

# COMPUTATIONAL PROTEOMICS AND METABOLOMICS

*Oliver Kohlbacher, Sven Nahnsen, Knut Reinert*

*4. Quantification I: General concepts, isobaric tags*



# Overview

- Quantification using mass spectrometry
  - Basic terms from analytical chemistry
  - Quantitative behavior of mass spectrometers
- Experimental quantification strategies
  - Absolute and relative quantification
  - Label-free vs. labeled techniques
  - Selected experimental techniques
  - Isobaric tags


# Analytical Chemistry

- “**Analytical chemistry** is the study of the separation, identification, and **quantification** of the chemical components of natural and artificial materials.”
- “**Quantification** [...] is the act of counting and measuring that maps human sense observations and experiences into members of some set of numbers.”
- **Quantitative Mass Spectrometry** :=  
use of a mass spectrometer to turn amounts of *analytes* into numbers

[http://en.wikipedia.org/wiki/Analytical\\_chemistry](http://en.wikipedia.org/wiki/Analytical_chemistry) [accessed 12.11.2011, 10:40 CET]

<http://en.wikipedia.org/wiki/Quantification> [accessed 12.11.2011, 10:45 CET]

# Some Terms

- **Analyte** – the stuff we want to analyze (proteins, peptides, metabolites)
  - **Matrix** – the components of the sample that are not analytes
  - The matrix can significantly impact the way the whole analysis is performed
  - **Example**
    - Proteomics analysis from urine
    - Urine contains
      - Proteins and peptides – the **analytes**
      - Water
      - Metabolites
      - Urea
- 

# Matrix Effects in LC-MS

- Components of the matrix are being separated just like the analytes
- Parts of the matrix can be ionized as well and then also show up as signals in the MS
- *A priori* it is unknown, which part of the signal stems from matrix or analytes
- Matrix can interfere with the analysis by
  - Competing with analytes for ionization -> reduce the number of analyte molecules ionized
  - Adsorb, precipitate or even react with the analyte

# Quantifying Analytes

- Analytes have to be in solution for proteomics and metabolomics
- We thus deal with concentrations: amounts per volume of sample  $V$

- Molar concentration

$$c_i = n_i / V \quad [\text{SI unit: mol/m}^3]$$

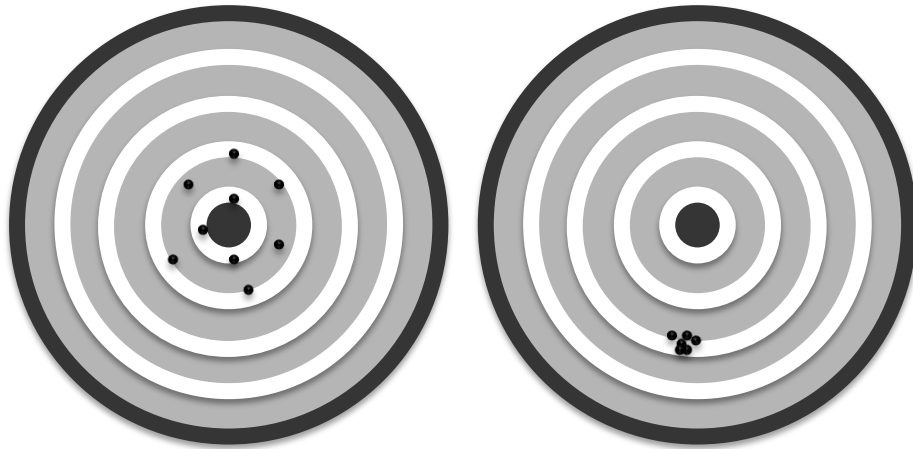
- Mass concentration

$$\rho_i = m_i / V \quad [\text{SI unit: kg/m}^3]$$

- Translating molar concentrations into mass concentrations can be done via the molecular weight  $M_i$  of the analyte

