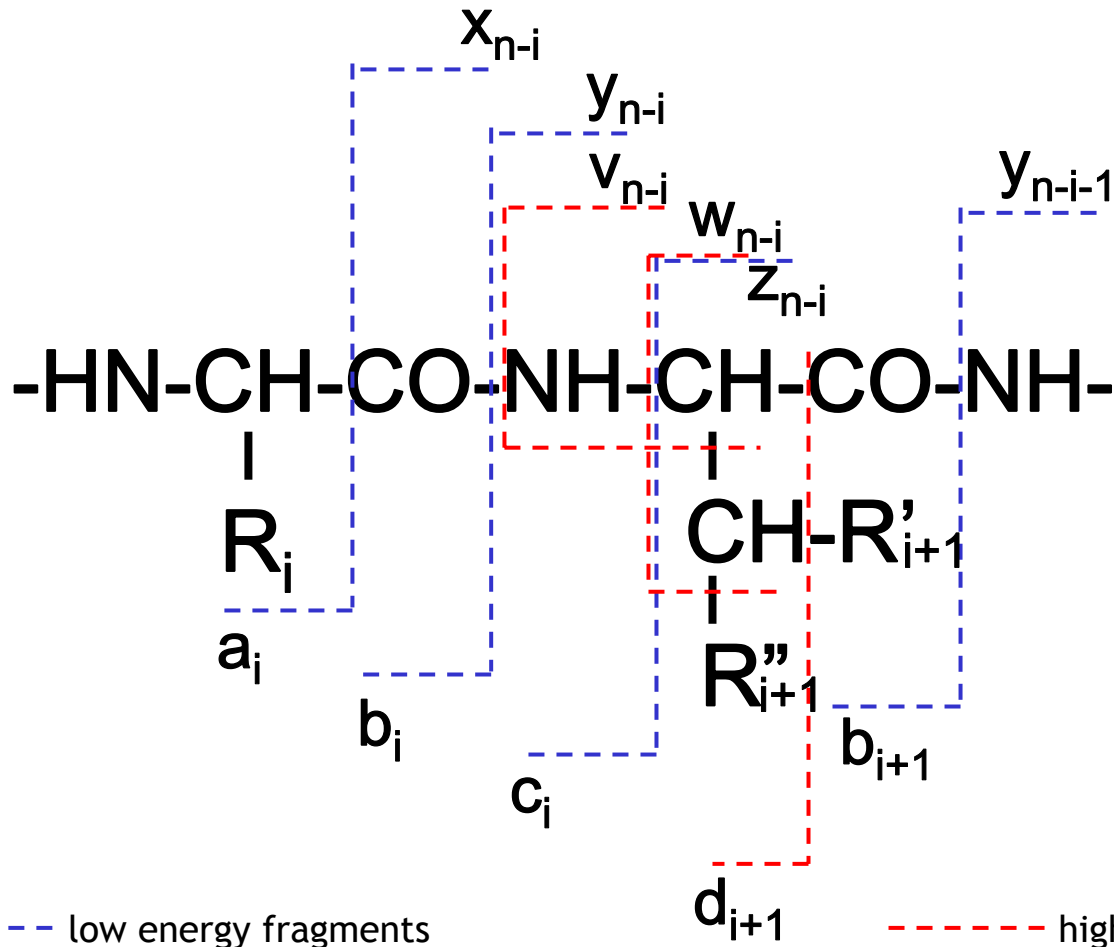


Peptide Ions in Spectrum



What's the problem?

