

# Analýza genomických a proteomických dat

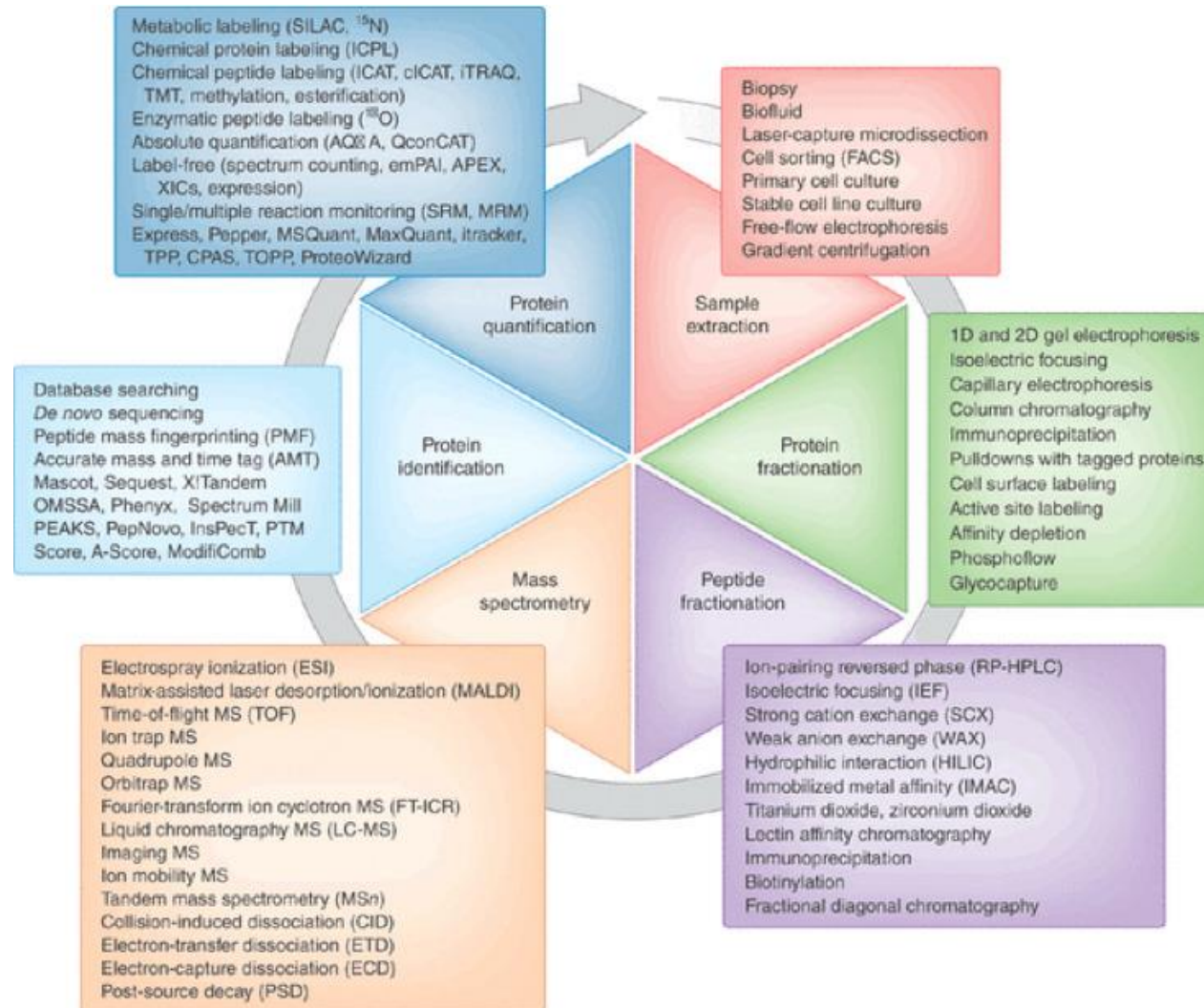
## Proteomické analýzy

Jaro 2022

19. duben 2022

Eva Budinská ([budinska@recetox.muni.cz](mailto:budinska@recetox.muni.cz))

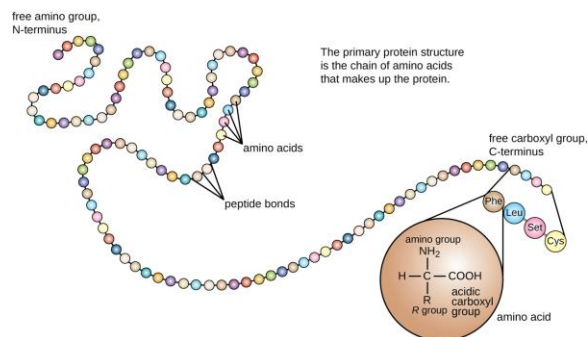
# Analýza proteomu



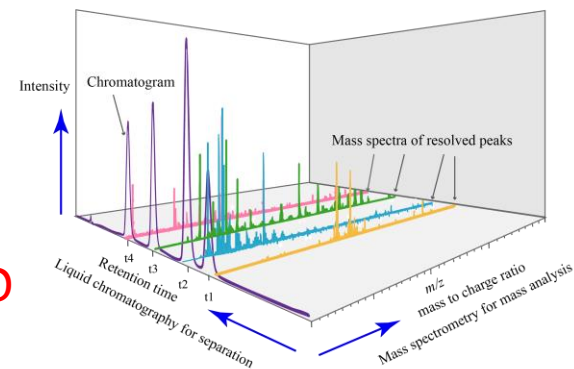
Zdroj: <https://www.creative-proteomics.com/services/proteomics-service.htm>

# Analýza proteomu

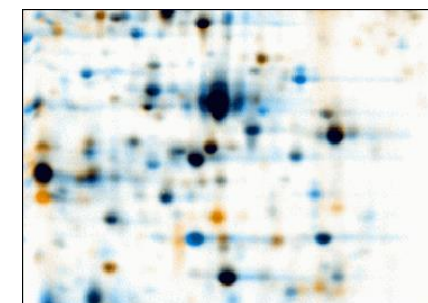
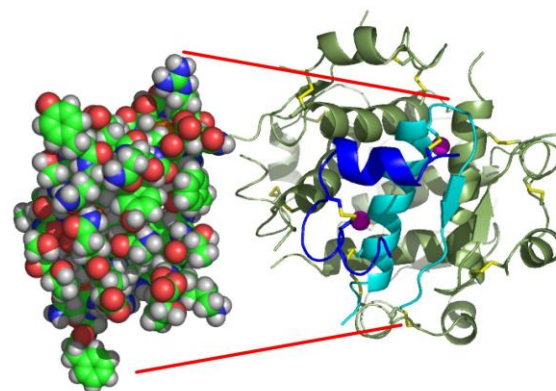
1. Analýza **struktury**:  
Proteínová sekvenace –  
Edmanova degradace,  
hmotnostní  
spektrometrie



2. Analýza **abundance**:  
**Hmotnostní**  
**spektrometrie**,  
proteínové mikročipy, **2D**  
**gelová elektroforéza**..



3. Analýza **funkce**: Modelování  
makromolekulárních systémů –  
odvozování vlastností z  
atomových interakcí



# Dělení proteomických experimentů

Dle toho **co zkoumají**: kvalita (struktura, funkce) nebo kvantita (abundance)

Dle **komplexity** vzorku: jeden, několik, tisíce?

Dle **úrovně frakcionace**: zkoumáme peptidy nebo proteiny?

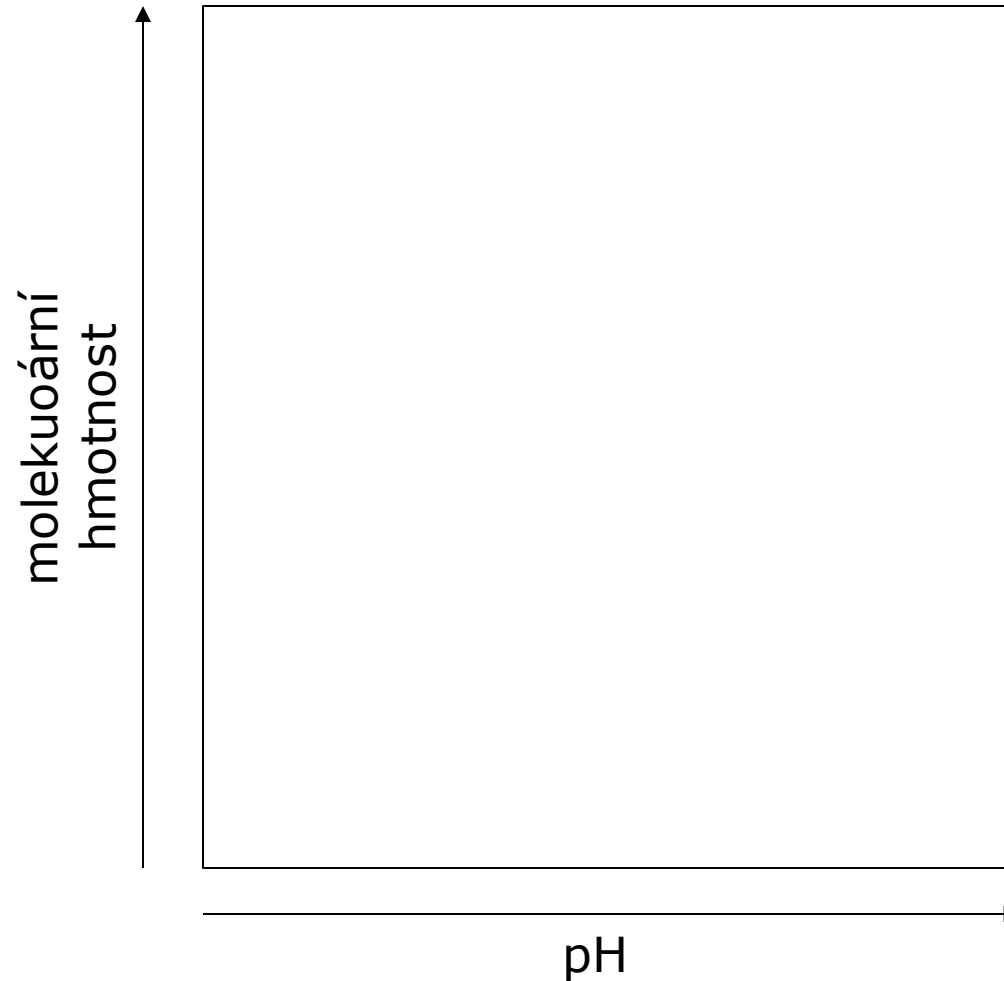
Dle **značení**: label-free, SILAC, iTRAQ, ...

Dle rozlišení: **nízké** vs **vysoké**

# 2D gelová elektroforéza

# Princip

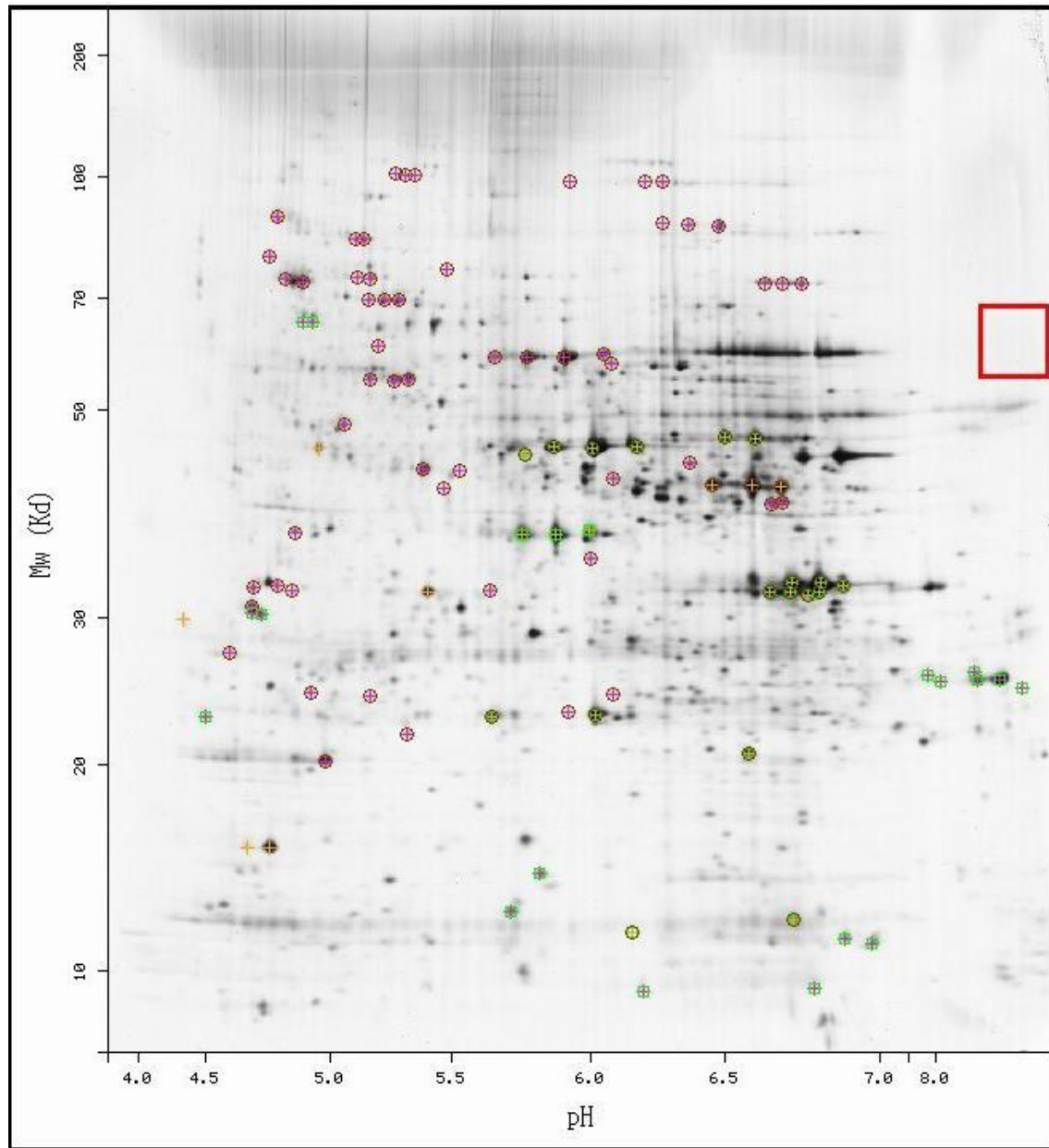
Proteiny jsou separované na gelu ve dvou dimenzích – na základě hmotnosti a na základě pH