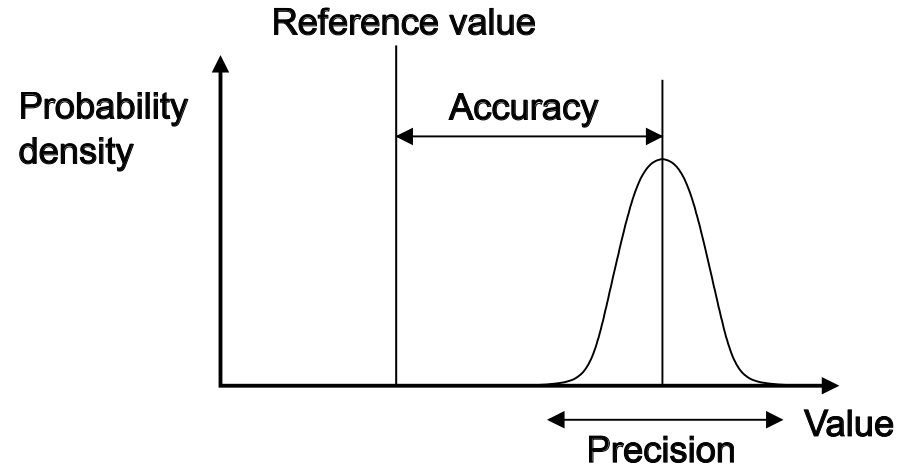
$$\rho_i = c_i M_i$$

# Precision and Accuracy



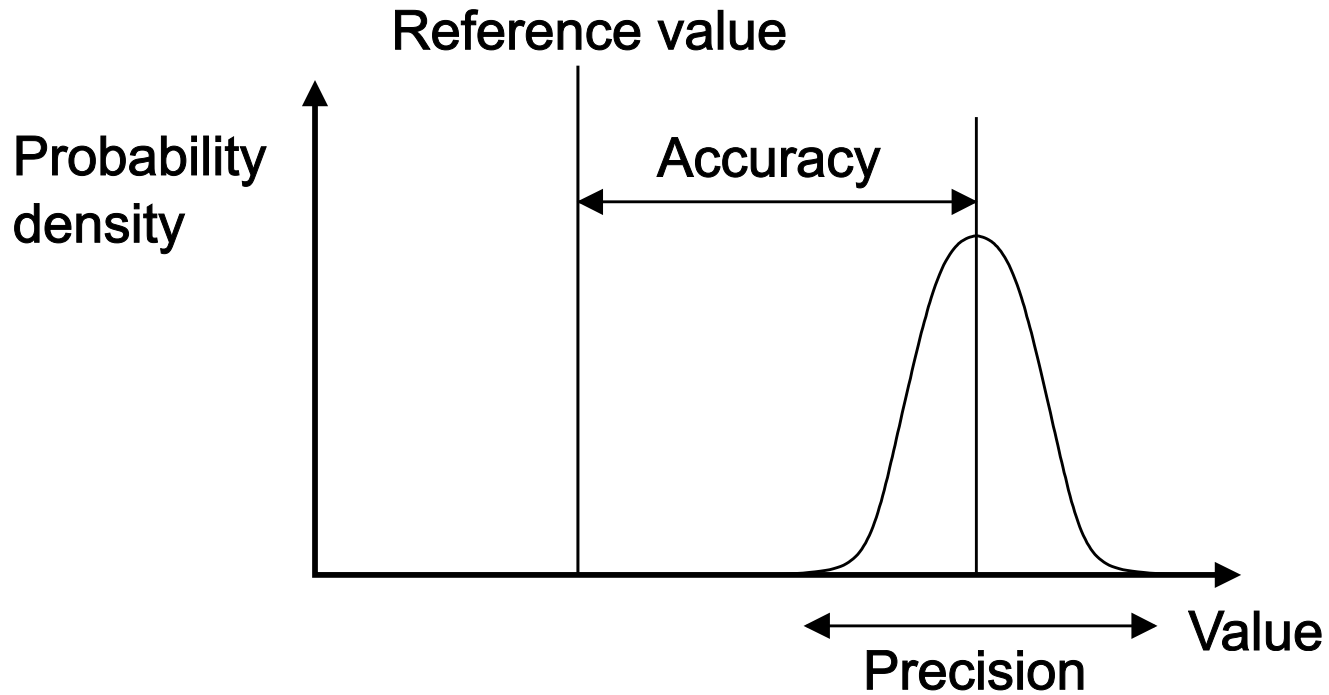
good accuracy,  
poor precision

good precision,  
poor accuracy



- **Accuracy:** closeness to the true value (mostly influenced by systematic error) – repetition of the experiment will not improve the result
- **Precision:** repeatability of the measurement (mostly influenced by random error) – repetition of the experiment will yield a value closer to the true value
- An ideal experiment combines high accuracy with high precision

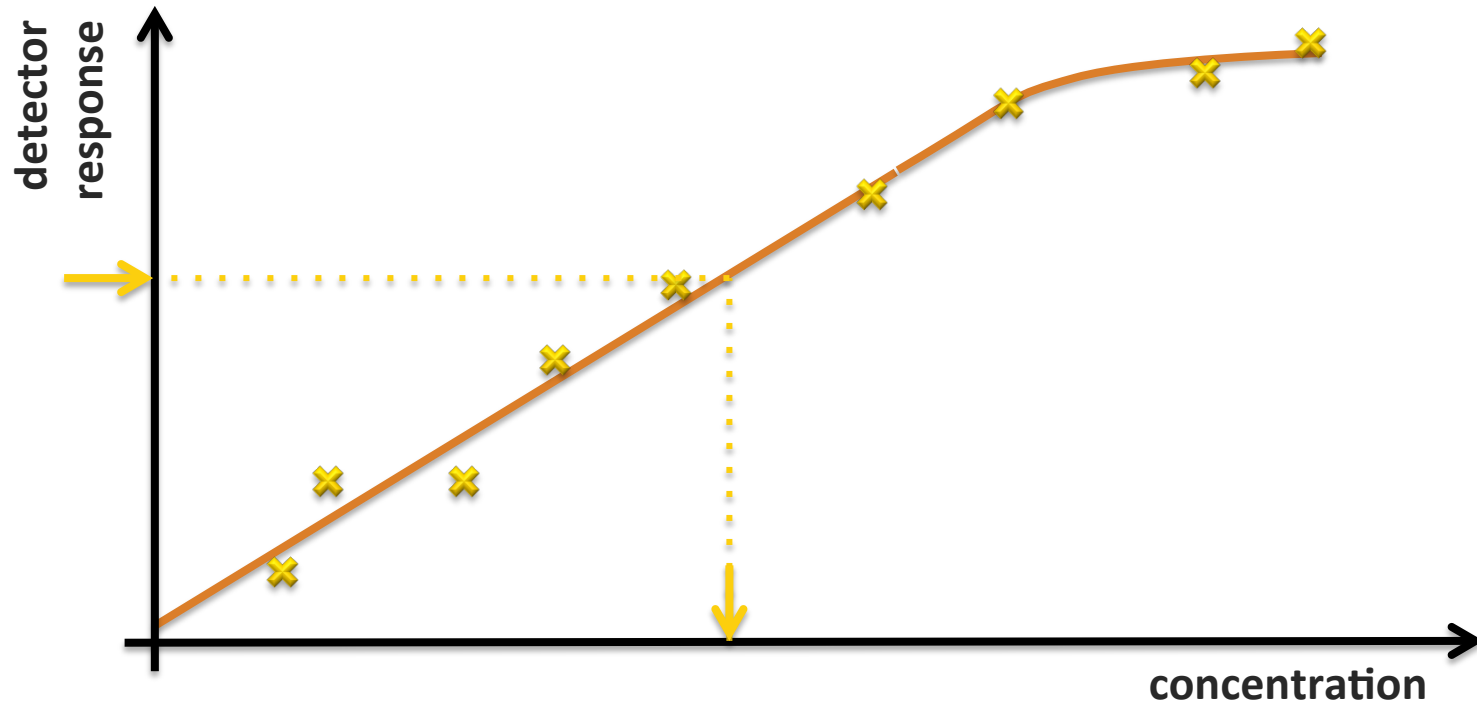
# Measurement Errors



- Each measurement is associated with an error
- There are two basic types of error:
  - **Random error**: defines the variance of repeated measurements (e.g., due to high noise level) – this is always present in every measurement
  - **Systematic error** (bias): shifts the mean of repeated experiments (e.g., due to an incorrect calibration)