Peptide fragmentation possibilities (ion types)



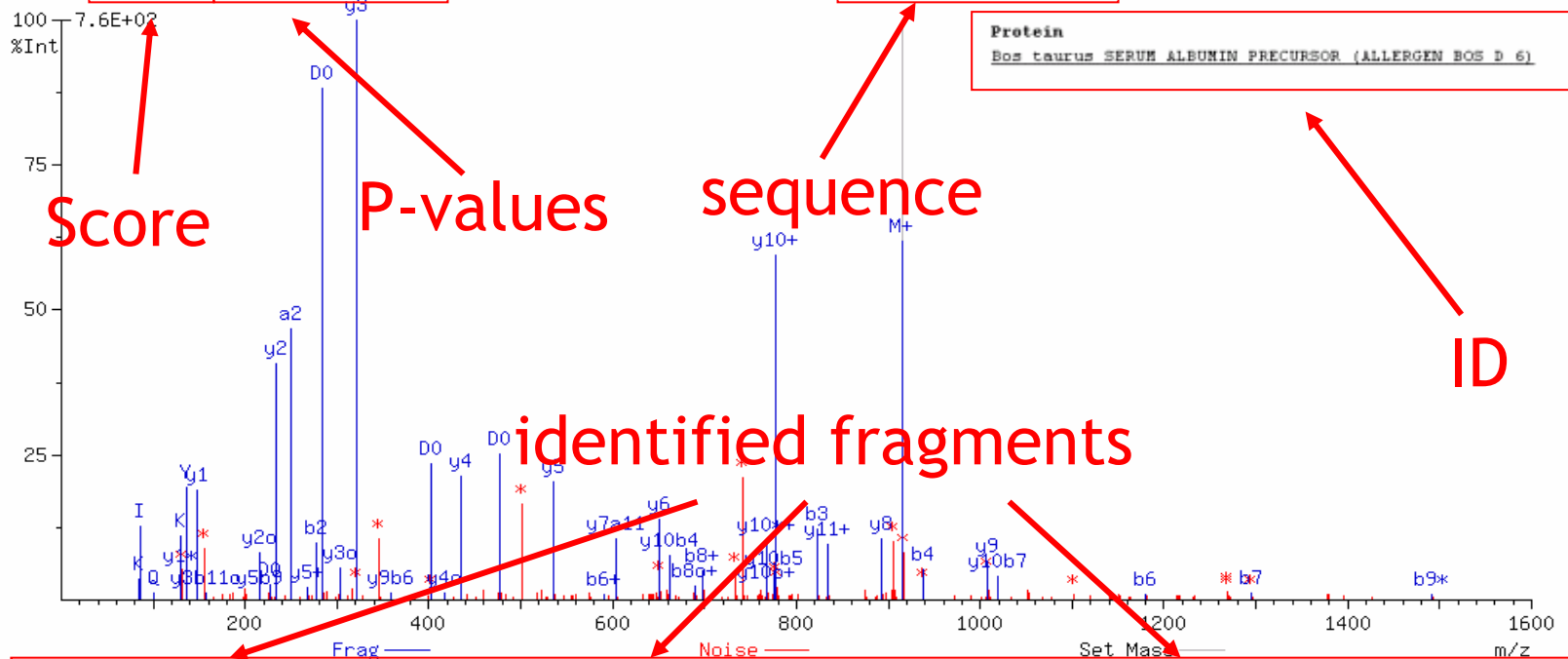
----- low energy fragments

----- high energy fragments

Identification - SCOPE



Rank	Score	DBpV	TSpV	Ions	%TIC	b/y Score	MH+	Sequence	#Inst	Prot ID	Prot Acc
1	192.9	0.015	0.020	44/63	68	58	1828.8	K.YICDNQDTISSK.L	1	ALBU_BOVIN	SP:P02769