# Postup experimentu

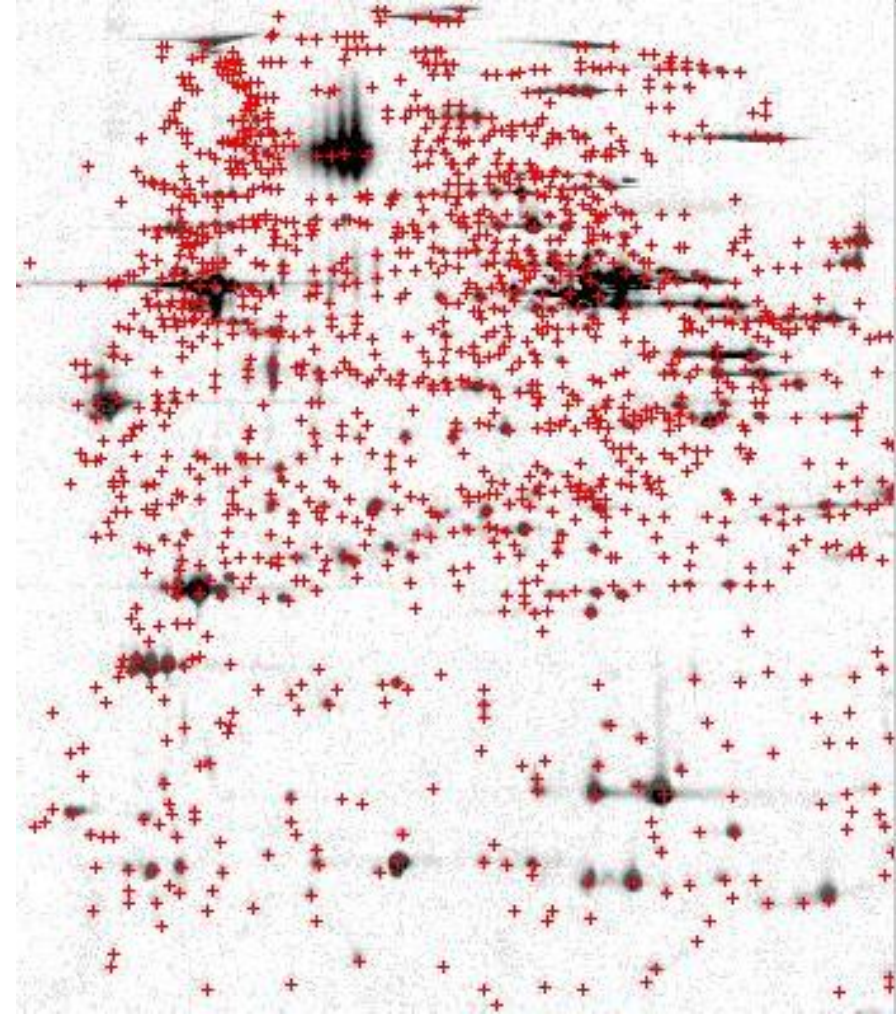
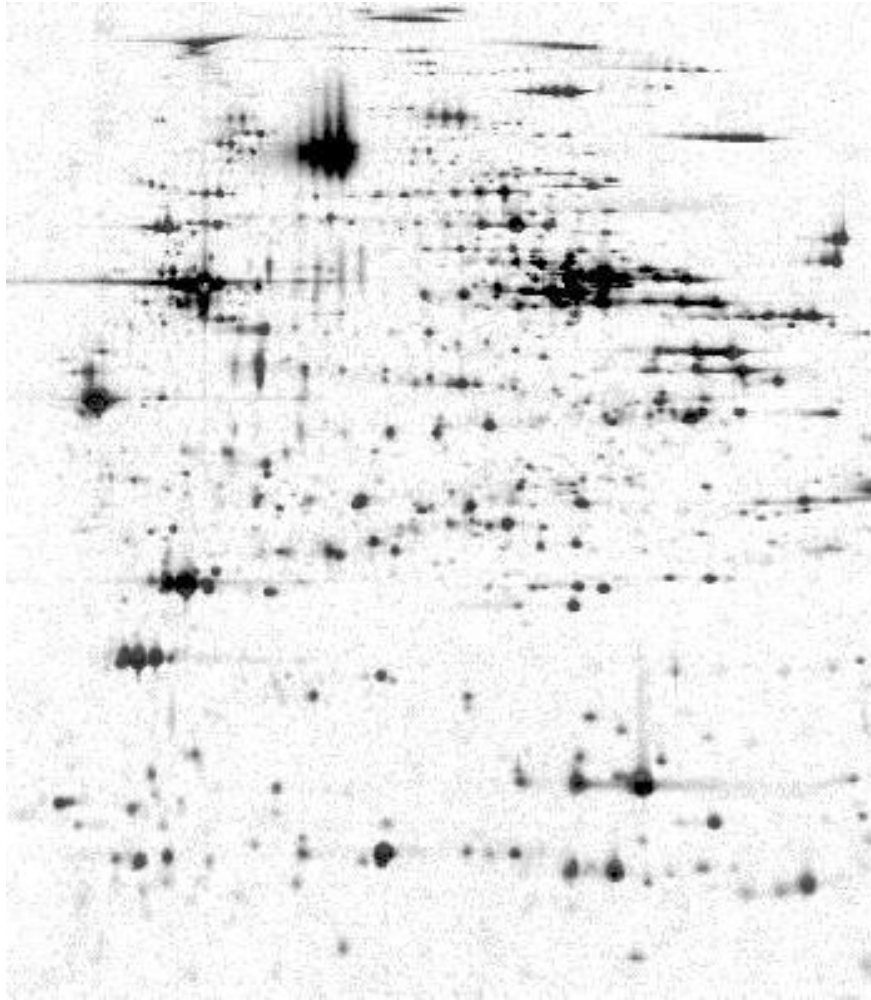
1. **Proteiny** jsou **extrahované** ze vzorku
2. **Vzorky se umístí na gel** a proteiny migrují až dosáhnou izoelektronický bod (kdy je jejich náboj nula) -  $pH(I)$ 
  - \* Důležitý je **výběr gelu**, který musí být dostatečně pórovitý, aby umožnil proteinům pohyb (agaróza nebo polyakrylamidový gel)
3. Takto se **proteiny oddělí** vzhledem ke svému izoelektrickému bodu
4. Následně proteiny necháme se pohybovat na základě **hmotnosti** ve druhé dimenzi
5. Nakonec je gel **zabarvený**, aby se detekovaly jednotlivé oblasti výskytu proteinů (spoty)
6. Zabarvený gel je pak digitalizován do obrazu (podobně jako mikročipy)
7. Intenzita pixelů koreluje s **množstvím proteinu**, používá se speciální SW pro analýzu spotů

# Jak vypadá obrázek





# 2-D gelová elektroforéza



# Jak vypadají data

← Vzorky →

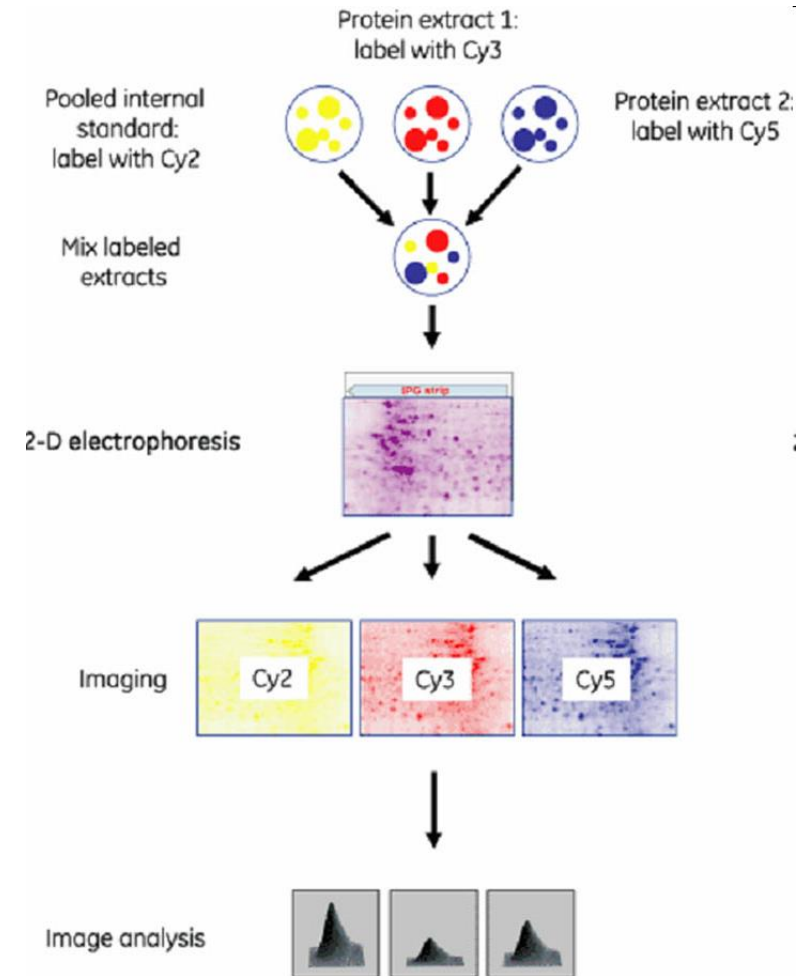
SSP	wt_A	wt_A	wt_A	wt_S
101	2338.84	2078.42	2625.1	2550.54
102	118.92	68.65	125.8	109.66
103	221.89	55.32	NA	NA
104	215.3	189.02	220.28	NA
105	106.56	NA	238.36	NA
202	328.32	226.46	522.52	1281.75
203	259.8	228.13	340.37	NA
205	1439.72	1213.28	1187.43	1353.14
206	1094.33	754.83	1291.89	1240.82
208	97.78	41.51	164.49	33.25
209	NA	NA	NA	22.42
301	212.63	92.12	307.19	317.67
302	1491.34	1703.79	1830.19	1976.66
304	71.25	72.72	127.87	199.31

Peptidy

# DIGE

\* Speciální typ 2-DE je 2-D Fluorescence Difference Gel Electrophoresis (DIGE).

- Proteiny se nejprve zabarví fluorescenčním barvivem
- Každé barvivo se skenuje pod jiným filtrem
- Takto se může porovnávat více vzorků

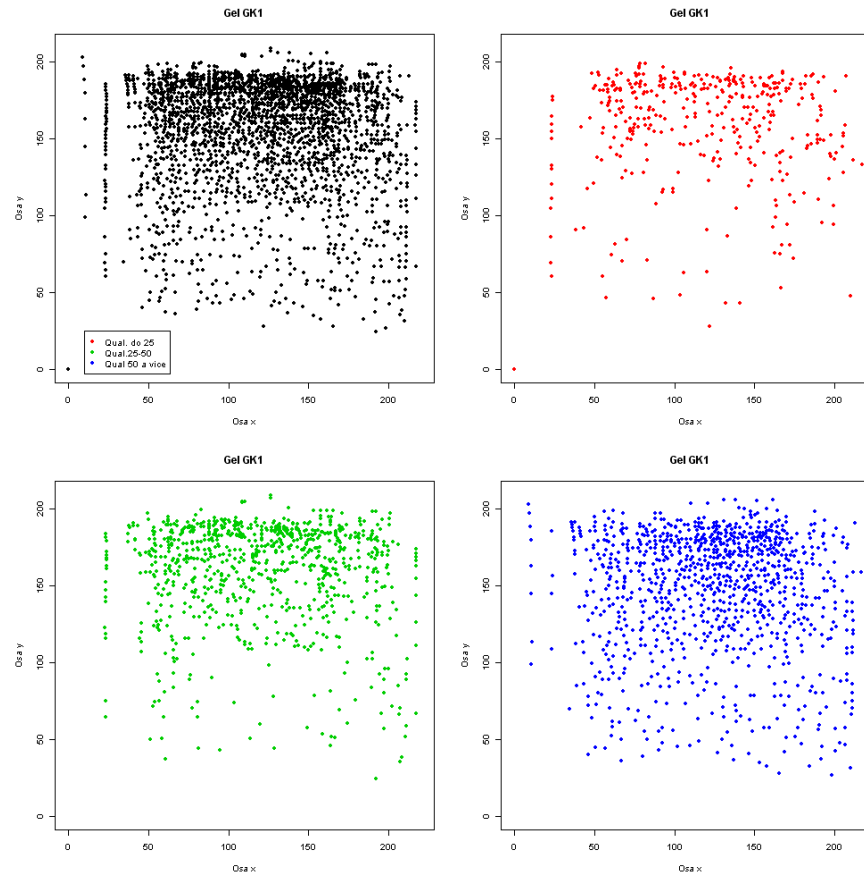


# Nutnost úpravy dat

- \* Tak jako mikročipový experiment i 2-DE je vystavená experimentálním chybám, které jsou zdrojem šumu
- \* Je nutná úprava a normalizace dat
- \* Neexistuje tu ale taková automatická kvantifikace spotů tak jako u mikročipů, protože spoty nejsou fixně dané předem!
  - \* existující automatická kvantifikace vyžaduje manuální úpravu
  - \* proměnné kvality spotů
- \* Data z 2DE nejsou normálně rozložená – je nutná transformace (log)

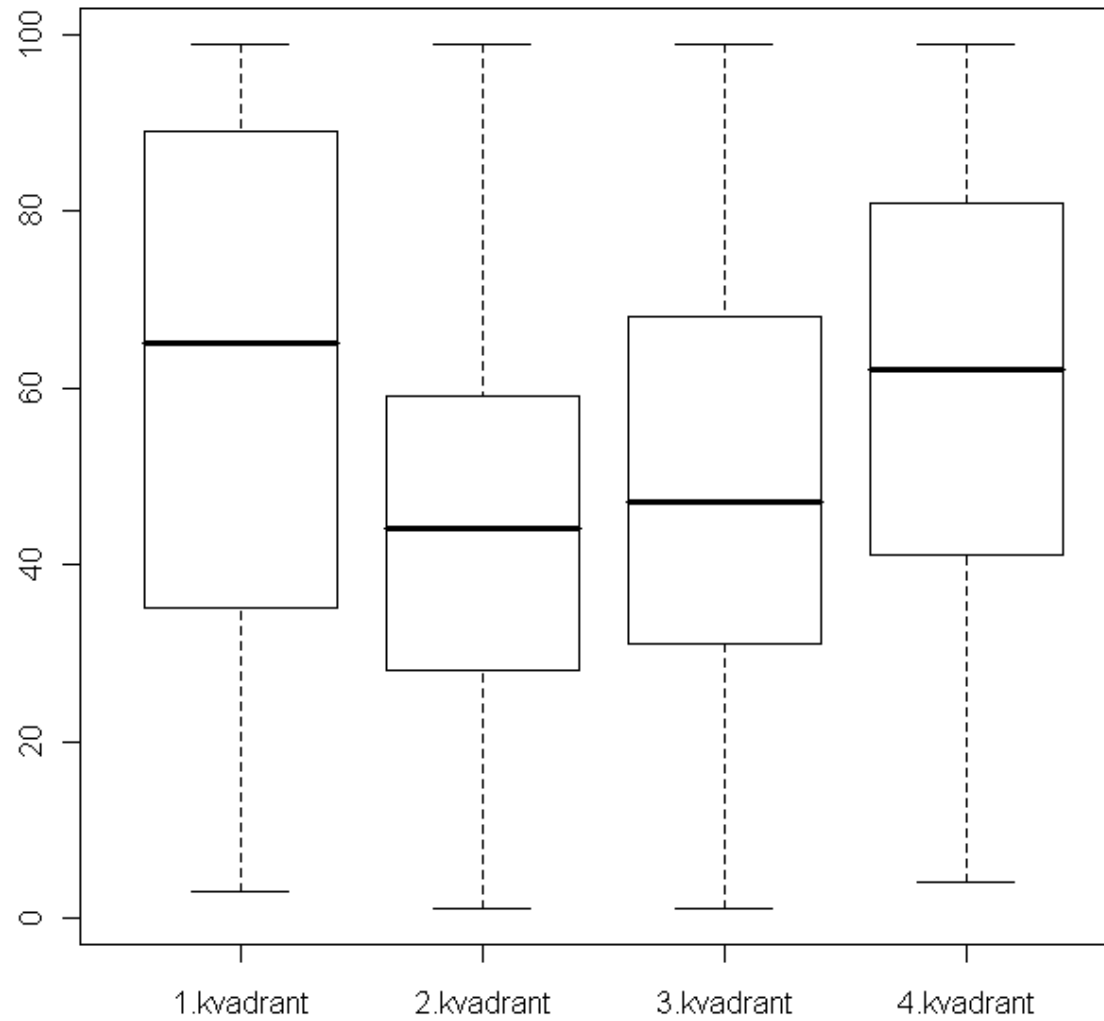
# Normalizace a úpravy dat

- \* Důležitým krokem v úpravě dat je **kalibrace všech expresních hodnot a gelů navzájem**
- \* V tomto procesu se odstraňuje prostorový efekt, i efekt barviva
- \* Na každém gelu jsou kontrolní proteiny, podle kterých se každý gel kalibruje (posouvá)



# Kontrola kvality spotů

Spot quality (N1=53, N2=598, N3=1217, N4=105)



# Hmotnostní spektrometrie

# Hmotnostní spektrometrie

Technika používaná pro charakterizaci (nejen) proteomu v biologickém vzorku (plasma, sérum, . . .)  
různé konzistence (pevná konzistence, tekutina, plyn)

Založená na rozdílném náboji a hmotnosti peptidů a proteinů (nebo jiných molekul)

Hmotnostní spektrometr je separuje na základě poměru **hmotnosti k náboji** (anglicky *mass to charge ratio* –  $m/z$ , jednotka Dalton), který je specifický pro každou molekulu.

Často používané systémy - TOF nebo Orbitrap



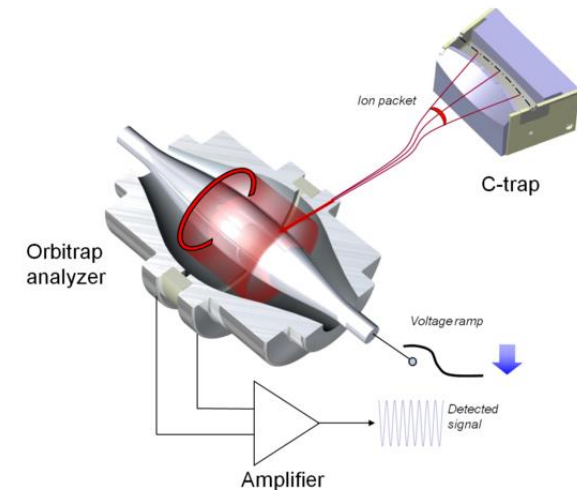
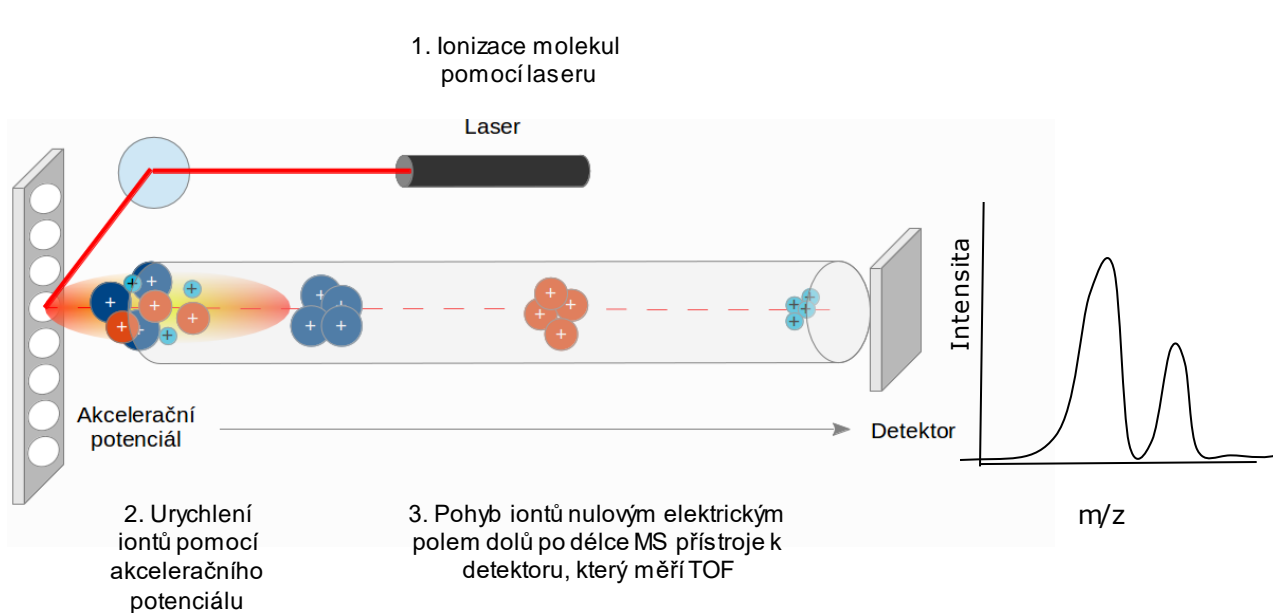
# Hmotnostní spektrometrie

Technika používaná pro charakterizaci (nejen) proteomu v biologickém vzorku (plasma, sérum, . . .) různé konzistence (pevná konzistence, tekutina, plyn)

Založená na rozdílném náboji a hmotnosti peptidů a proteinů (nebo jiných molekul)

Hmotnostní spektrometr je separuje na základě poměru **hmotnosti k náboji** (anglicky *mass to charge ratio* –  $m/z$ , jednotka Dalton), který je specifický pro každou molekulu.