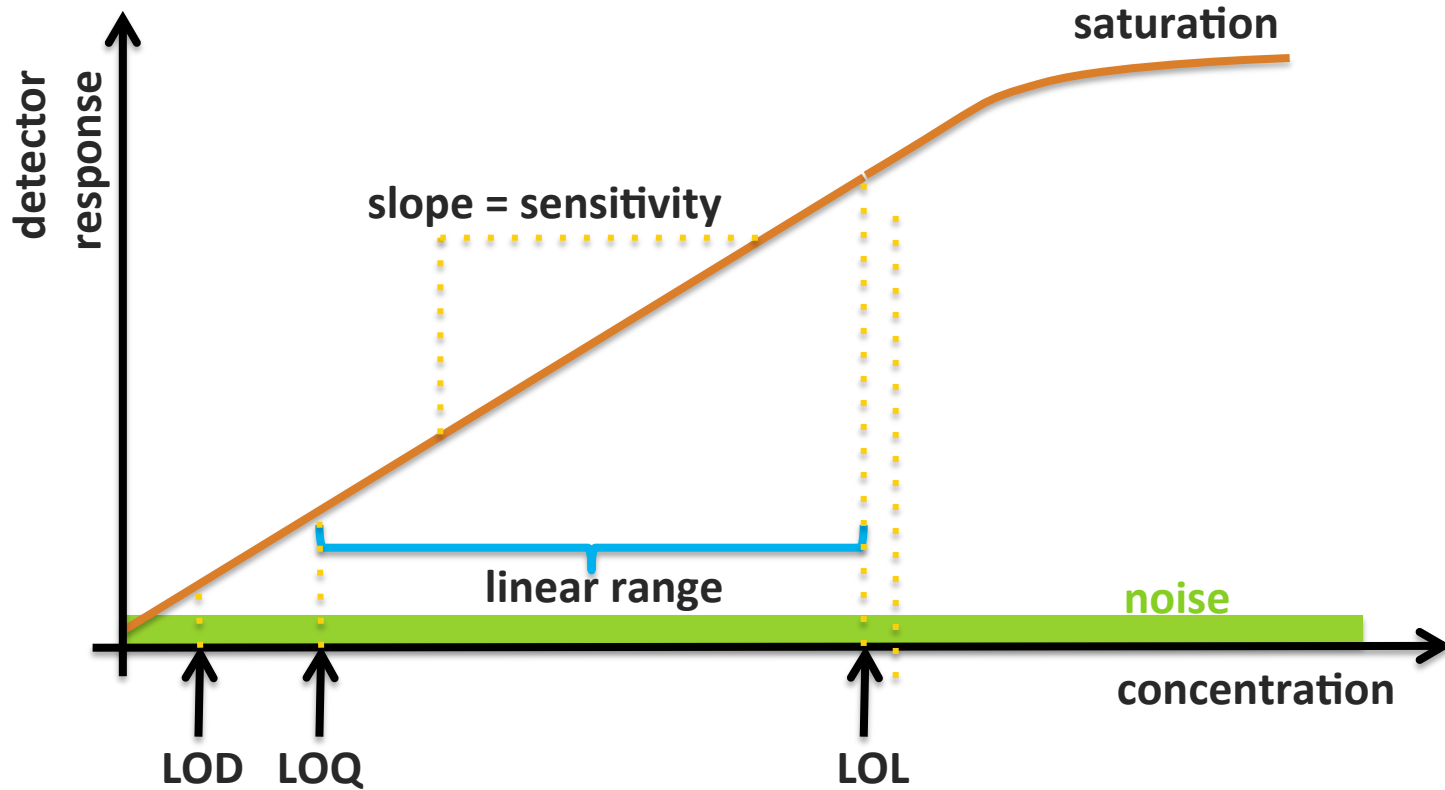


# Calibration Curve



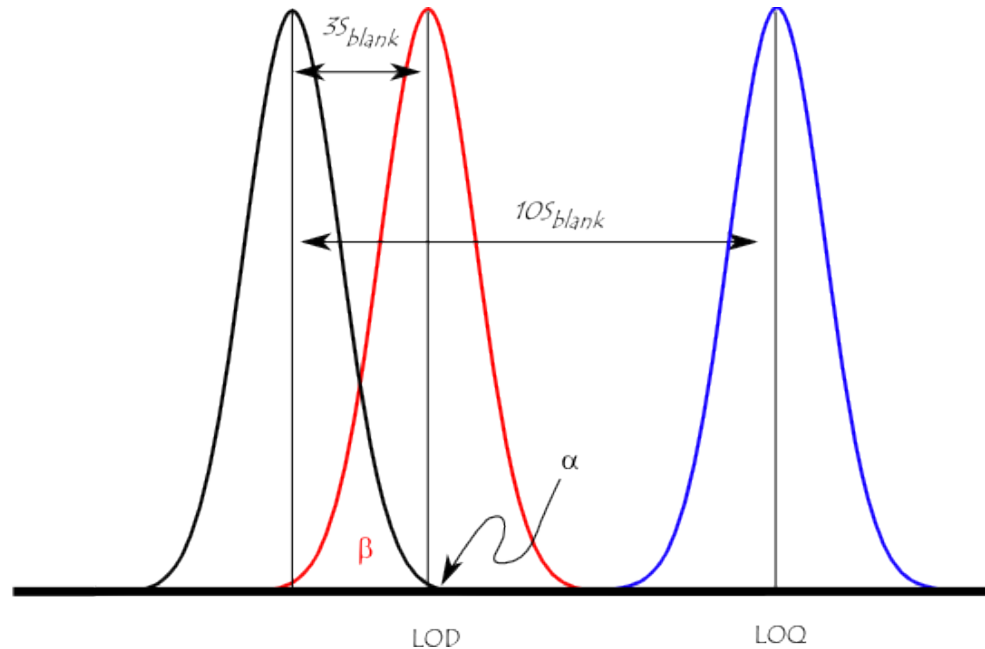
- Measurement of the detector response for various (known) concentrations allows the construction of a calibration curve
- Most detector responses are chosen in a way that the response changes linearly with the concentration
- Once the calibration curve has been measured, it allows the determination of the concentration of an unknown sample

# Response



- **LOD**: level of detection – at what concentration can we decide that the analyte is present
- **LOQ**: level of quantification – at what concentration can we accurately quantify it
- **LOL**: limit of linearity – saturation effects start here
- **Linear range (dynamic range)**: the concentration range where we get a response that is linear in the concentration