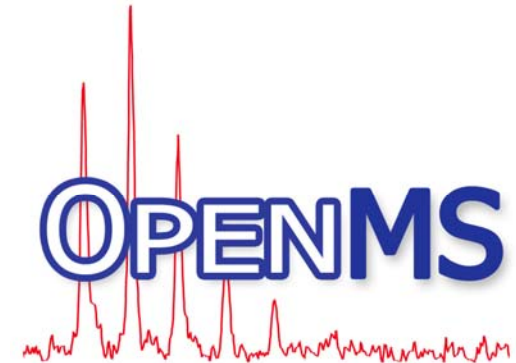
N-Terminal			C-Terminal			Parent Ions			
249.2	YI	a2	+0.02	833.4	ICDNQDTISSK y11+	-0.02	914.9	YICDNQDTISSK M+	-0.01
277.2	YI	b2	+0.02	767.8	CDNQDTISSK y10o+	-0.02			
822.4	YIC	b3	-0.01	768.3	CDNQDTISSK y10*+	-0.01			
937.4	YICD	b4	+0.00	776.8	CDNQDTISSK y10+	-0.02			
590.3	YICDNQ	b6+	-0.01	1007.5	DNQDTISSK y9	+0.03			
1179.5	YICDNQ	b6	-0.01	892.4	NQDTISSK y8	-0.02			
1294.5	YICDNQD	b7	+0.15	650.3	DTISSK y6	-0.00			
689.3	YICDNQDT	b8o+	-0.02	268.2	TISSK y5+	+0.01			
698.3	YICDNQDT	b8+	-0.01	535.3	TISSK y5	+0.03			
1491.6	YICDNQDTI	b9*	-0.05	416.3	ISSK y4o	-0.01			
				434.3	ISSK y4	-0.00			
				383.2	SSK y3o	+0.01			
				284.1	C	0.00			

Vineet Bafna, Nathan Edwards, Proc. ISMB 2001



- Motivation for OpenMS
- Bioinformatics Issues in quantitative Proteomics
 - Signal Processing
 - Feature Finding
 - Map Mapping
 - Differential Quantitation
 - Identification, Clustering
 - **Software Engineering, Databases**

- Software **framework** for shotgun proteomics
- **ISO/ANSI C++ compliant**
- **Features**
 - **Open Source (LGPL)**
 - **Data structures & algorithms**
 - *d*-dimensional core
 - Signal processing for raw data conversion
 - Data **import/export** in standard formats
 - **Database** storage of all data
(using an extended PEDRO/MIAPE model)
 - DB-driven **workflow** management
 - **Visualization** for spectra, maps