Často používané systémy - TOF nebo Orbitrap



# Hmotnostní spektrometr TOF - princip

- TOF (time-of-flight) závisí na hmotnosti proteinů nebo přesněji na jejich  $m/z$  a představuje sumu těchto časů:

$$TOF = t_a + t_D + t_d$$

je čas letu v akcelerační oblasti,  $t_a$

je čas přeletu v oblasti s nulovým elektrickým polem  $t_D$

je čas detekce  $t_d$

- $TOF$  lze aproximovat pouze pomocí  $t_D$

*mass-to-charge ratio* je vypočteno podle:

$$m/z = B(t_D - A)^2$$

- $A$  a  $B$  jsou stanoveny pomocí kalibrace

# Hmotnostní spektrometr TOF - druhy

## ■ Příklady TOF spektrometrů:

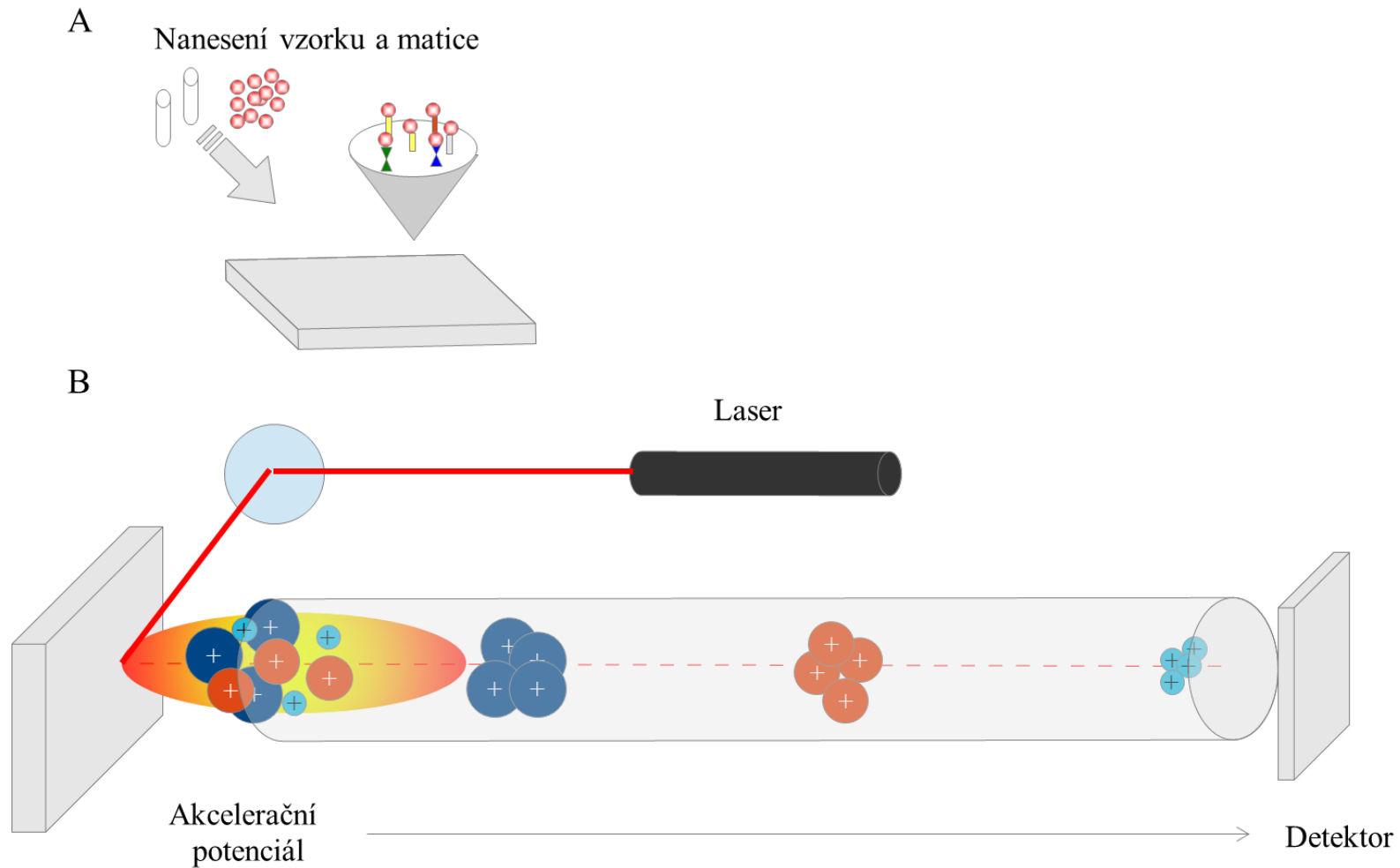
- Matrix-Assisted Laser Desorption-Ionisation (MALDI)-TOF
- Surface-Enhanced Laser Desorption-Ionisation (SELDI)-TOF

## ■ Způsob uchycení proteinů a ionizace

- Proteiny vzorku jsou před samotnou analýzou upevněny na podklad, který se v závislosti od typu hmotnostní spektrometrie liší.
- Jeho úkolem je také absorbovat energii v ionizátoru a předat ji vzorku a tak usnadnit jeho ionizaci.
- u MALDI se jedná o energii-absorbující matrici (matrix), co je nejčastěji organická kyselina s aromatickým jádrem
- SELDI využívá proteinový čip (s několika - obvykle osmi - spoty), opatřen speciálním chromatografickým povrchem, takže se na povrch váží různé proteiny v závislosti na svých chemických vlastnostech a vlastnostech čipu. A až potom dojde k nanesení matrice, která se vzorkem vytvoří krystaly.

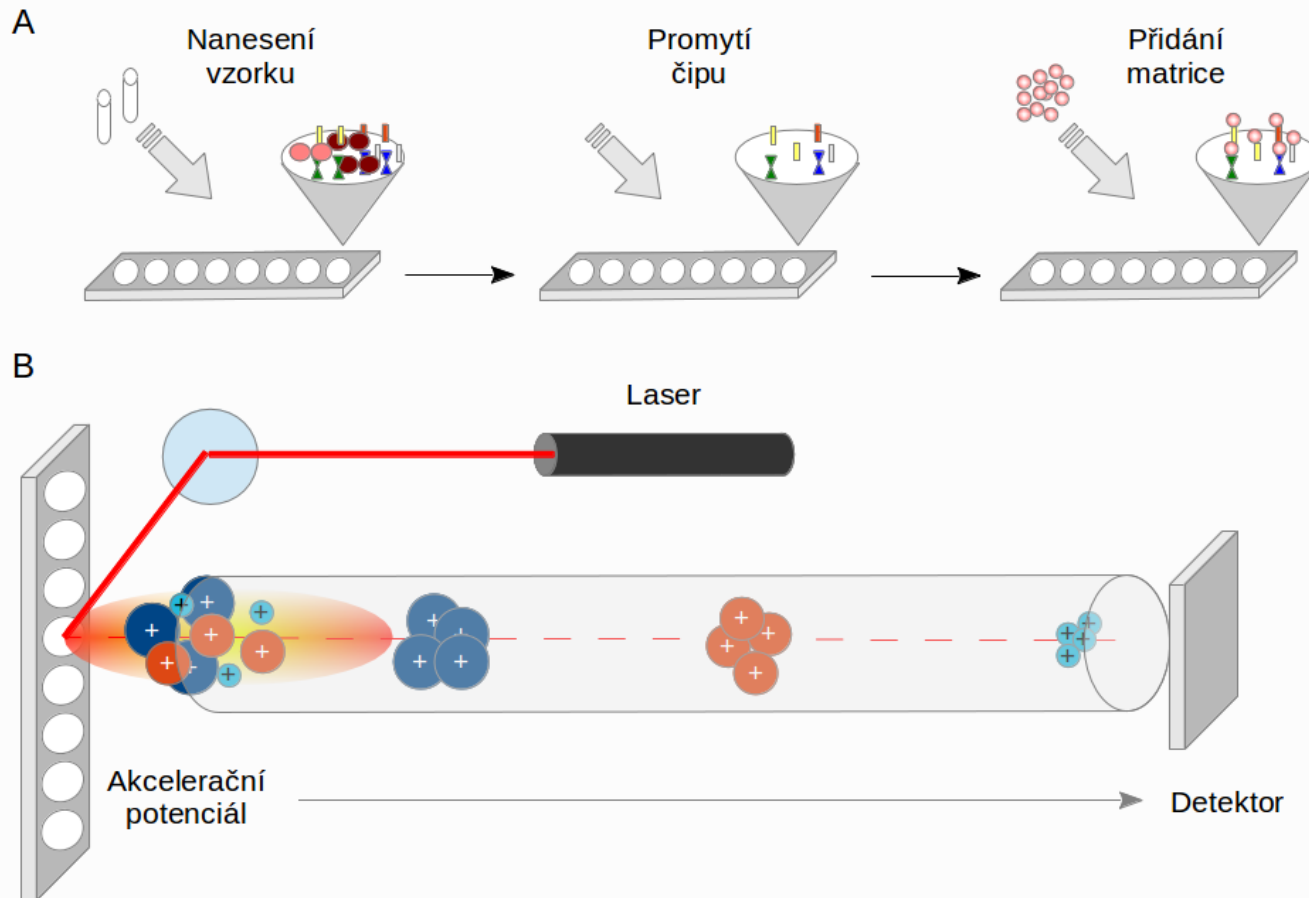
# MALDI-TOF

- Matrix-Assisted Laser Desorption-Ionisation - TOF



# SELDI-TOF

- **Surface-Enhanced Laser Desorption-Ionisation – TOF**
- Existuje několik druhů čipů (IMAC30, H50, NP20...), které se liší svým aktivním povrchem (anionický, kationický, kovový, normální fáze, hydrofobický, ...) a proto také přednostně vážou jiné molekuly.

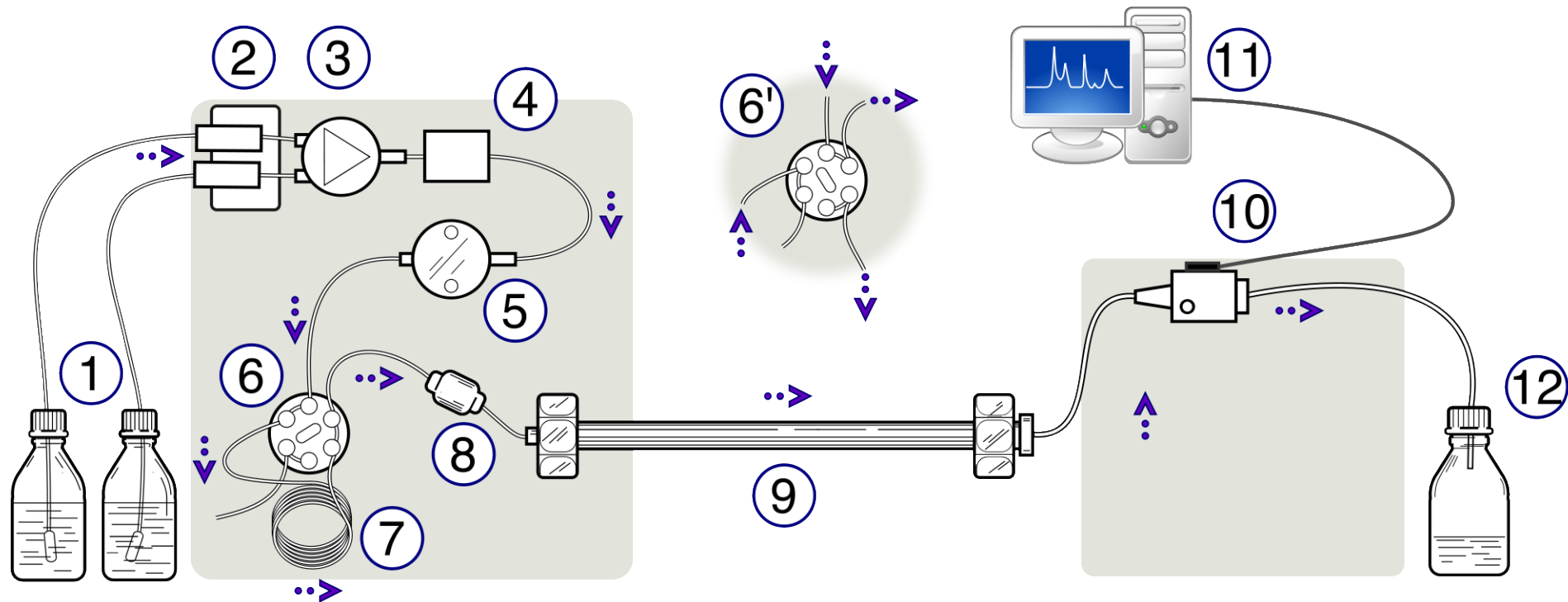


## Výhoda SELDI:

Možnost odmyt látky, které by jinak ovlivňovali spektrum vzorku (např. močovina používaná k přípravě vzorku, nebo  $\text{Na}^+$  ionty přítomné fyziologicky ve vzorcích).

# Kapalinová chromatografie – LC-MS/MS

- Další druh hmotnostní spektrometrie pro identifikaci proteinů
- Vzorky nejsou na matrici jako u MALDI nebo SELDI, ale v kapalině
- MS/MS - tandemová hmotnostní spektrometrie



# Tandemová hmotnostní spektrometrie (MS/MS)

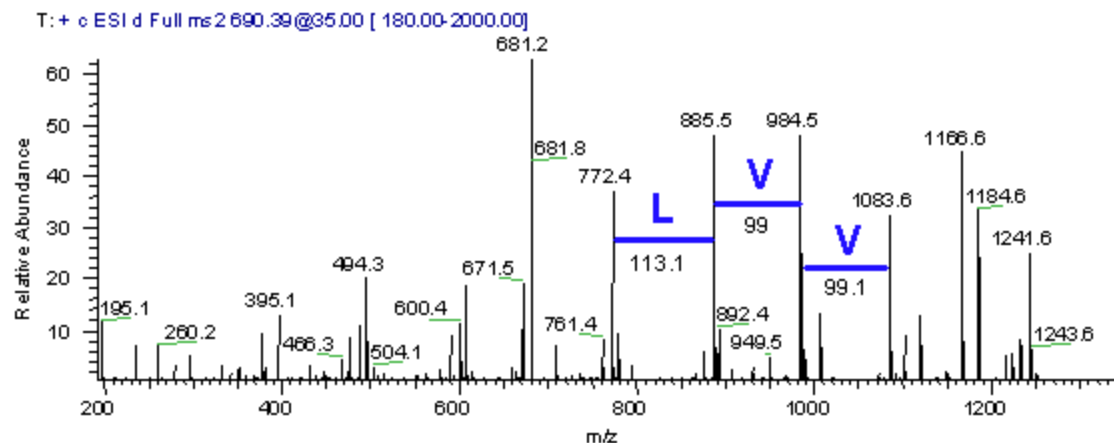
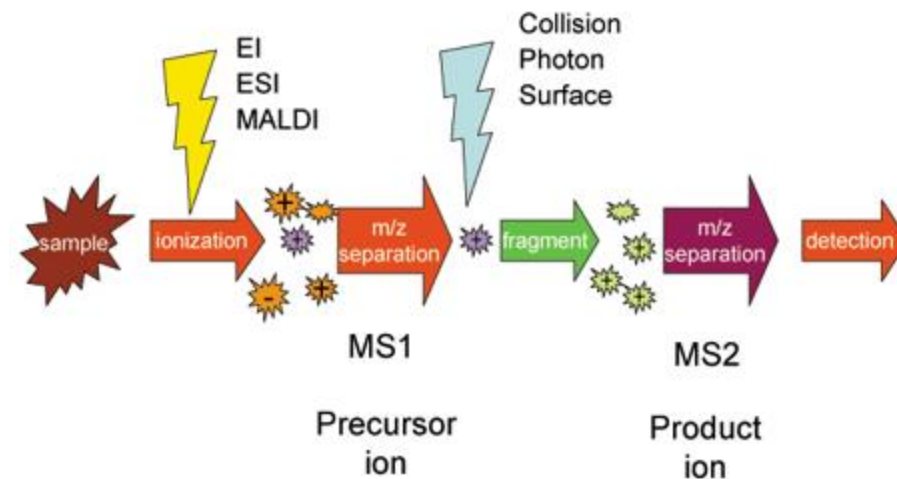
Jde o použití dvou spektrometrií jednu po druhé.