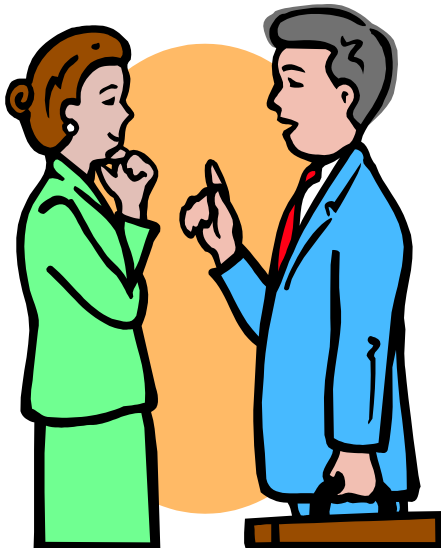
# Detection Limit



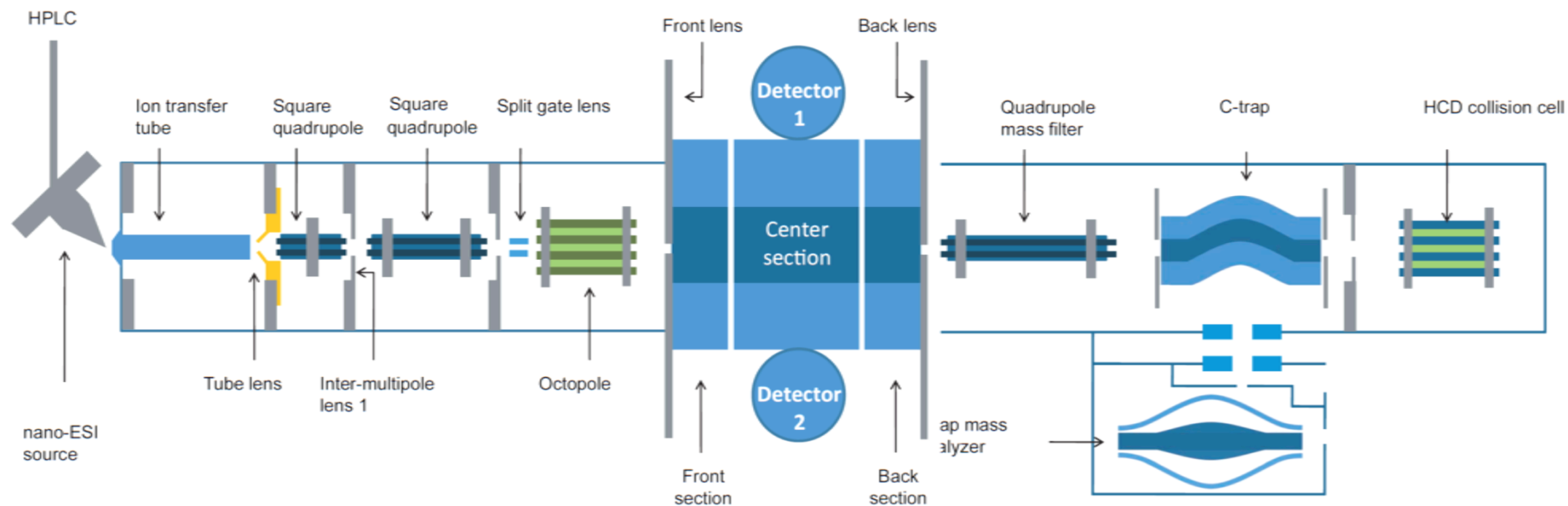
- **Limit of detection** (detection limit) -- LOD: the lowest analyte concentration that can be distinguished from the absence of the analyte (blank) within a stated confidence limit (generally 99% confidence)
- **Limit of quantification** – LOQ: the concentration at which we can distinguish two values with reasonable confidence
- Both depend on the noise level, the matrix, the instrument, the sensitivity for a specific analyte, etc.

# LOD/LOQ

*“Suppose you are at an airport with lots of noise from jets taking off. If the person next to you speaks softly, you will probably not hear them. Their voice is less than the LOD. If they speak a bit louder, you may hear them but it is not possible to be certain of what they are saying and there is still a good chance you may not hear them. Their voice is  $>LOD$  but  $<LOQ$ . If they speak even louder, then you can understand them and take action on what they are saying and there is little chance you will not hear them. Their voice is then  $>LOD$  and  $>LOQ$ . Likewise, their voice may stay at the same loudness, but the noise from jets may be reduced allowing their voice to become  $>LOD$ . Detection limits are dependent on both the signal intensity (voice) and the noise (jet noise).”*



# Quantitative Mass Spectrometry



- Ionization: number of ionized analyte molecules proportional to the total amount present
- MS detector: proportional to the number of ions (the ion current)
- Caveats:
  - Saturation: there is an upper limit to the response
  - Noise: does the signal really come from the analyte?

# Quantitative LC-MS

- **Fixed volume** of the sample is injected
- Total amount of analyte eluting from the column is the same amount as the amount injected (normally, **nothing gets 'lost'** on the column)
- Analyte spreads out, elutes over a certain timespan from the column: maximum **concentrations at the end of the column depend on retention time** (peak broadening)
- Only a fraction of the analyte really enters the MS (skimmer!)
- **Ionization efficiency differs** between analytes

# Quantitative LC-MS

- MS signal intensity for peptide  $i$  at time  $t$  is proportional to **concentration  $c_i(t)$**  eluting off the column.

$$I_i(t) = f_i \cdot c_i(t)$$

- The **area under the (chromatographic) peak is proportional to the total amount  $c_i^{\text{tot}}$**  of analyte eluting and thus to the amount in the sample. Hence we want to integrate over time.