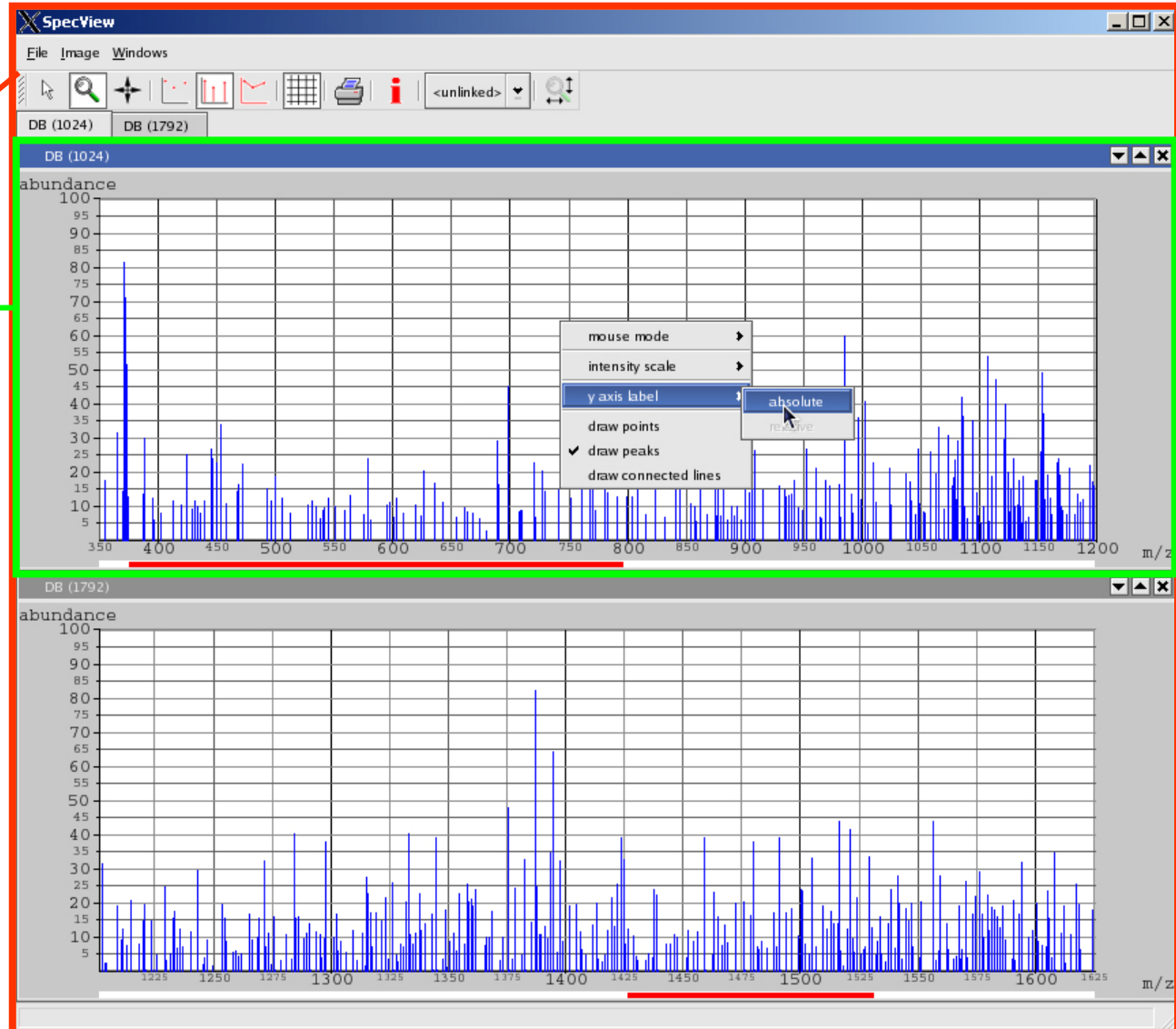
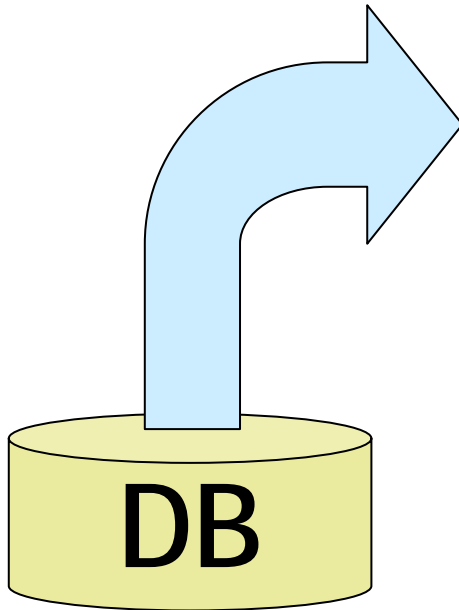