1. **Molekuly vzorku** jsou ionizovány a první spektrometr (označený MS1) odděluje tyto ionty podle jejich poměru hmotnosti k náboji ( $m/z$ ).
2. Ionty s konkrétním poměrem  $m/z$  pocházející z MS1 jsou vybrány a poté se **rozštěpí (fragmentují)** na menší ionty (kolizí indukovanou disociací, reakcí iontů a molekul nebo fotodisociací).
3. Tyto **fragmenty** dále putují do druhého hmotnostního spektrometru (MS2), který dále odděluje fragmenty podle jejich poměru  $m/z$  a detekuje je.

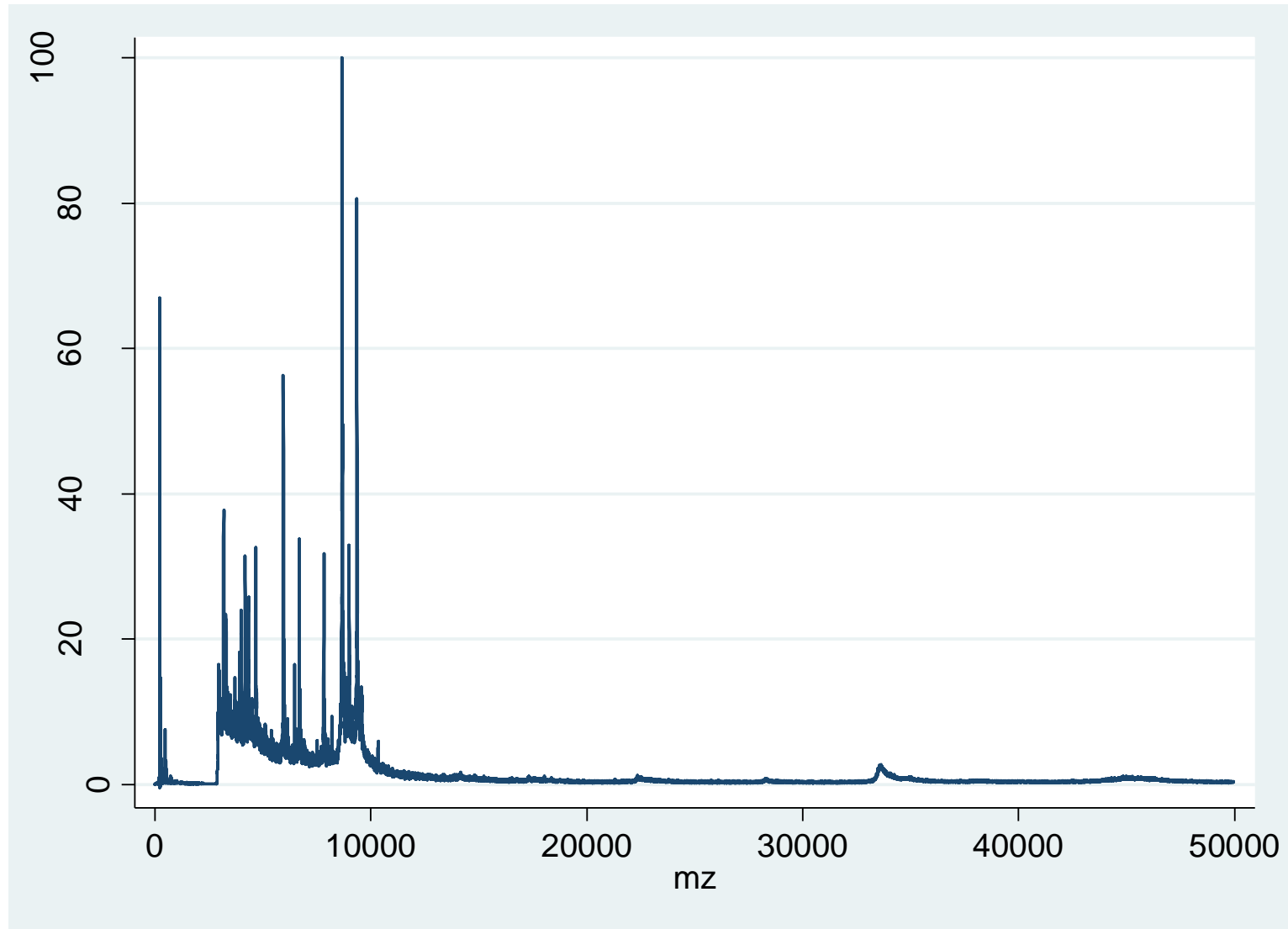
Fragmentační krok umožňuje identifikovat a separovat ionty, které mají velmi podobné  $m/z$ -poměry v běžných hmotnostních spektrometrech.

Používá se v proteínovém sekvencování.

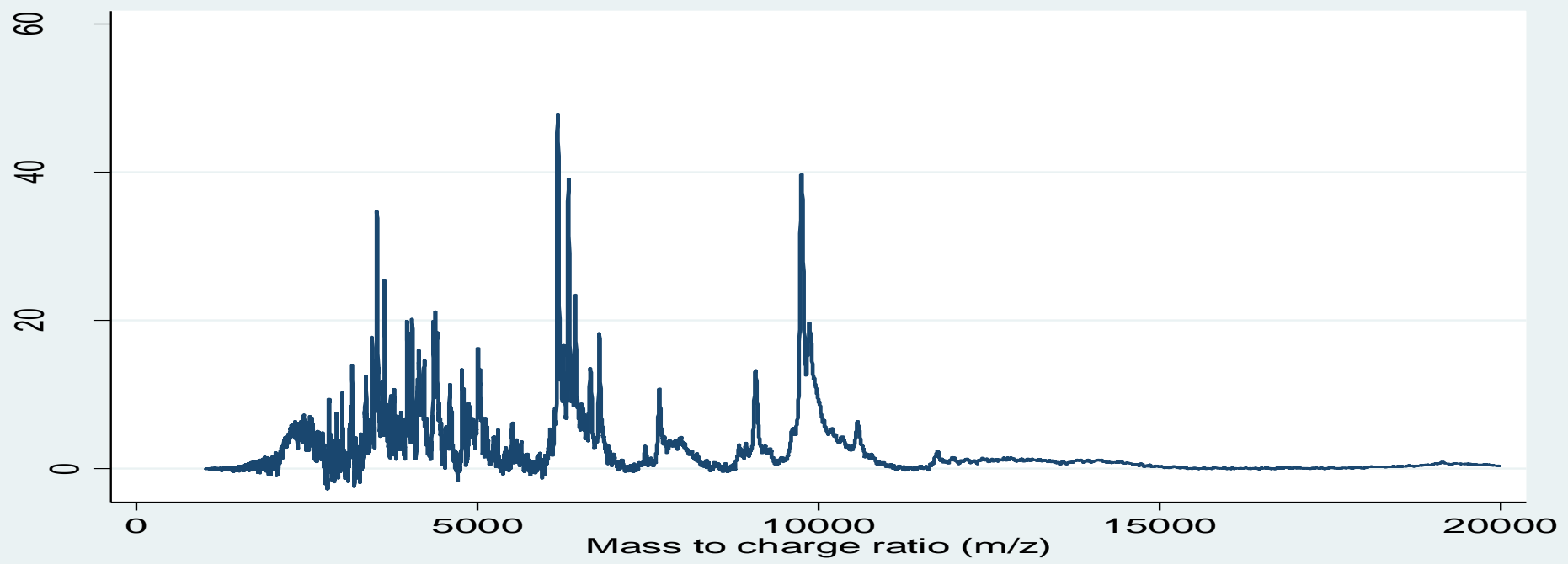
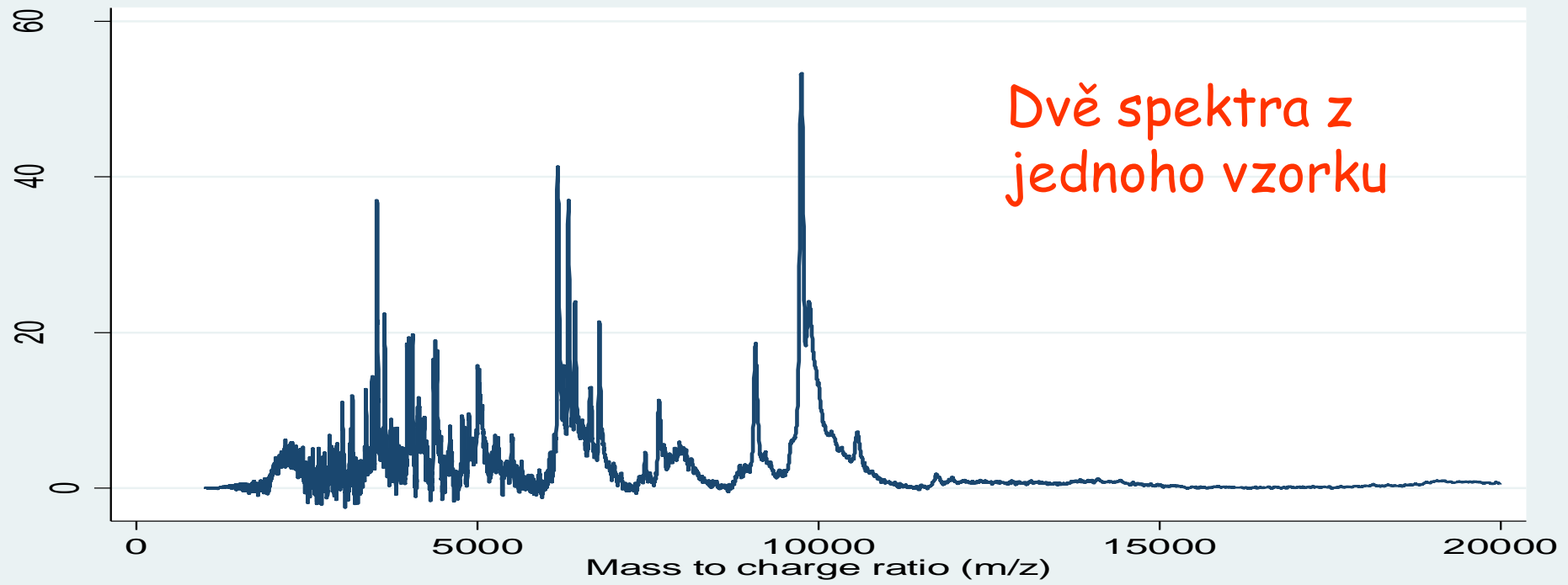
Výsledkem je **peptide sequence tag**



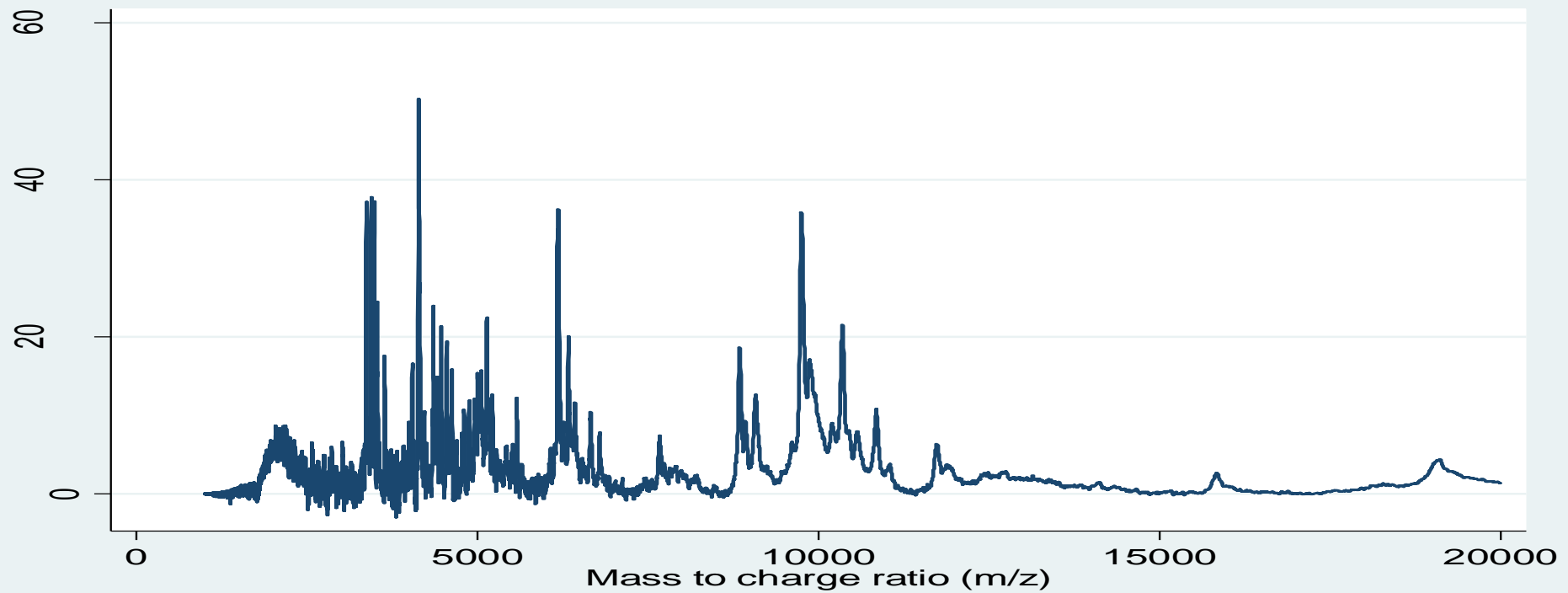
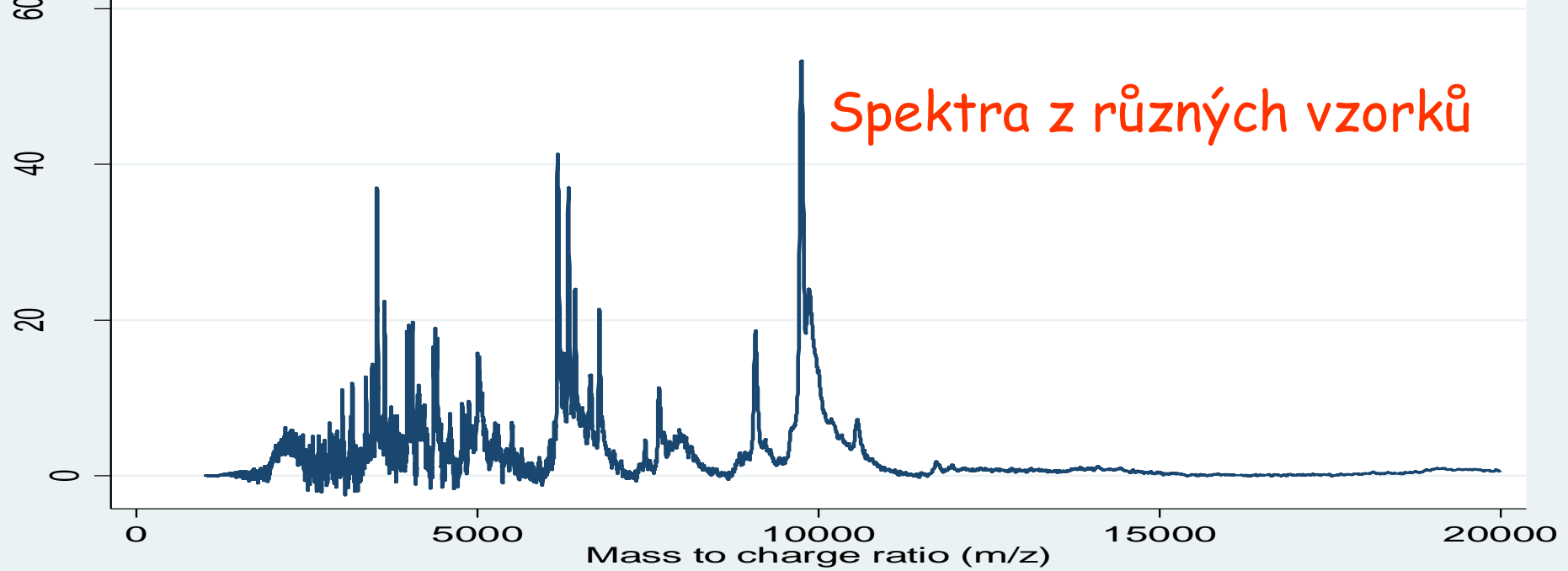
# Jak vypadají data vzorku z hmotnostního spektrometru







# Spektra z různých vzorků



# Zpracování dat

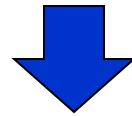
---

1. úprava hrubých dat (MS/MS i MS), normalizace, identifikace píků

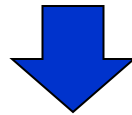
# Vznik a úprava dat

- **Kalibrace**

- Signál je přeměněný na škálu m/z pomocí množství kalibračních proteinů ze známou m/z hodnotou. Toto se děje ještě v přístroji



Základní data (formáty raw, mzXML, mzML, ...)