$$\int_t I_i(t) = f_i \cdot \int_t c_i(t)$$

# Quantitative LC-MS

- Elution profiles are (roughly) Gaussians. Hence we can model the the elution as a product of the total concentration spread by a retention time model

$$c_i(t) = g(rt_i, \sigma_i, t) c_i^{tot}$$

- Strategy
  - Integrate over the MS signal (intensity  $I_i(t)$ ) caused by the analyte  $i$  over the total elution time of an analyte (centered around  $rt_i$ , peak width defined by standard deviation of the Gaussian)
  - Response factor  $f_i$  is unknown

$$\int_t I_i(t) = f_i \cdot c_i^{tot} \cdot \int_t g(rt_i, \sigma_i, t)$$

$$\int_t I_i(t) = f_i \cdot c_i^{tot} \cdot 1$$



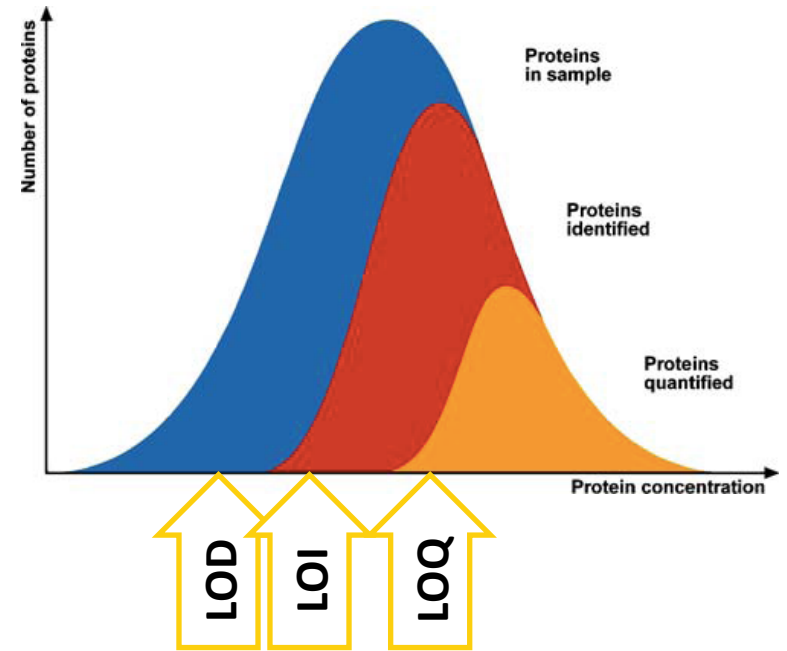
# Detection, Identification, Quantification

- Proteomics

- More peptides/proteins are usually identified than quantified
- Identification: MS/MS, quantification usually by MS -> independent processes
- Many things can be seen (detected) but cannot be identified or quantified

- Metabolomics

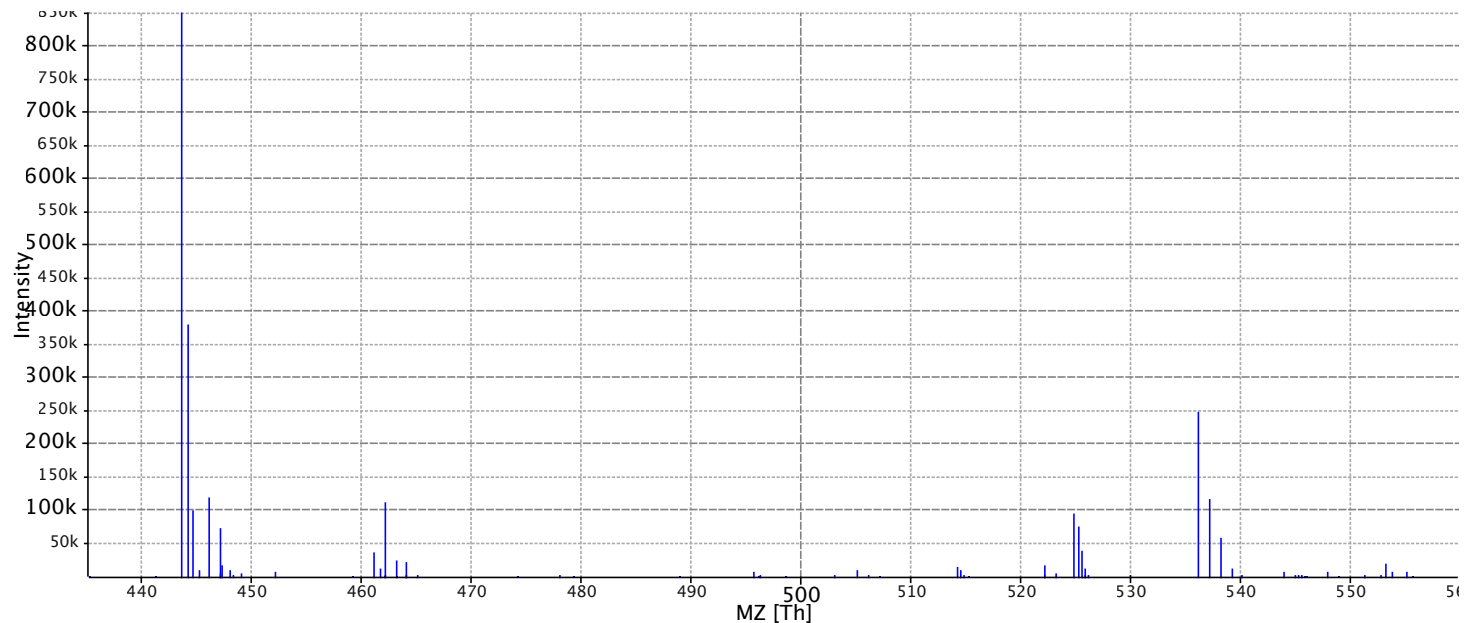
- Identification here is particularly difficult
- We can identify only a fraction of what we can quantify