SpecView - 1D Window

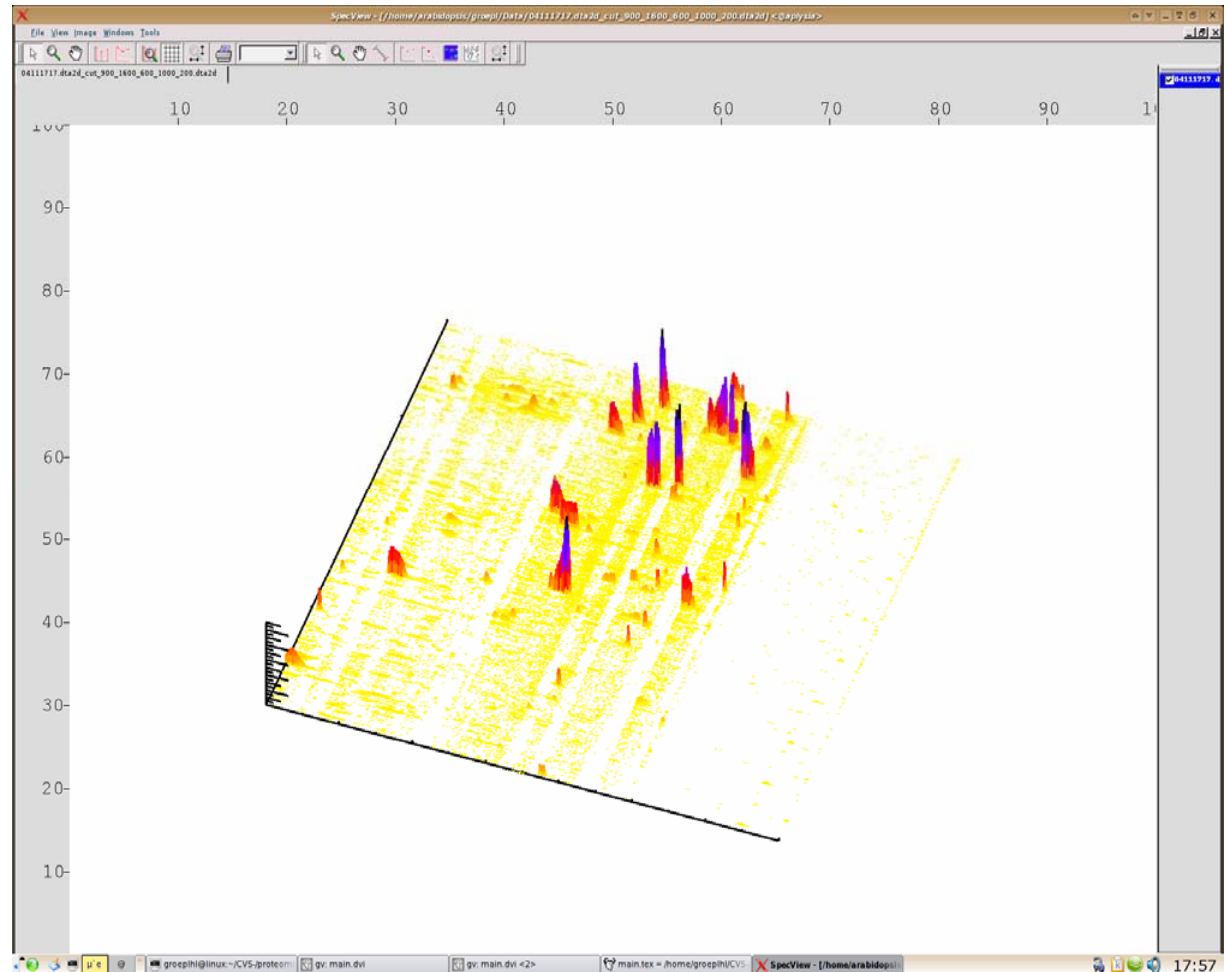


MDI Window

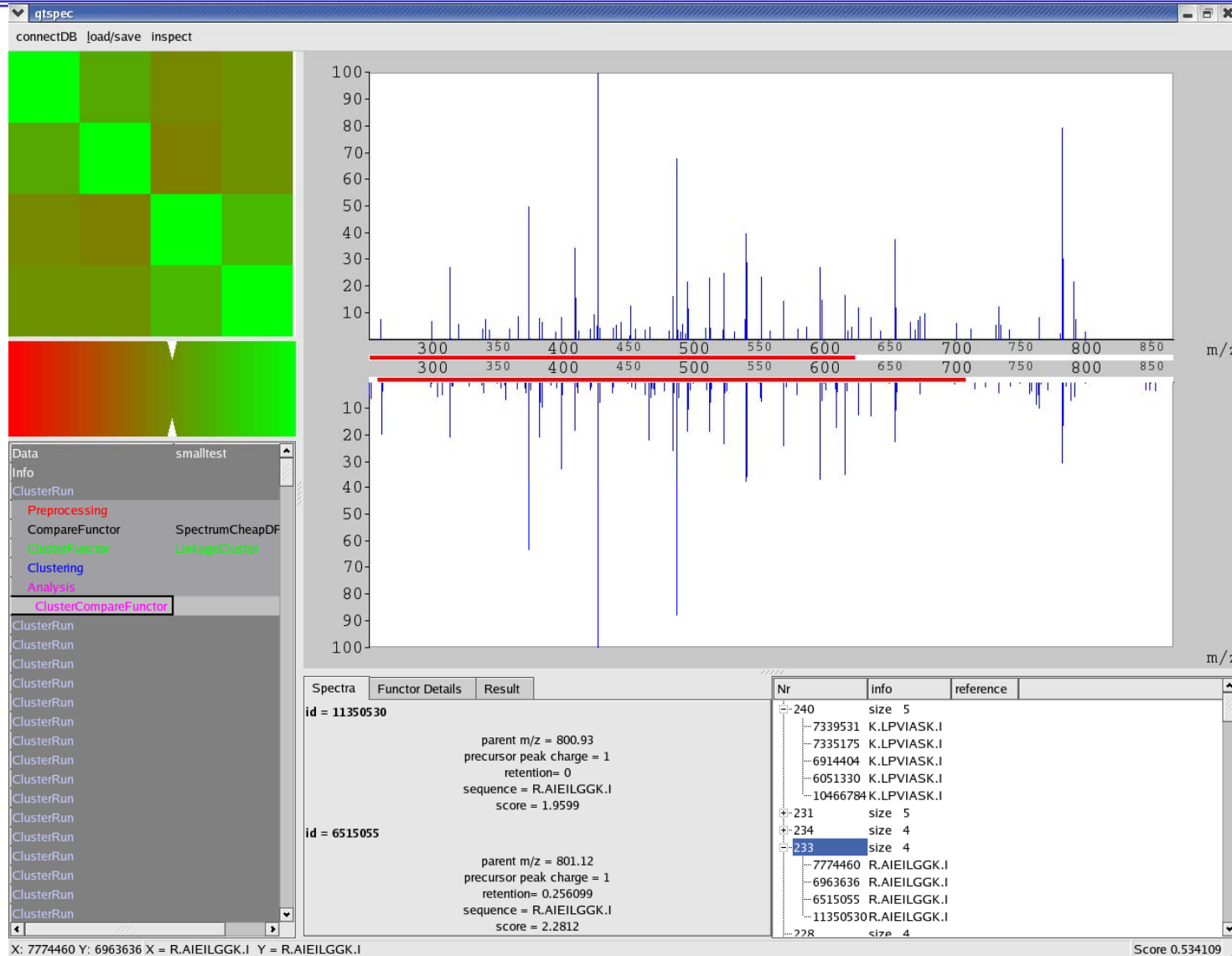
1D Window



- Visualization of raw, stick, and feature maps
- Rapid interactive zoom in/out through efficient quad tree implementation
- Data import from files or direct visualization from the DB



Clustering/ID View



Visualization Code Example

```
// create application with a single widget
QApplication app(argc, argv);
Spectrum1DWidget* wid = new Spectrum1DWidget;
app.setMainWidget(wid);

// load data from DB
DBAdapter dba;
dba.connectDB("OpenMS", "<user>", "<password>");
dba.executeQuery("SELECT id FROM PeakList ... ", true);
wid->setDataSet(dba.objects()[0]);

// execute application
wid->show();
return app.exec();
```

Acknowledgements

Dr. Clemens Gröpl
Eva Lange, Tim Conrad,
Ole Schulz-Trieglaff
*(Algorithmische Bioinformatik,
FU Berlin)*

Prof. Hartmut Schlüter
(Universitätsmedizin Berlin, Charité)

Prof. Dr. Oliver Kohlbacher
Marc Sturm, Andreas Bertsch
Jens Joachim
(SBS/WSI, Tübingen)

*Andreas Hildebrandt
(Uni Saarbrücken)*

Prof. Dr. Christian Huber
Bettina Mayr et al.
*(Instr. Analytik & Bioanalytik,
Univ. des Saarlandes, Saarbrücken)*

Dr. Albert Sickmann
(Virchow-Zentrum, Würzburg)

Herbert Thiele
Jens Decker
(Bruker Daltonics, Bremen)

Dr. Christoph Klein
(IRMM, Geel; now IHCP Ispra)

1. Überblick : DNA Sequenzierung
2. Überblick : BAC-by-BAC Assembler
3. Überblick : WGS Assembler
4. Überblick : Compartmentalized Assembler
5. Details der CSA Pipeline

Danke für die
Aufmerksamkeit !!!

Shotgun DNA Sequenzierung (Technologie)