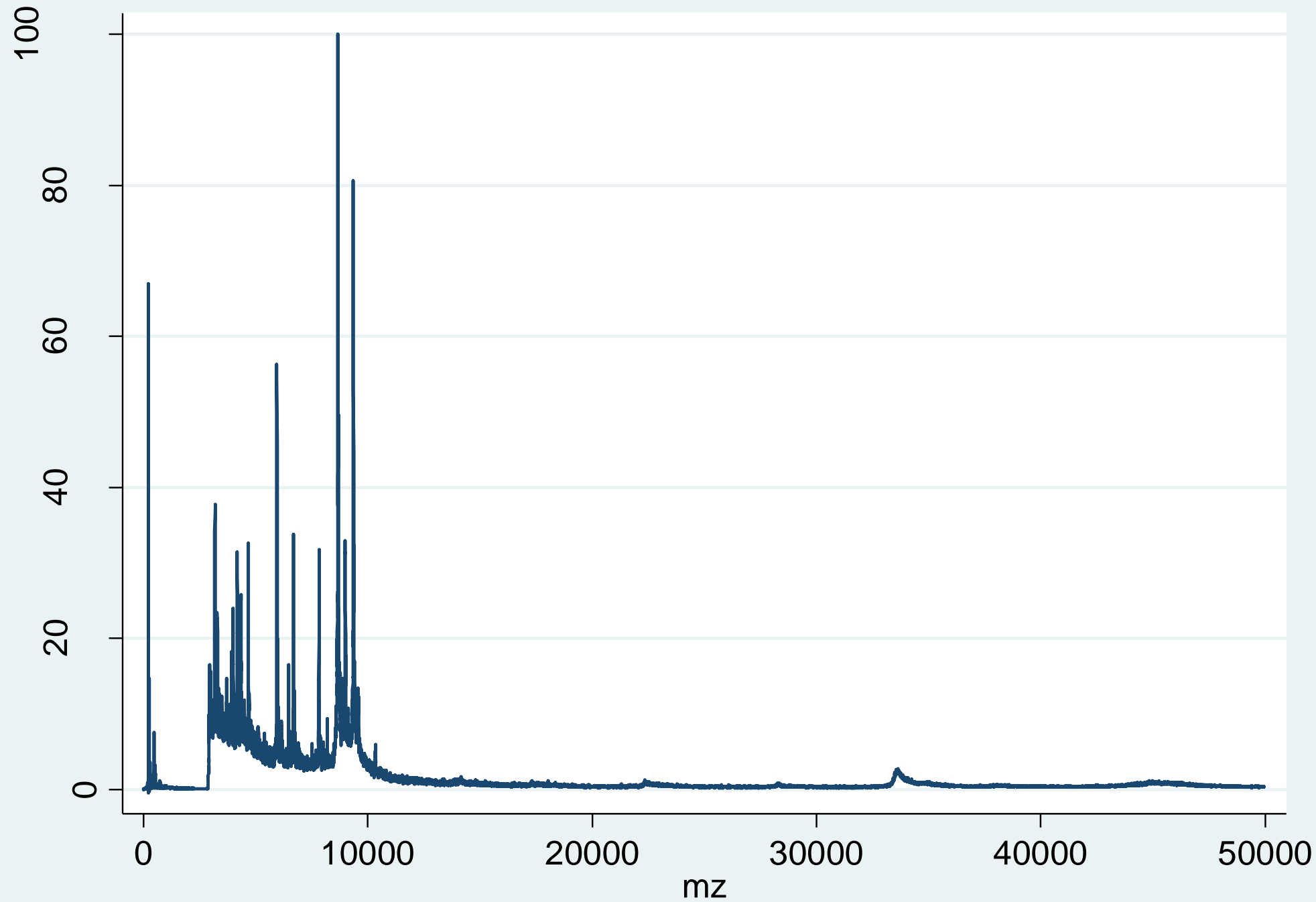


- **Odstranění baseline**

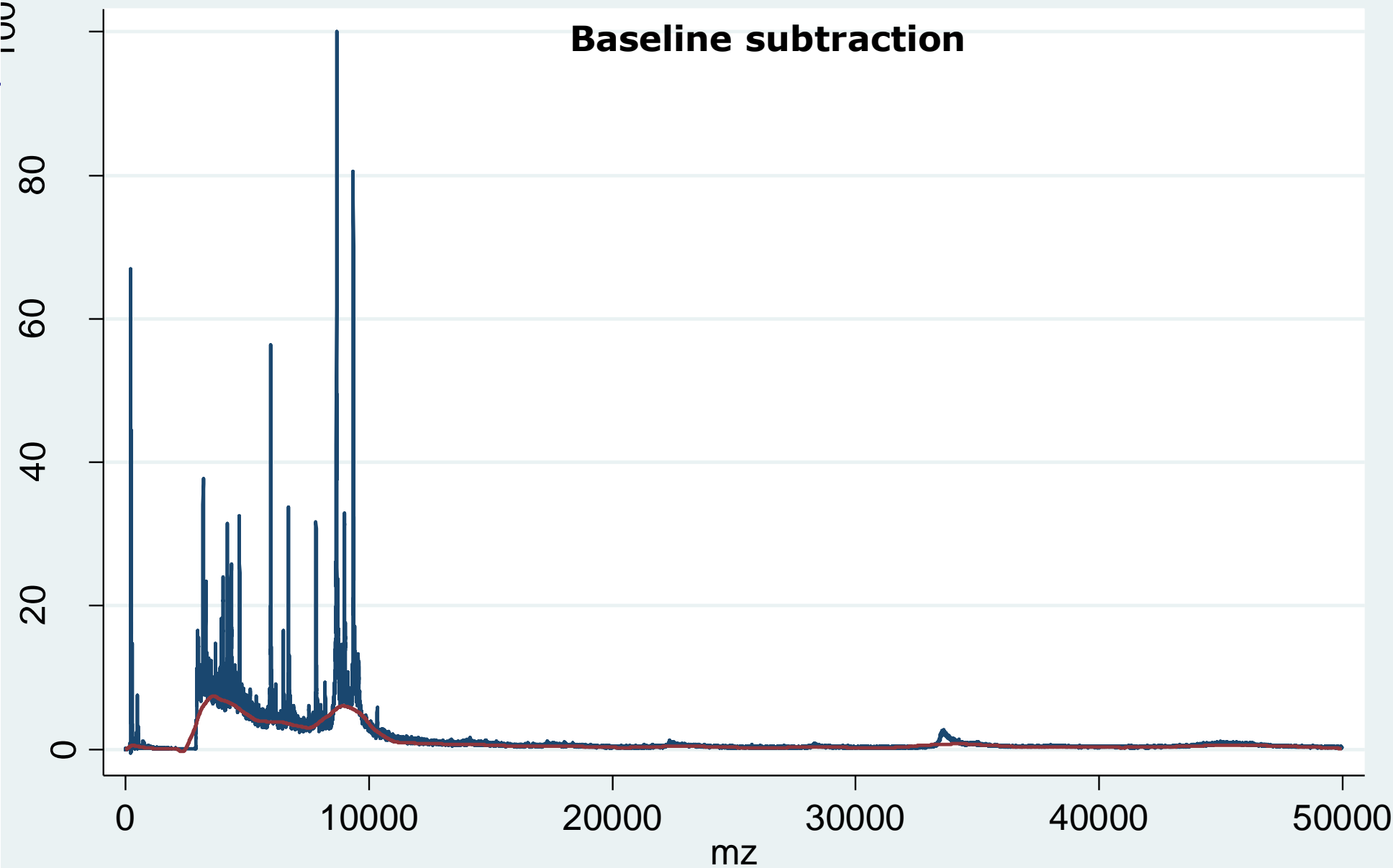
- Odstranění baseline šumu z profilu, například pomocí **loess**

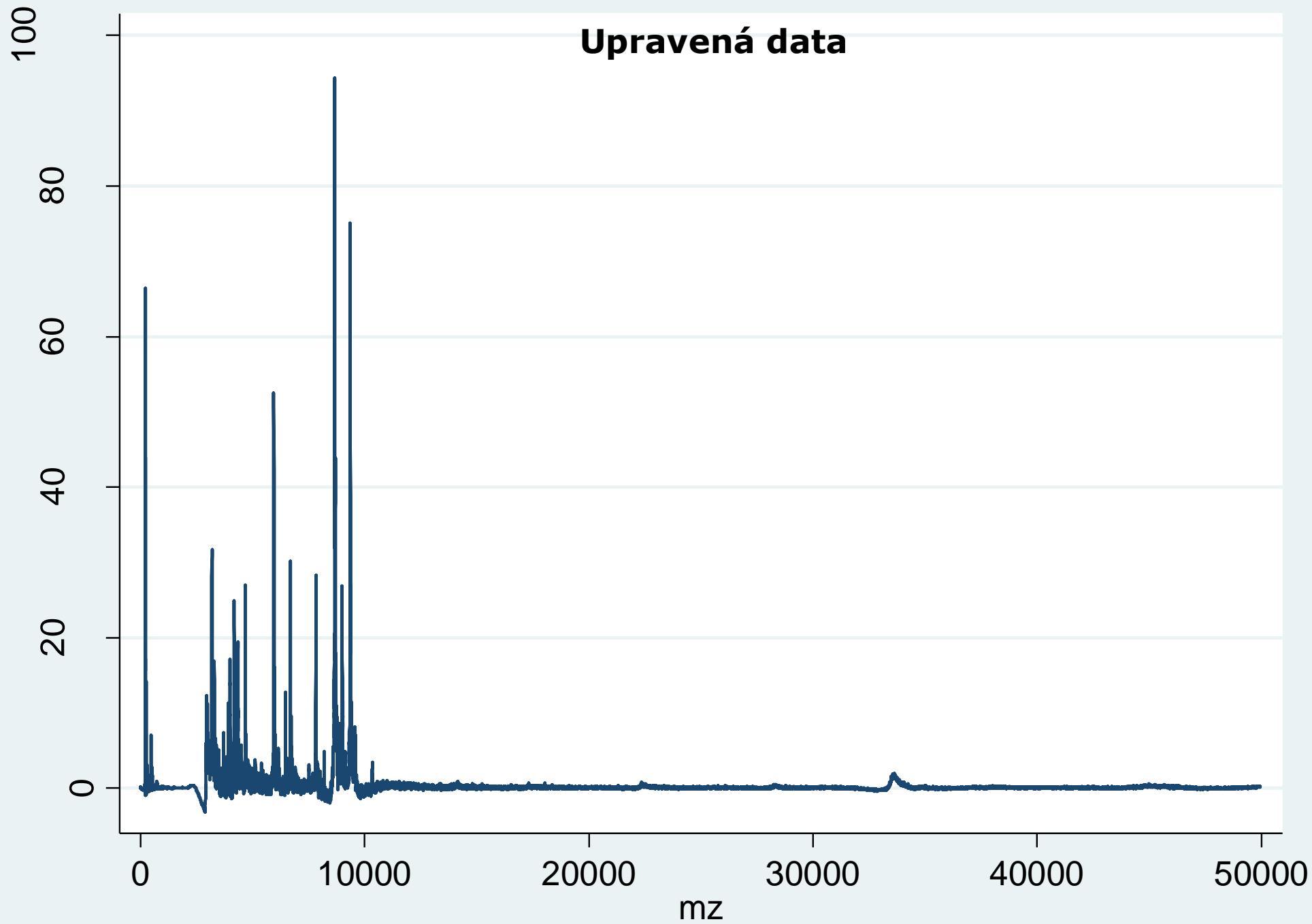
- **Normalizace**

- Abychom mohli porovnat spektra mezi vzorky



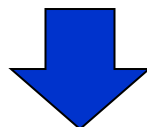
# Baseline subtraction





# Normalizace

- Odstraňujeme technickou variabilitu (přístrojové chyby, odlišné množství vzorku)
- Koncentrace proteinu se odhaduje jako plocha pod píkem (Area Under Curve – AUC)

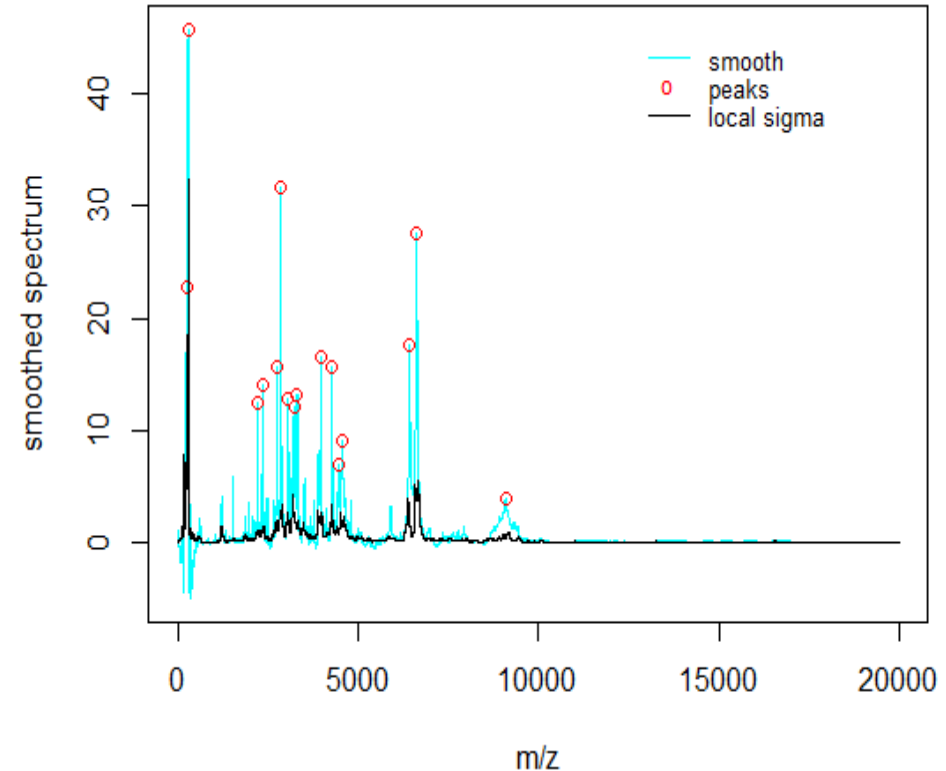


- Normalizace pomocí *průměrné AUC (TIC – total ion current)*  
AUC celého spektra / průměrná AUC všech spekter



# Detekce píků a jejich zarovnání

- Pík ~ peptid/proteín, definuje se jako lokální maximum na základě porovnání variability v okolí
- Existují nepřesnosti na x ( $m/z$ ) a y (signál) osách
- Píky každého spektra můžou být definované jako body které jsou maximálně  $\pm N$  bodů v okolí  $m/z$ 
  - first, second, estimated..
- Důležité je brát do úvahy *signal-to-noise* ratio – píky musí překročit nějakou běžnou hranici šumu



# Jak vypadají data po zarovnání a detekci píků

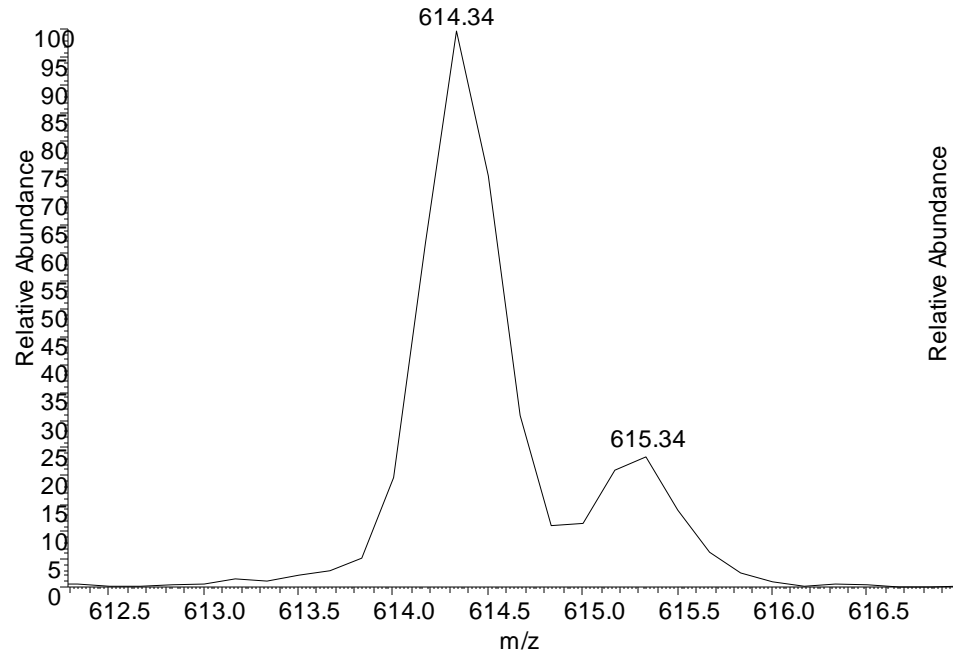
## ■ SELDI-TOF

Cluster	Group	Norm. Log Intensity	M/Z	Intensity	Norm. Linear Intensity	Type	Mass Dev.
1	chemoresistentni	0.581550	2392.84	3.058176	30.578211	estimated	0.000007
1	chemoresistentni	-0.072123	2392.84	1.943959	12.984676	estimated	0.000007
1	chemoresistentni	0.023116	2392.84	2.076621	15.079403	estimated	0.000007
1	chemoresistentni	0.160910	2392.84	2.284742	18.365652	estimated	0.000007
1	chemoresistentni	0.199591	2392.84	2.346828	19.345988	estimated	0.000007
1	chemoresistentni	0.161331	2392.82	2.285410	18.376190	first	-0.000004