**LOI: “Level of identification”**

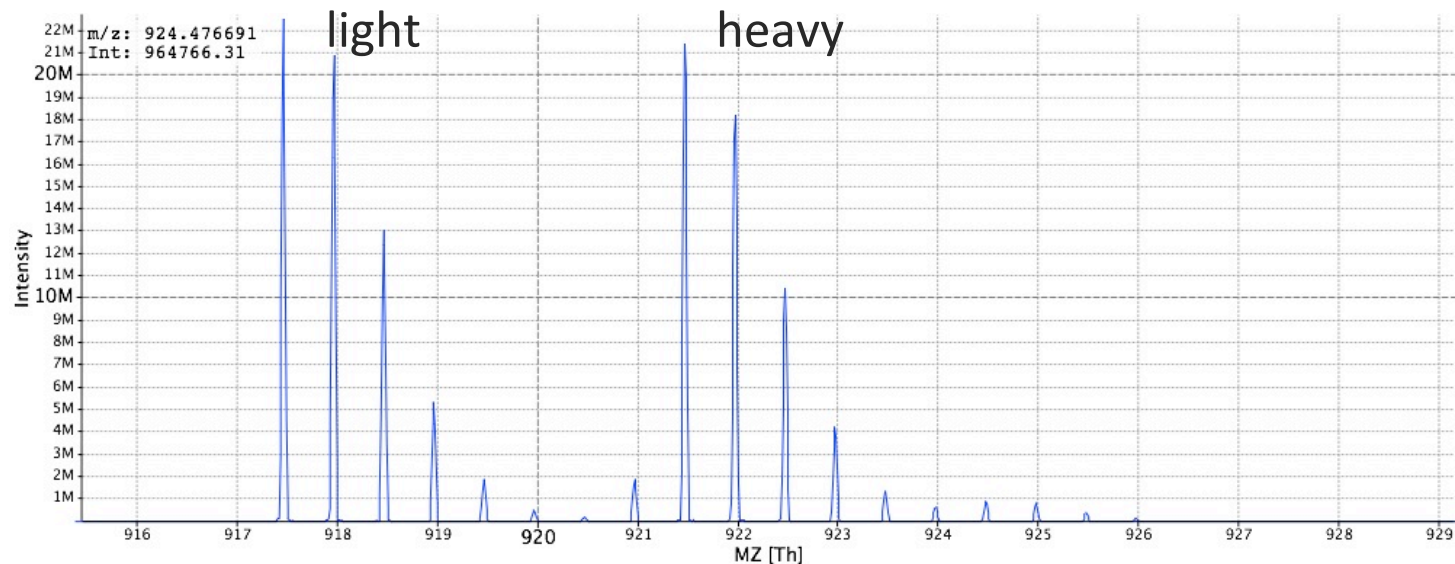
# Quantitative Data – MS Spectra

- Different ionized species in the same MS spectrum result in different peaks
- Example
  - Each peptide leads to a distinct set of peaks (isotope patterns!)
  - Intensity of each peak is proportional to the concentration at the time of elution



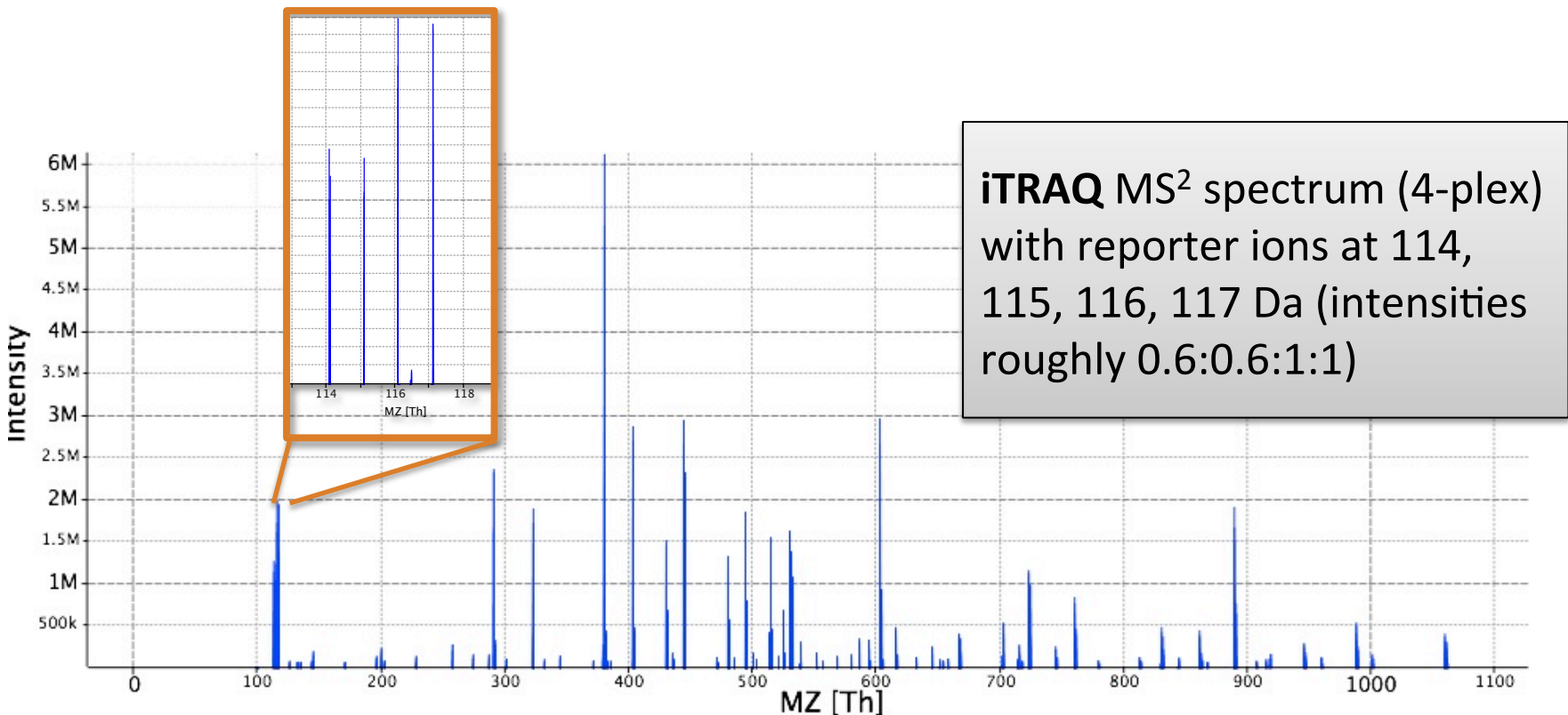
# Quantitative Data – MS Spectra

- **Direct comparison of intensities of different analytes** in the same spectrum **is not possible** because they have different response factors!
- Exception: peptides/metabolites that differ only by a stable isotope label will have identical response factors – their intensities can be compared **within the same spectrum**! This is the basis for isotopic labels.



# Quantitative Data – MS<sup>2</sup> Spectra

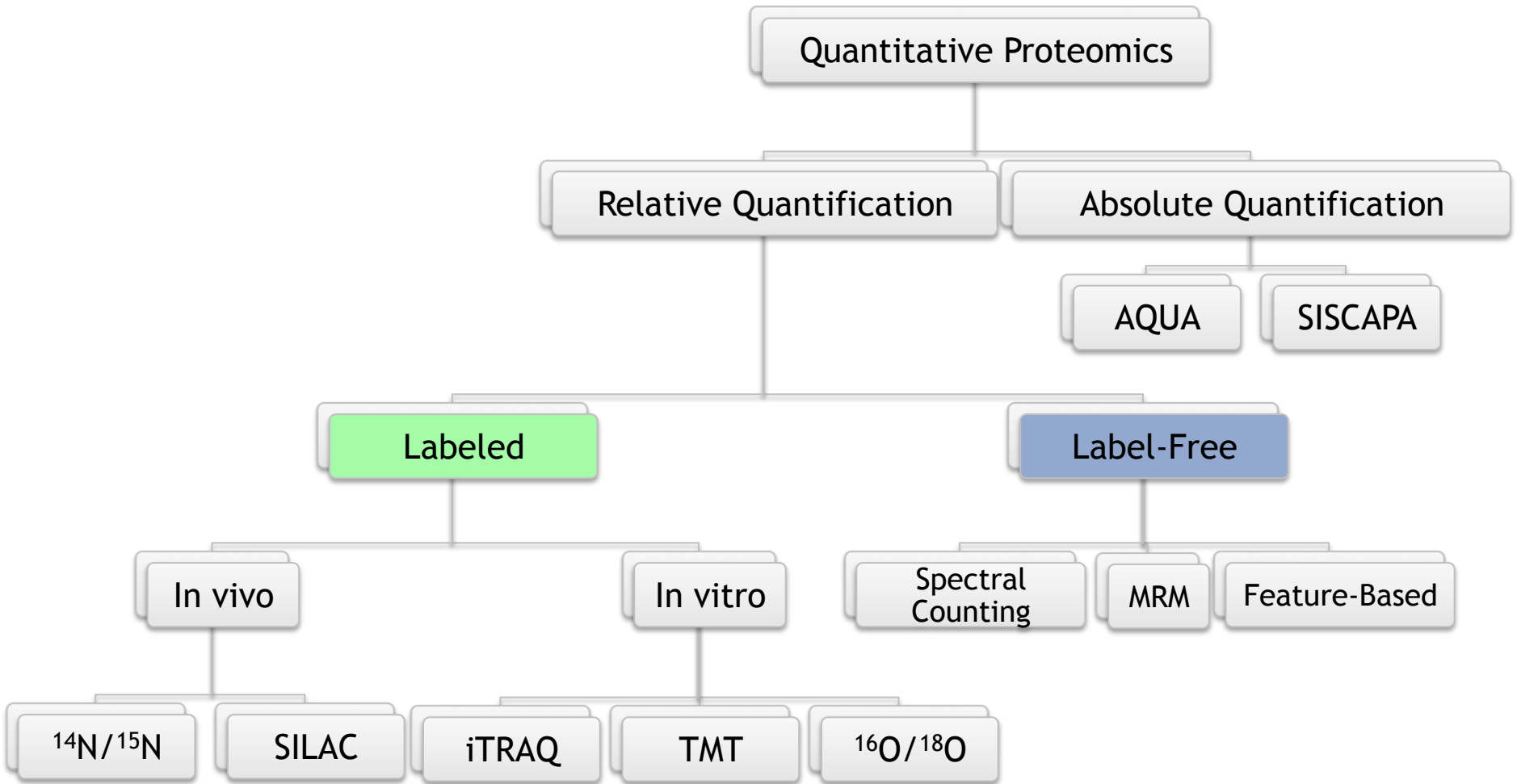
- Fragment spectra can be used for quantification as well
  - Under identical fragmentation conditions, the fragment ion intensity is proportional to the parent ion concentration/intensity
  - Key methods: MRM, iTRAQ



# Chromatograms

- Except for quantification techniques where a direct comparison is made within the same spectrum (iTRAQ, SILAC), elution profiles have to be considered
- Accurate quantification requires accurate **integration over the retention time profile**
- Since the peak area remains the same, this means the quantification will be independent of changes in the peak shape and width
- Elution profiles are often assumed to be Gaussian, but in reality they can deviate significantly (tailing/heading leads to asymmetric peak shapes – in the model of theoretical plates, this corresponds to incomplete equilibration)
- For details, see Learning Unit 2A

# Quantification Strategies



# Labeling Techniques

- Many labeling techniques exploit stable isotope labeling
  - Different isotopes of the same element behave chemically basically identically (often used:  $^1/2\text{H}$ ,  $^{12}/^{13}\text{C}$ ,  $^{14}/^{15}\text{N}$ ,  $^{16}/^{18}\text{O}$ )
  - Their masses differ, however, so the MS can distinguish them
- Introducing a label in one sample and a different (or no label) in another, mixing allows a relative quantification between two (or more) samples
- **Advantages**
  - Both samples are treated identically, systematic errors affect them in the same way
  - Can be easily annotated manually (e.g., look for pairs of peaks)
- **Disadvantages**
  - Labels can be expensive, difficult, unreliable to introduce
  - Labeling *in vivo* is not always possible, not all techniques support *in vitro* labeling