# Úprava hrubých dat - dvě možnosti

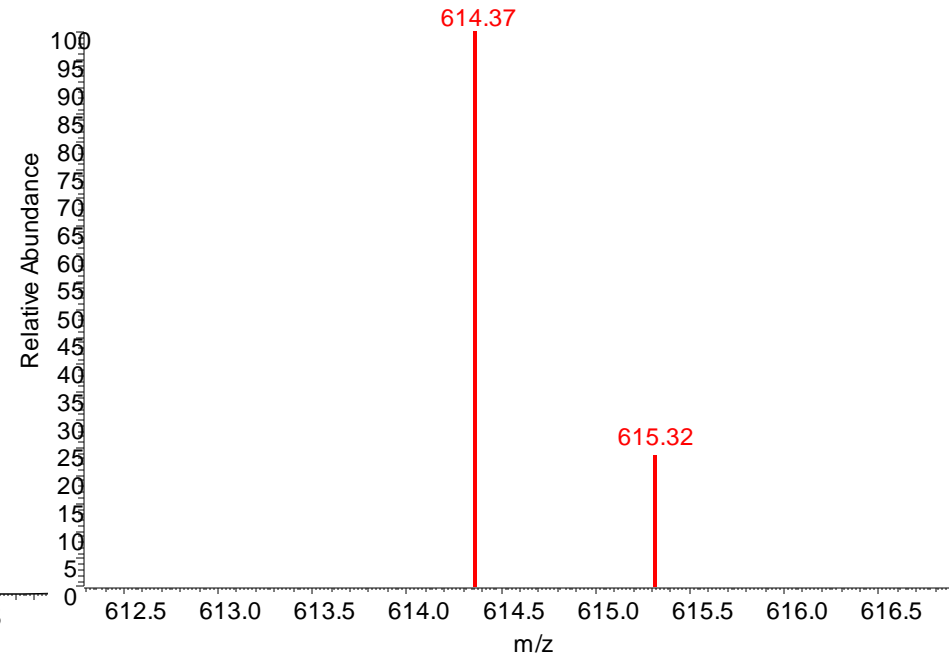
## Profilové spektrum

- Získané z experimentu



## Čárové spektrum

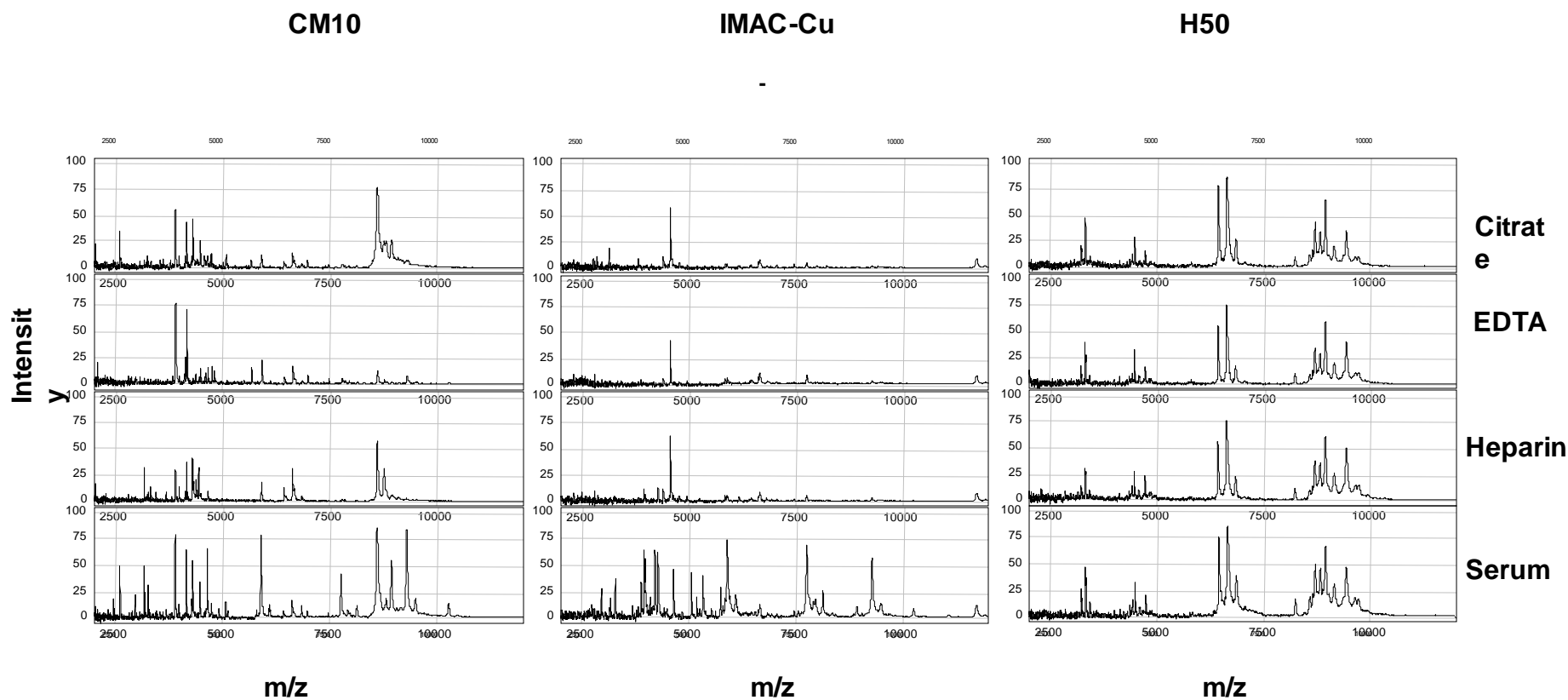
- Vypočítané z profilového



# Čipy pro SELDI hmotnostní spektrometrii

- Kvantitativní hodnoty proteomu jsou také ovlivněné různými zdroji variability (experimentální i biologické)
- Velmi velké rozdíly mezi typy použitého čipu!

	<u>H50</u>		<u>IMAC30</u>		<u>NP20acid</u>		<u>NP20alkaline</u>	
	<u>N</u>	<u>%</u>	<u>N</u>	<u>%</u>	<u>N</u>	<u>%</u>	<u>N</u>	<u>%</u>
<u>H50</u>	<b>75</b>	<b>100.0</b>	<b>19</b>	<b>47.5</b>	<b>24</b>	<b>52.2</b>	<b>56</b>	<b>59.6</b>
<u>IMAC30</u>	<b>19</b>	<b>25.3</b>	<b>40</b>	<b>100.0</b>	<b>19</b>	<b>41.3</b>	<b>21</b>	<b>22.3</b>
<u>NP20acid</u>	<b>24</b>	<b>32.0</b>	<b>19</b>	<b>47.5</b>	<b>46</b>	<b>100.0</b>	<b>30</b>	<b>31.9</b>
<u>NP20alkaline</u>	<b>56</b>	<b>74.7</b>	<b>21</b>	<b>52.5</b>	<b>30</b>	<b>65.2</b>	<b>94</b>	<b>100.0</b>
<u>separate M/Z</u>	<b>15</b>	<b>20.0%</b>	<b>15</b>	<b>37.5%</b>	<b>12</b>	<b>30.0%</b>	<b>28</b>	<b>29.8%</b>

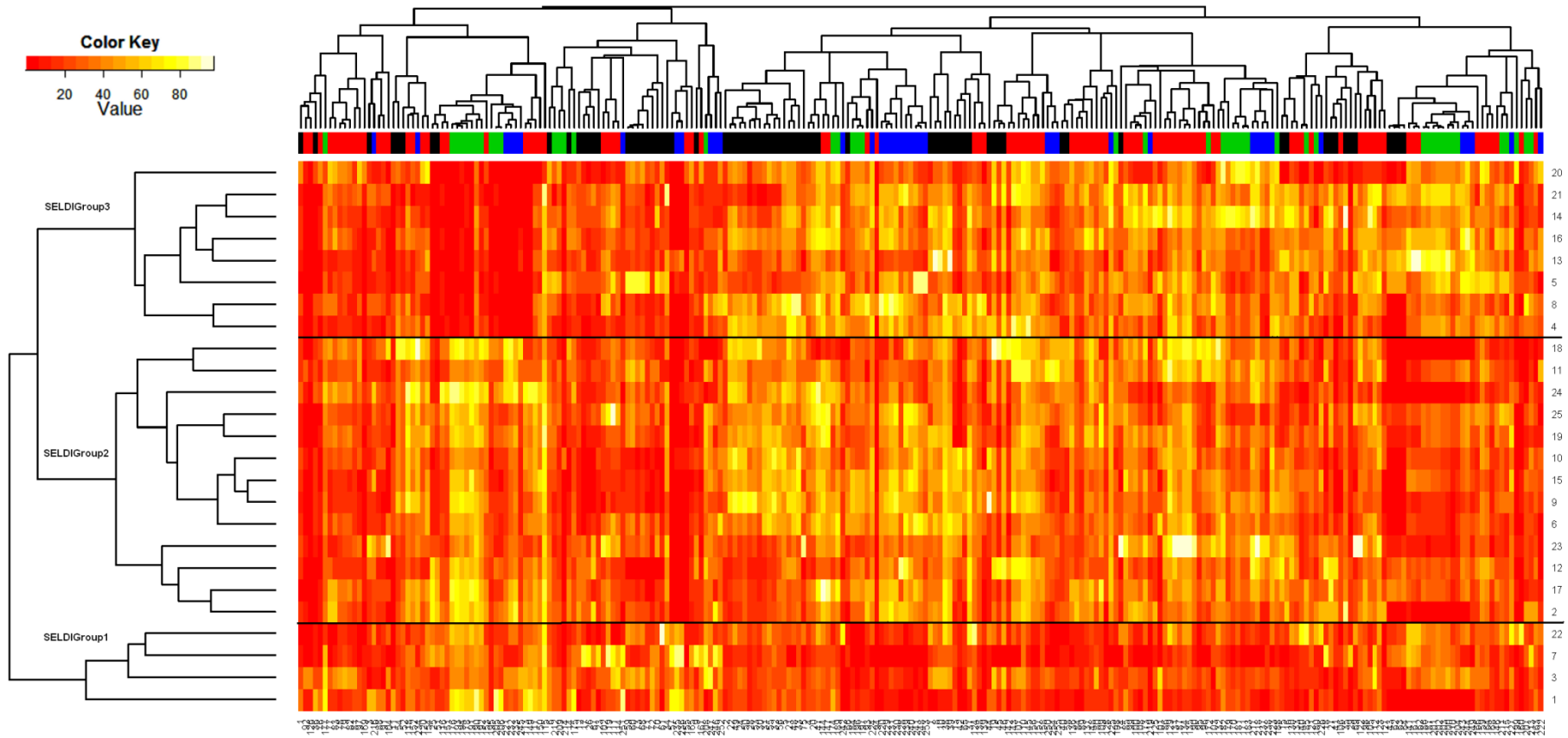


Peaky profilů 3 odlišných SELDI čipů  
 Vzorky zpracované 4 různými způsoby  
 Banks et al, Clinical Chemistry 2005

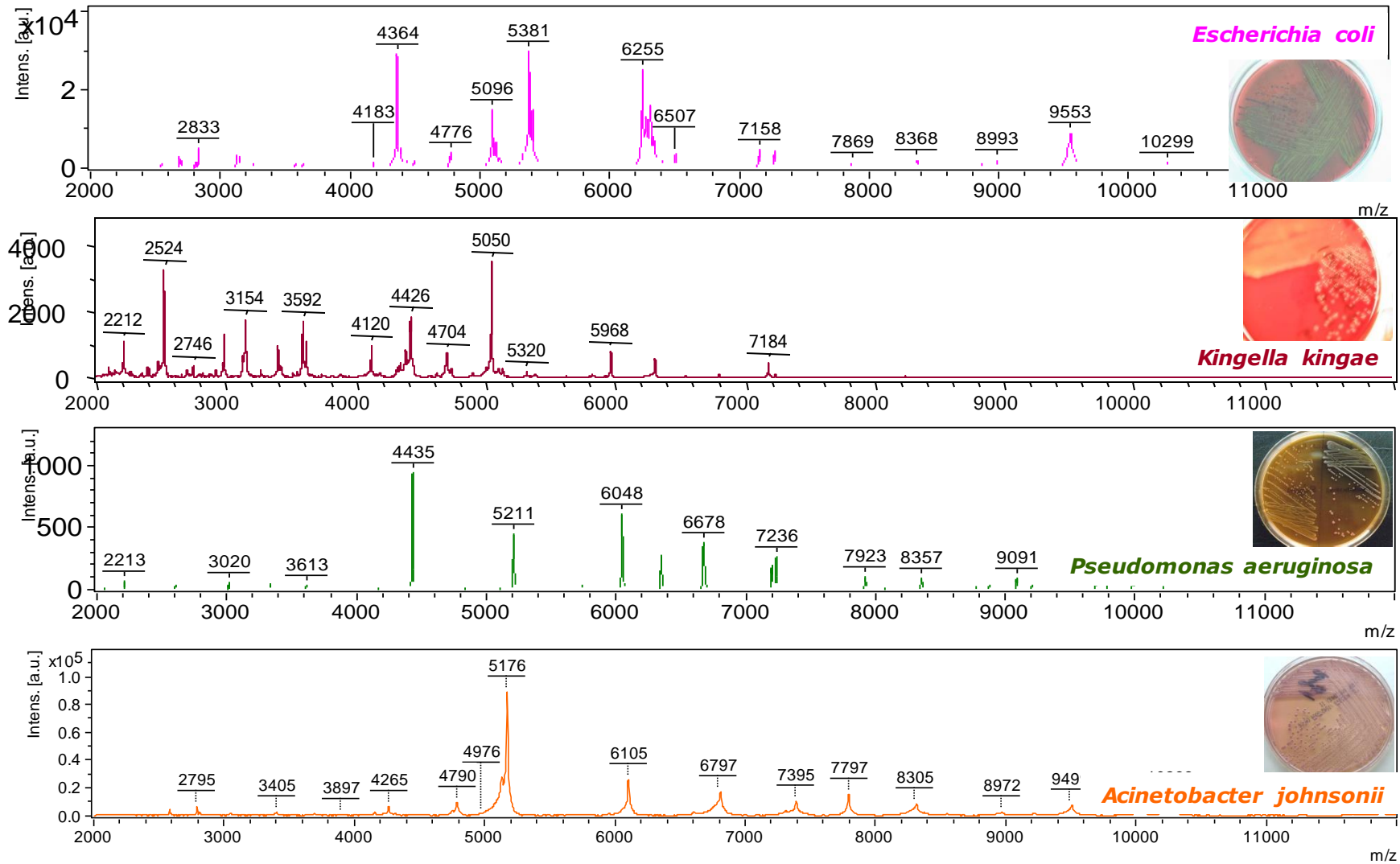
# Příklad

- Shlukování profilů stejných vzorků ze 4 typů SELDI sklíček:

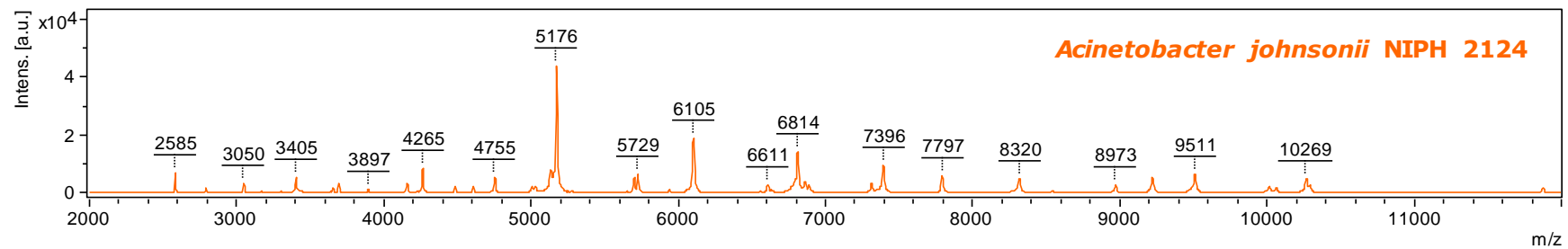
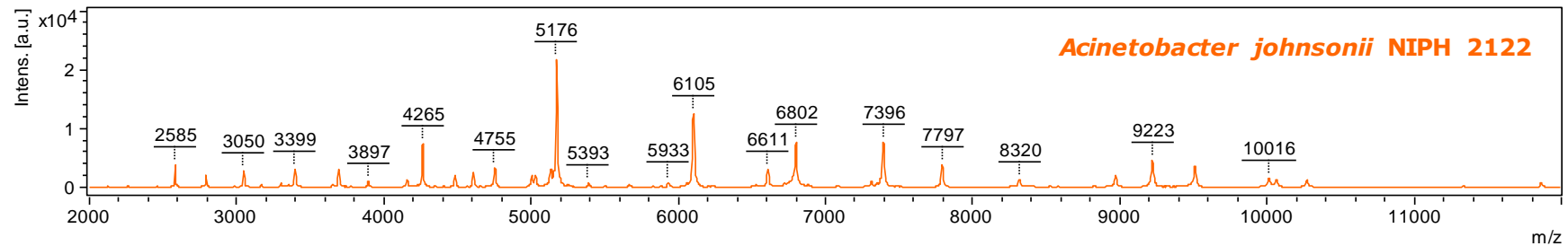
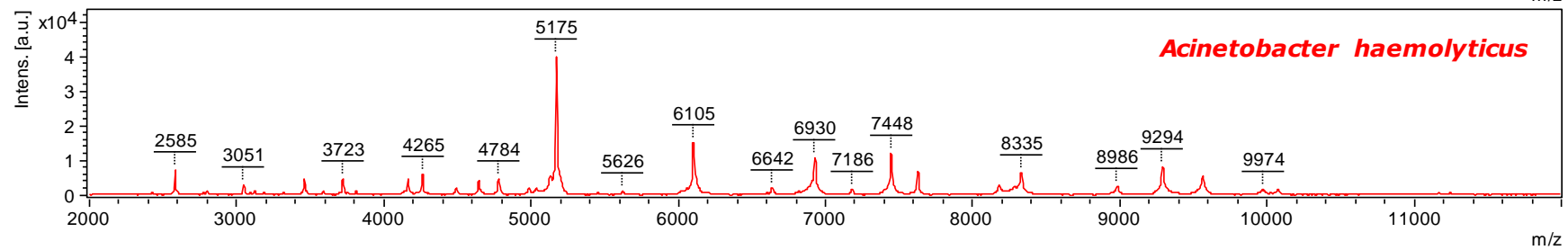
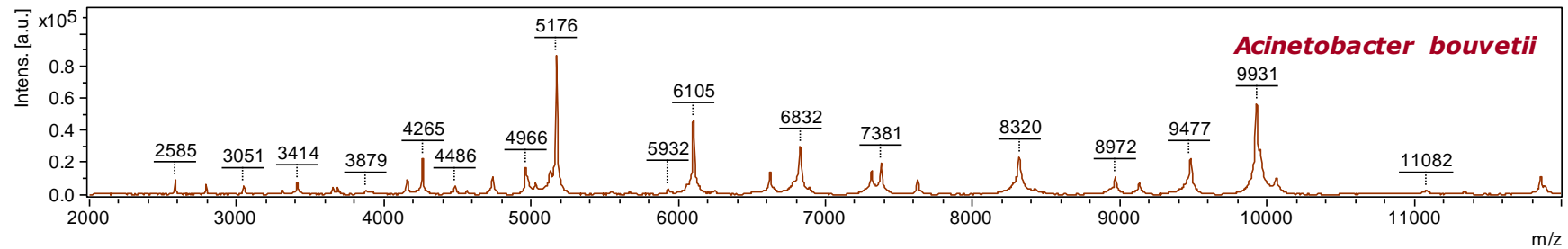
IMAC30, H50, NP20zas, NP20kys



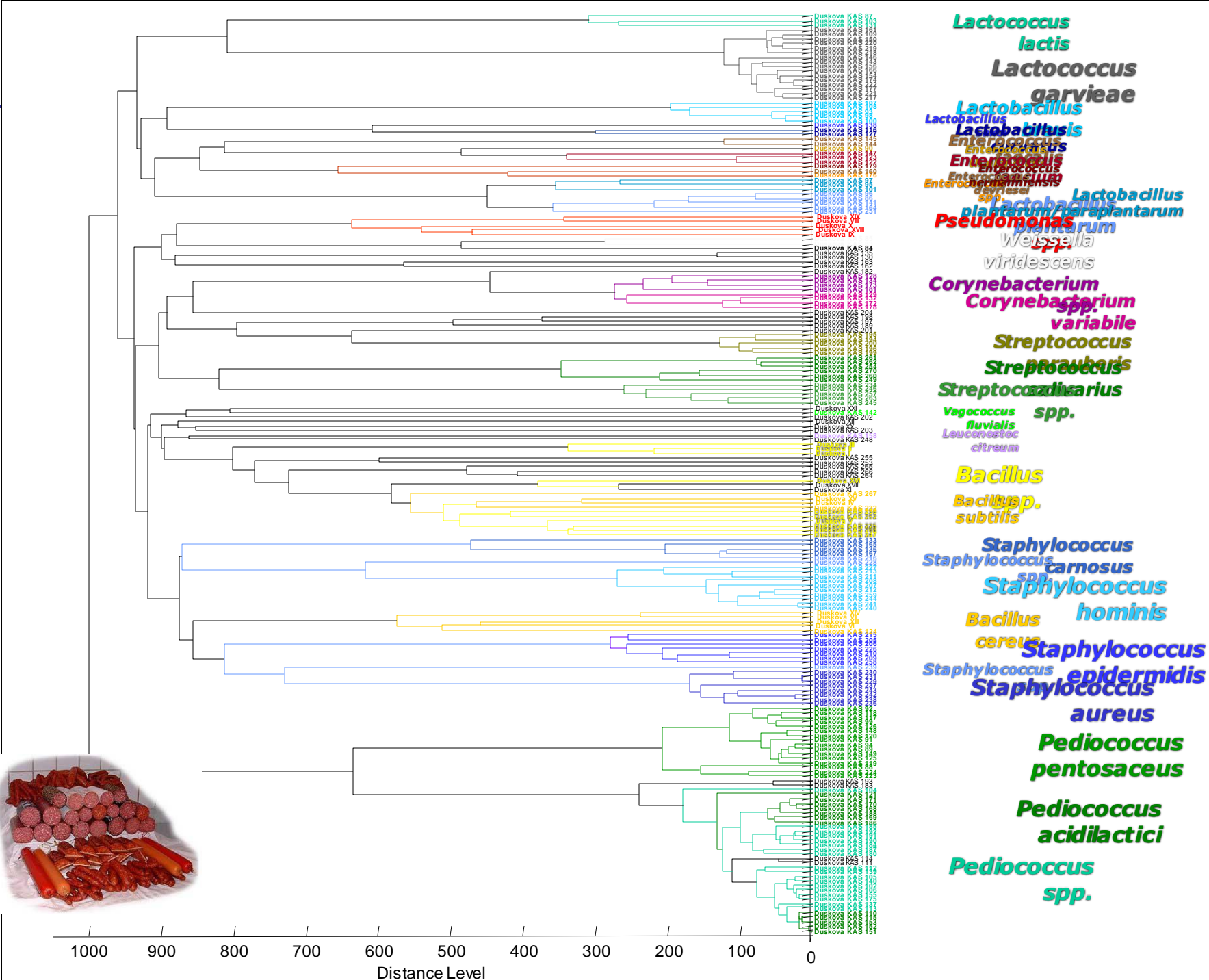
# Aplikace I – identifikace bakterií



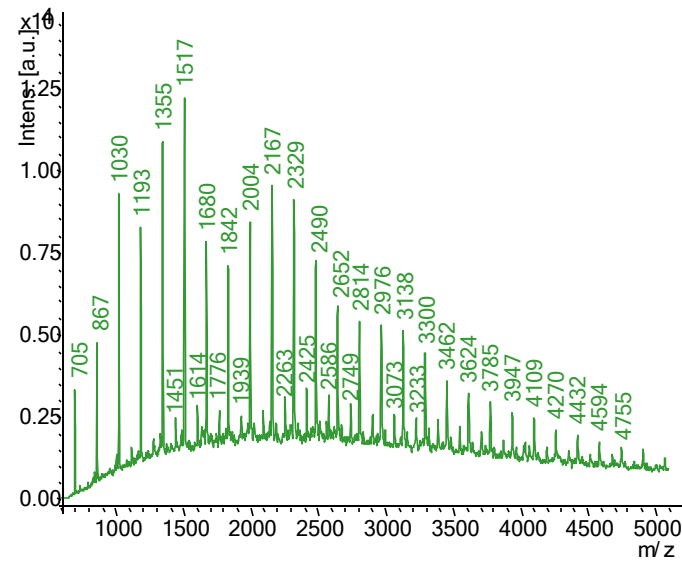
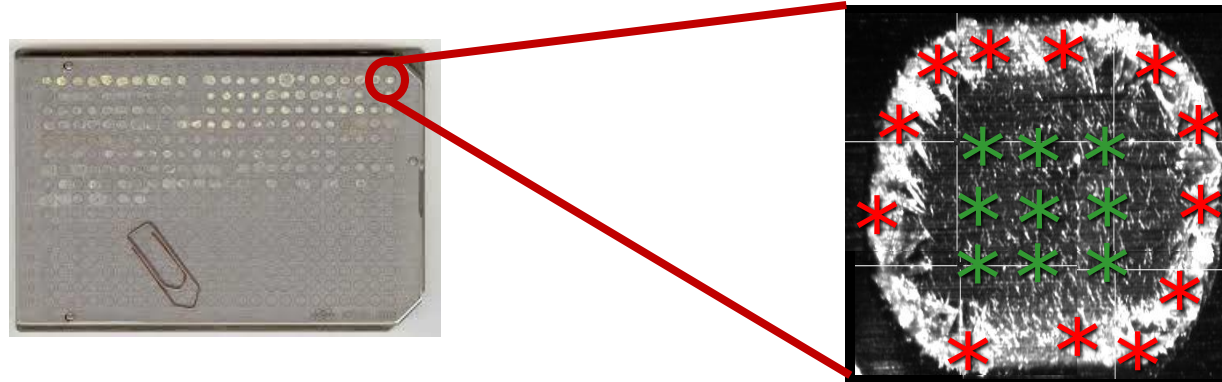
# Aplikace I – identifikace bakterií



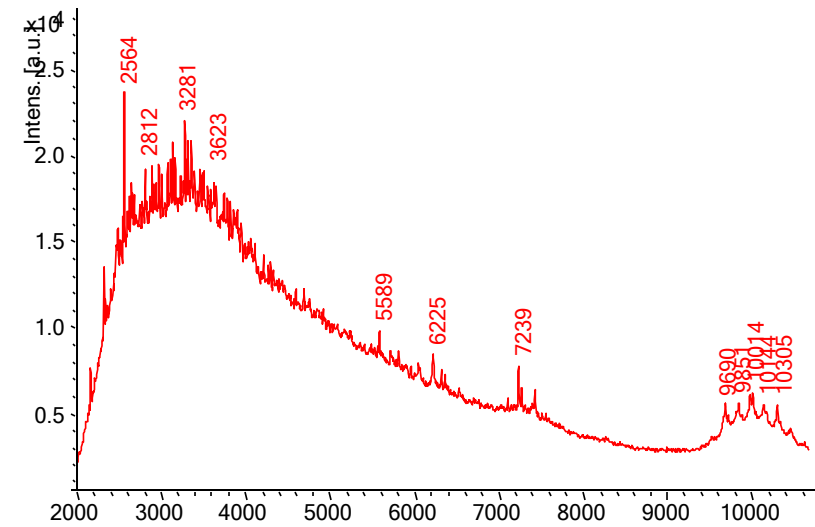




# Aplikace II

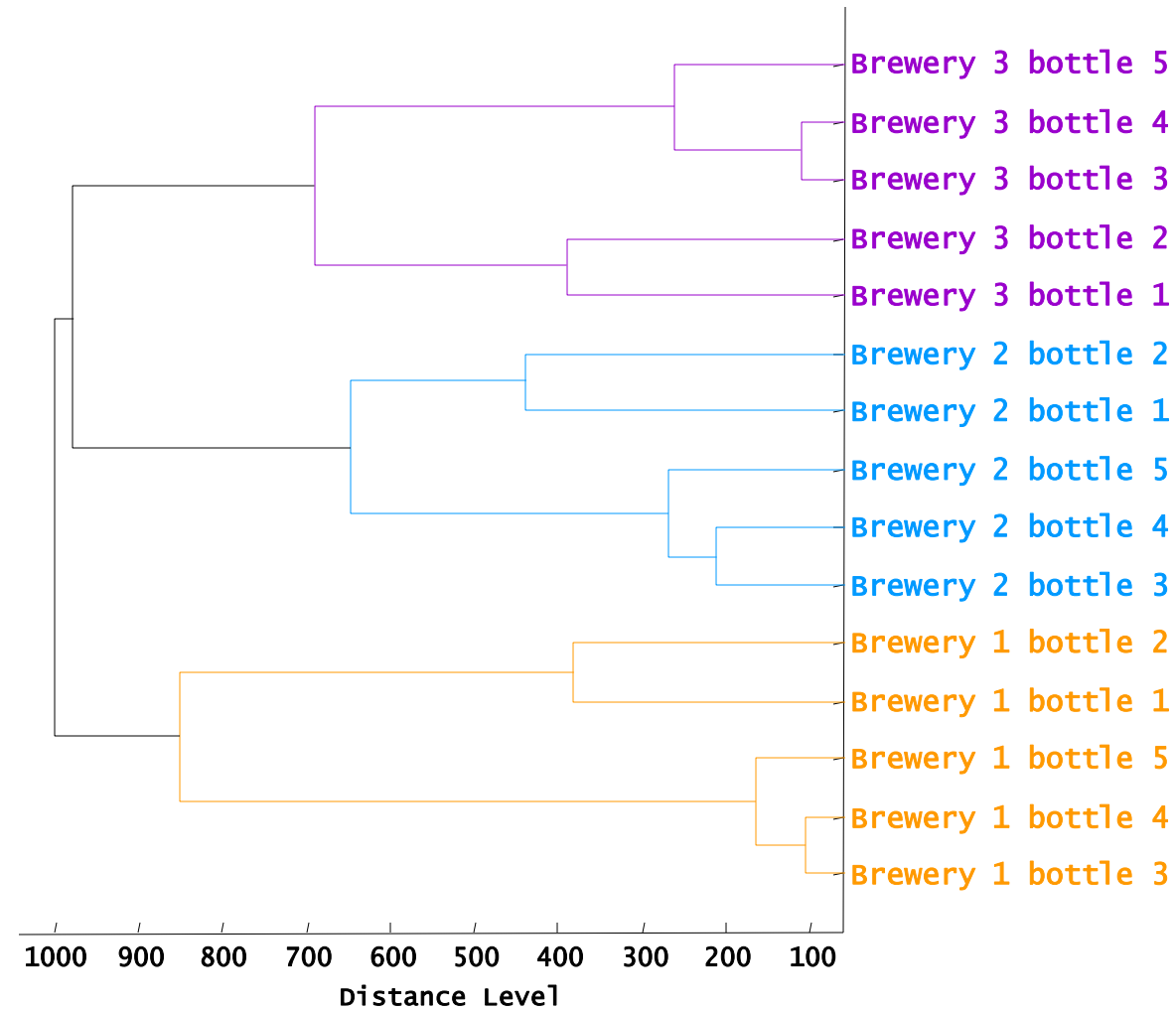


**MALDI-TOF MS fingerprint containing maltooligosaccharides**

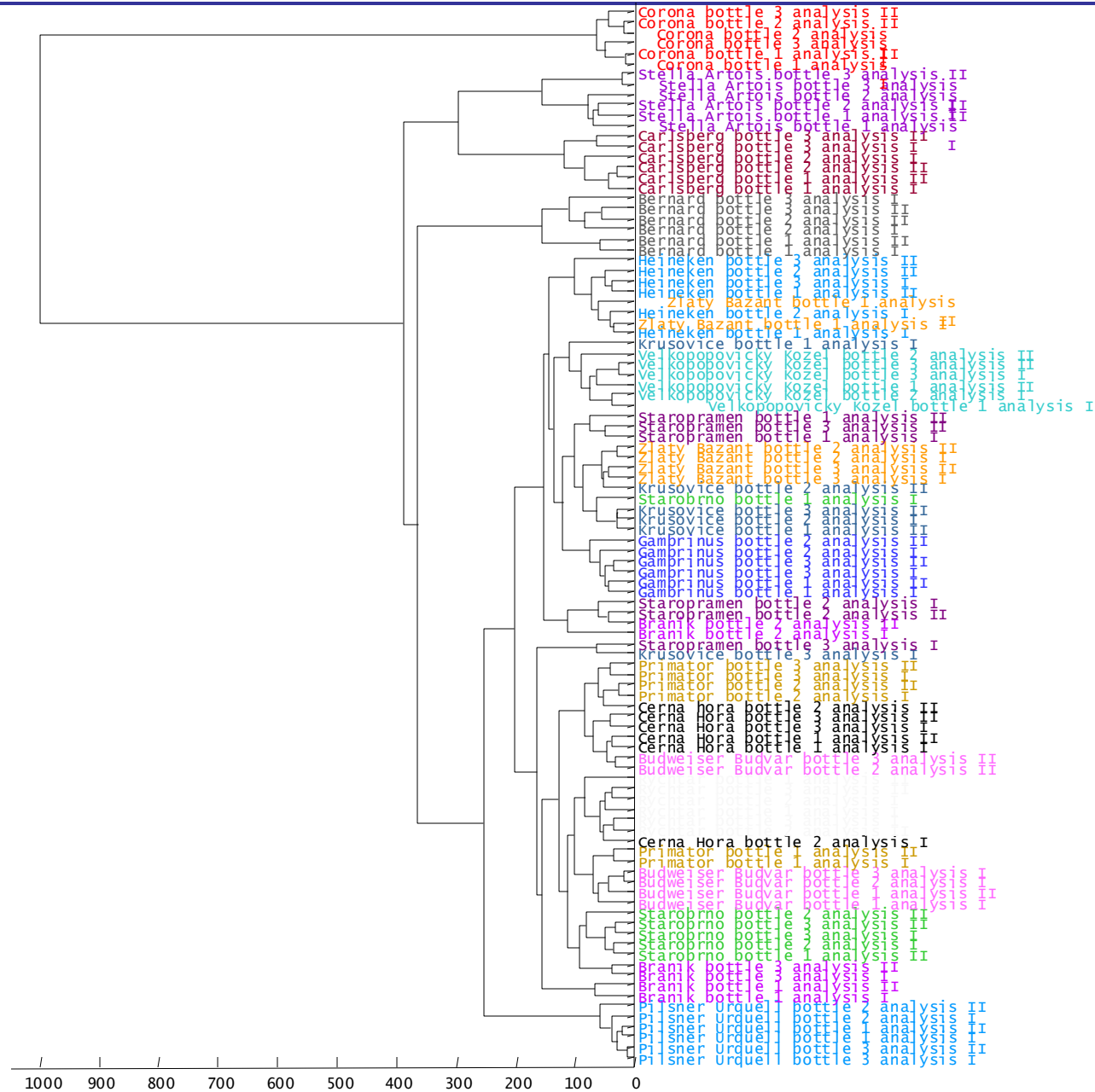


**MALDI-TOF MS fingerprint containing proteins**

# Aplikace II



# Aplikace II





**ALE MY CHCEME JEŠTĚ URČIT  
PROTEINY...**

# Zpracování dat

---

1. Úprava hrubých dat (MS/MS i MS), normalizace, identifikace píků
2. Identifikace proteinů s pomocí databáze

# Identifikace proteinů

**Princip:** Porovnáváme získaná spektra s cílovou databází pomocí databázových vyhledávačů (Sequest a Mascot), výsledkem je seznam shod spekter vzorku se spektry proteinů v db (peptide sequence matches - PSMs) => identifikace

**Problém:** Posoudit přesnost těchto identifikací však není triviální.

**Řešení:** Statistické přístupy a machine learning

# Tři základní kroky identifikace proteinů

## 1. Příprava dat

- Výběr „representativních“ signálů MS/MS
- Odstranění „méně kvalitních“ spekter MS/MS
  - Top N (z okna), dekonvoluce signálu a šumu
- Získáme tabulku m/z hodnot a intenzit

## 2. Příprava databáze

- *in silico* štěpení sekvencí z databáze
- Přiřazení jednoho a nebo více peptidů k jednomu spektru (s pomocí statistiky a machine-learning přístupů)

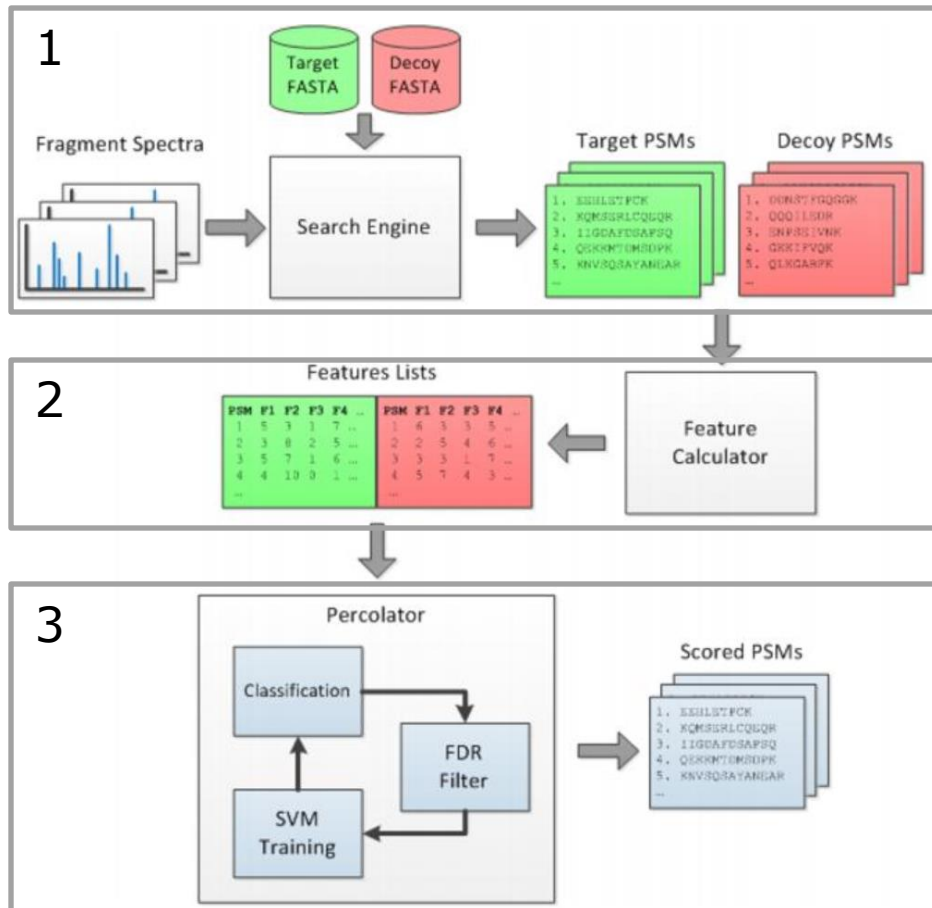
## 3. Výběr peptidových identifikací (kam patří, přiřazení k proteinu)



# Percolator (<http://percolator.ms>)



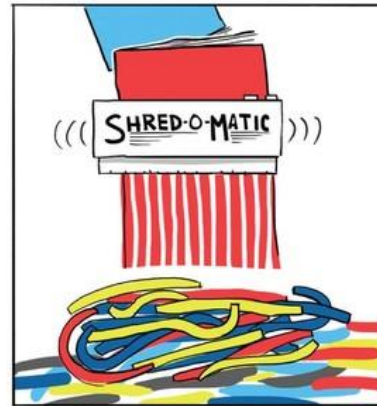
1. prohledání dat MS/MS
2. výpočet „vlastností“ peptidů
3. propočítání skóre peptidů



## Propočítání skóre peptidů

- Použití *support vector machines* (SVM)
- sady identifikací
  - falešně pozitivní – *decoy* databáze
  - pozitivní – původní databáze (skóre)
- přiřazení vah vlastnostem v SVM
  - např. skóre; chyba hmotnosti intenzita, modifikace, ...
- **víc identifikovaných peptidů**

# Rekonstrukce sady proteinů



Analogie puzzle, ALE:

- Tisíce kousků:
  - Stejné
  - Poškozené
  - Chybějící
  - Z jiných skládaček
- Pasují na stejná místa

# Metody rekonstrukce sady proteinů

- \* **Cíl:** zjistit, které peptidy patří kterým proteinům s větší pravděpodobností

- \* Dva základní přístupy:

## 1. N – peptidové pravidlo

- \* Proteiny, u kterých pozorujeme alespoň N peptidů

- \* Vysoká falešná pozitivita

- \* Používané na sekvenční homologické proteiny

## 2. Pravděpodobnostní přístupy

- \* *ProteinProphet, Nested mixtures, Fido*

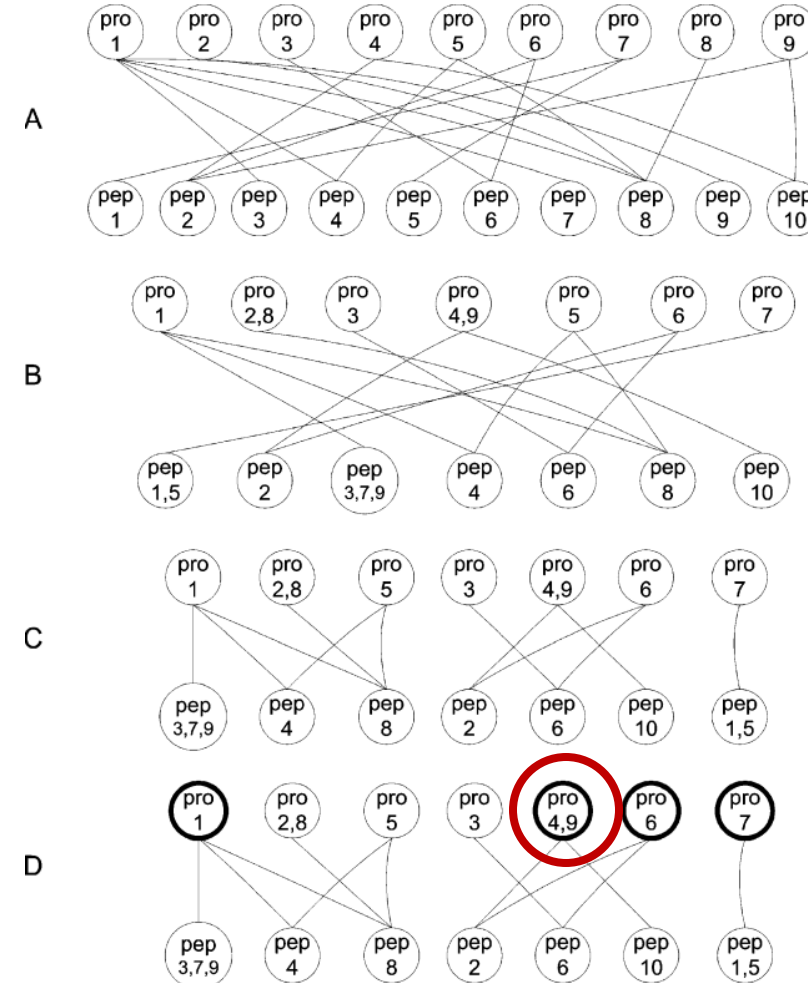
# Princip parsimonie a Occamové břitvy

A. Vytvoření bipartitního grafu:  
peptidy - možné proteiny

B. Sloučení proteinů a peptidů do skupin (např. pep 3,7,9; pro 4,9)

C. Rozdělení skupin

D. Výběr minimální sady proteinů



Důsledek: falešná negativita výsledků

# Co s identifikovanými proteiny?

---

- \* Závisí od původního experimentu

- \* Typicky doplnění anotace proteinů z databáze (GO, KEGG, TAIR) a použití metod analýzy genových sad (další přednáška)

# Identifikace proteinů

- NCBI Protein - <http://www.ncbi.nlm.nih.gov/protein>
  - \* jen pro proteinové sekvence odvozené translací nukleotidových sekvencí
- RefSeq - <http://www.ncbi.nlm.nih.gov/RefSeq/>
- UniProt– administrovaná databáze; kompozit SwissProt, TrEMBL a PIR-PSD–  
<http://www.uniprot.org>

# Cvičení a R balíky

- \* Provedeme cvičení MassSpectrometry.R
- \* Využívá bioconductor balík PROcess
  
- \* Další balíky
  - rTANDEM (an R/Bioconductor package for MS/MS protein identification)
  - dagLogo: An R/Bioconductor package for identifying and visualizing differential amino acid group usage in proteomics data

# Databáze dat

---



# Veřejně přístupné databáze

- Velké experimenty mají až stovky, a nebo tisíce vzorků, v každé se studují desetitisíce až stovky genů
- Pro publikaci výsledků je vyžadované vložit data ve standardizovaném formátu (MIAME – Minimal Information About a Microarray Experiment) do jedné z veřejně přístupných databází tak, aby kdokoliv byl schopný výsledky zreprodukovat
- Toto přináší velkou výhodu:
  - Můžeme data podrobit meta-analýze (simultánně porovnat data z různých experimentů)
  - Díky standardnímu formátu můžeme vyhledávat soubory s parametry, které potřebujeme
  - Data můžeme automaticky stahovat

# GEO na NCBI

The screenshot shows a Netscape browser window displaying the "GEO Database Design Brief" page. The browser's address bar shows the URL `http://www.ncbi.nlm.nih.gov/geo/info/scheme.cgi`. The page header includes the NCBI logo, the text "Gene Expression Omnibus", and the "geo" logo. A navigation bar contains links for "Entrez", "ProbeSet", "SAGEmap", "PubMed", "UniGene", and "LocusLink". Below this is a search bar with the text "Database Design Brief" and a "Query:" field with a "go" button.

The main content area features a sidebar on the left with various links and a main text area on the right. The sidebar includes sections for "Paper | FAQ | News", "Feedback **NEW**", "Retrieval tools" (with sub-links for GEO accession and attribute), "Deposit tools" (with sub-links for web, direct deposit, and new account), "Brief info" (with sub-links for current holdings, retrieving data, depositing data, and database design), and "Ad nauseam" (with sub-links for SOFT guide, examples, web deposit guide, entry fields, data tables, and SQL implementation).

The main text area contains the following text:

Please fill out our [feedback suggestionnaire](#) **NEW**

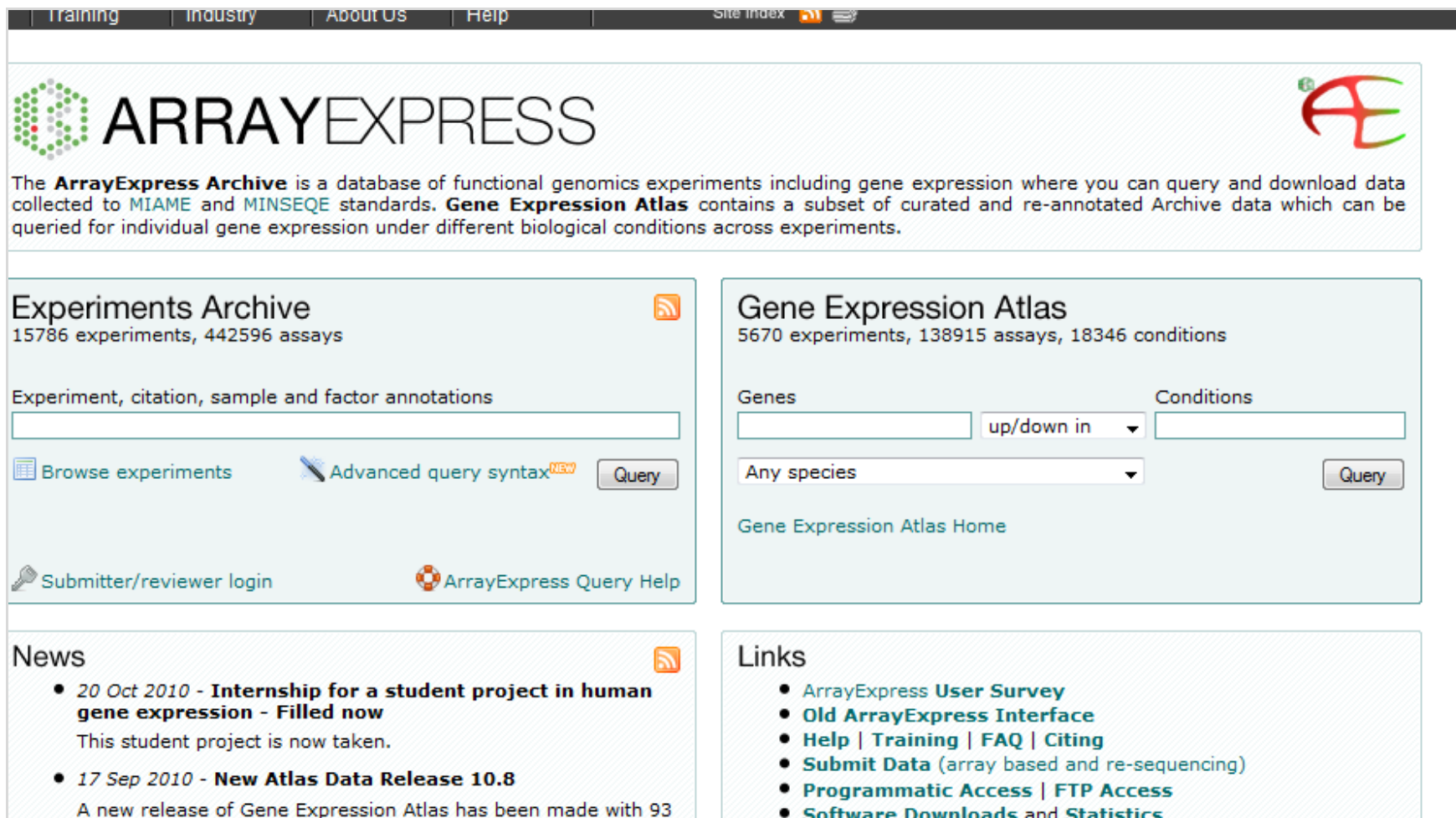
At the most basic level of organization of GEO there are four entities.

Below this text is a diagram illustrating the four entities: *Submitter*, *Platform*, *Series*, and *Sample*. The diagram shows a central box labeled *Submitter* at the top, connected by lines to three boxes below it: *Platform* on the left, *Series* on the right, and *Sample* at the bottom. The *Platform* and *Series* boxes are also connected to the *Sample* box, forming a diamond shape with a vertical line from the top to the bottom.

```
graph TD;
  Submitter[Submitter] --- Platform[Platform];
  Submitter --- Series[Series];
  Submitter --- Sample[Sample];
  Platform --- Sample;
  Series --- Sample;
```



# Array Express na EBI

<http://www.ebi.ac.uk/arrayexpress/>




The screenshot shows the ArrayExpress website interface. At the top, there is a navigation bar with links for Training, Industry, About Us, Help, and Site Index. The main header features the ArrayExpress logo and a description of the ArrayExpress Archive and Gene Expression Atlas. Below the header, there are two main sections: Experiments Archive and Gene Expression Atlas. The Experiments Archive section includes a search bar, a 'Query' button, and links for 'Browse experiments' and 'Advanced query syntax'. The Gene Expression Atlas section includes a search bar with a 'Query' button and a 'Gene Expression Atlas Home' link. At the bottom, there are sections for News and Links. The News section contains two items: '20 Oct 2010 - Internship for a student project in human gene expression - Filled now' and '17 Sep 2010 - New Atlas Data Release 10.8'. The Links section contains a list of links including 'ArrayExpress User Survey', 'Old ArrayExpress Interface', 'Help | Training | FAQ | Citing', 'Submit Data', 'Programmatic Access | FTP Access', and 'Software Downloads and Statistics'.



Training Industry About Us Help Site Index



 **ARRAYEXPRESS** 

The **ArrayExpress Archive** is a database of functional genomics experiments including gene expression where you can query and download data collected to **MIAME** and **MINSEQE** standards. **Gene Expression Atlas** contains a subset of curated and re-annotated Archive data which can be queried for individual gene expression under different biological conditions across experiments.

**Experiments Archive**   
15786 experiments, 442596 assays

Experiment, citation, sample and factor annotations

 Browse experiments  Advanced query syntax <sup>NEW</sup>


 Submitter/reviewer login  ArrayExpress Query Help

**Gene Expression Atlas**  
5670 experiments, 138915 assays, 18346 conditions

Genes  up/down in  Conditions

Any species

[Gene Expression Atlas Home](#)

**News** 

- **20 Oct 2010 - Internship for a student project in human gene expression - Filled now**  
This student project is now taken.
- **17 Sep 2010 - New Atlas Data Release 10.8**  
A new release of Gene Expression Atlas has been made with 93

**Links**

- [ArrayExpress User Survey](#)
- [Old ArrayExpress Interface](#)
- [Help | Training | FAQ | Citing](#)
- [Submit Data](#) (array based and re-sequencing)
- [Programmatic Access | FTP Access](#)
- [Software Downloads and Statistics](#)

# Další čtení

- E-learningová skripta analýzy dat IBA
- <http://portal.matematickabiologie.cz/index.php?pg=analiza-genomickych-a-proteomickych-dat--analiza-genomickych-a-proteomickych-dat>