# Labeling Techniques

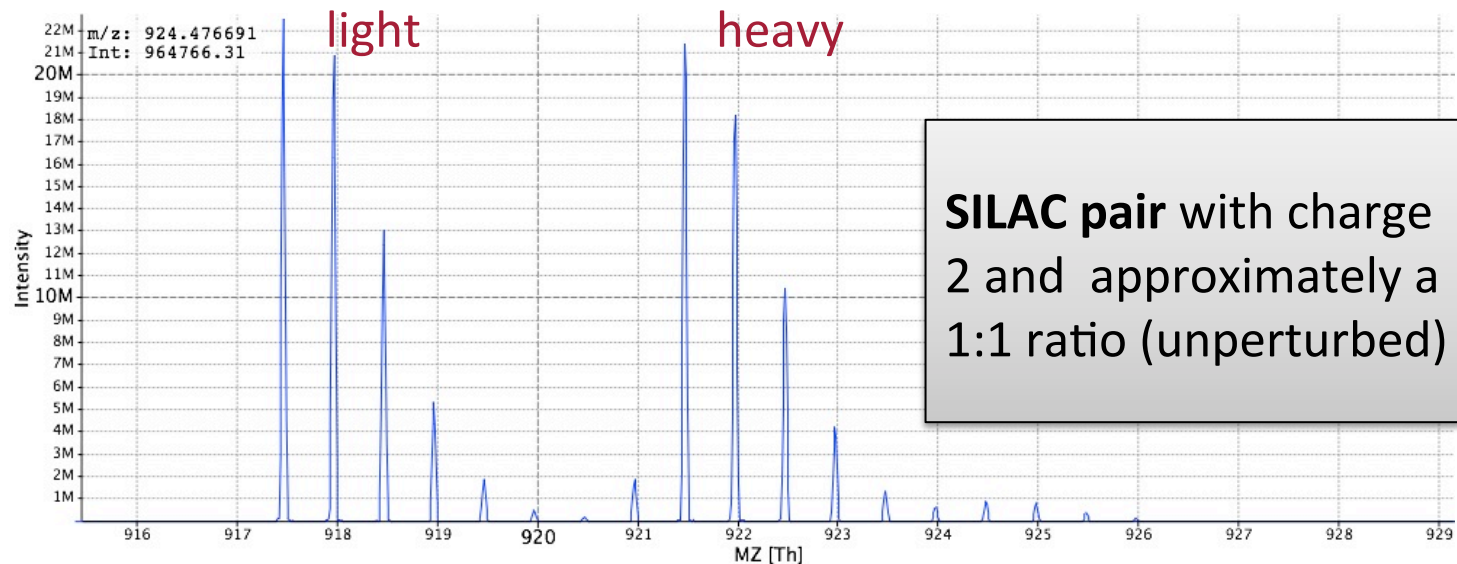
- **Chemical labeling**
  - Peptides are modified chemically after extraction
  - Label is usually attached covalently at specific functional groups (N-terminus, specific side chains, ...)
  - Does not involve a perturbation of the in vivo system
  - Labeling occurs late (during sample preparation) and thus does not account for variance introduced in the early steps
- **Metabolic labeling**
  - Stable isotope labels are integrated by 'feeding' the organism with labeled metabolites (amino acids, nitrogen sources, glucose, ...)
  - Full incorporation of the label can take a while
  - Requires perturbation of the in vivo system, depending on the size quite expensive
  - Labeling occurs early in the study, results in higher reproducibility



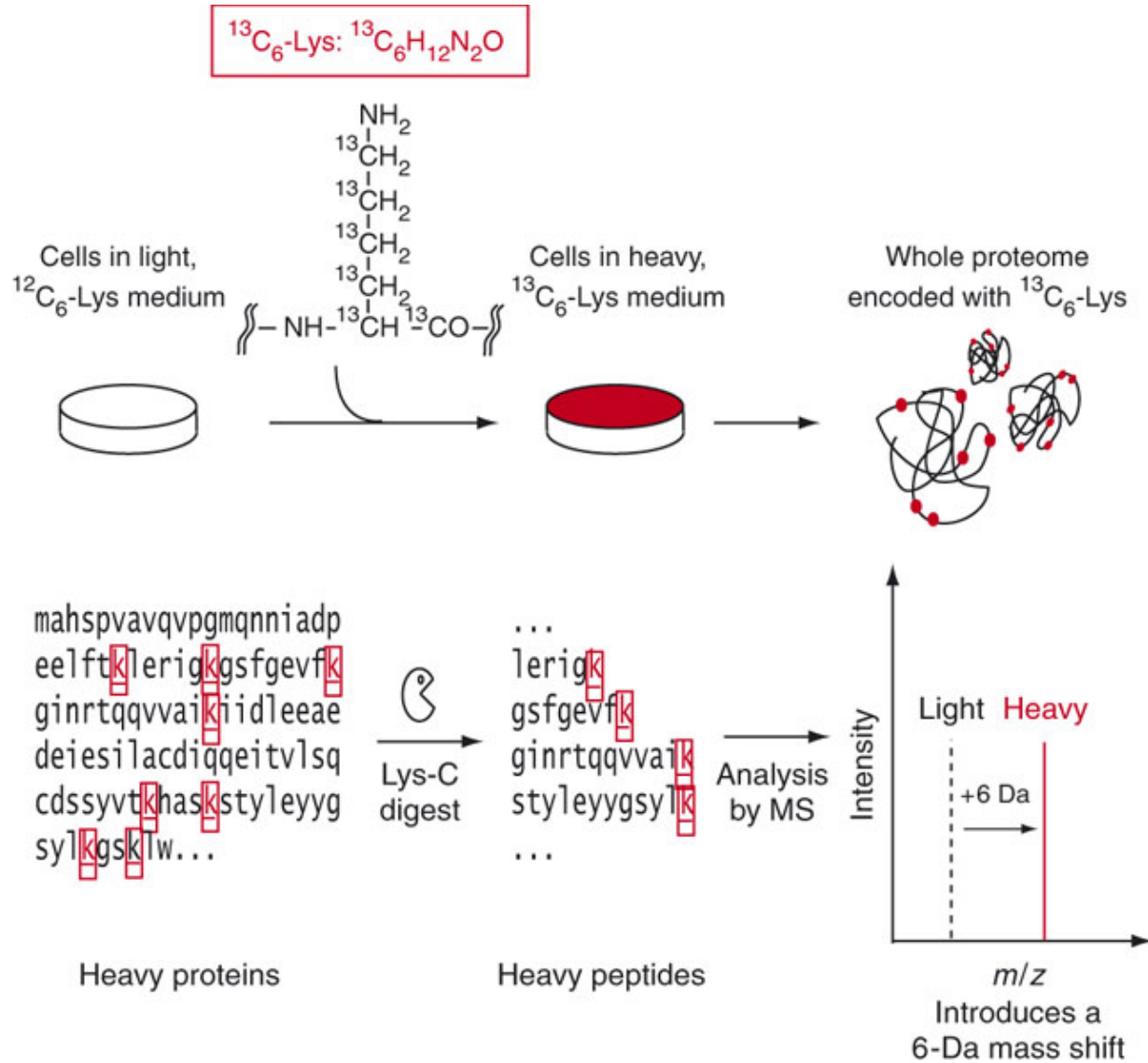
# SILAC

- **SILAC – Stable Isotope Labeling with Amino Acids in Cell Culture**

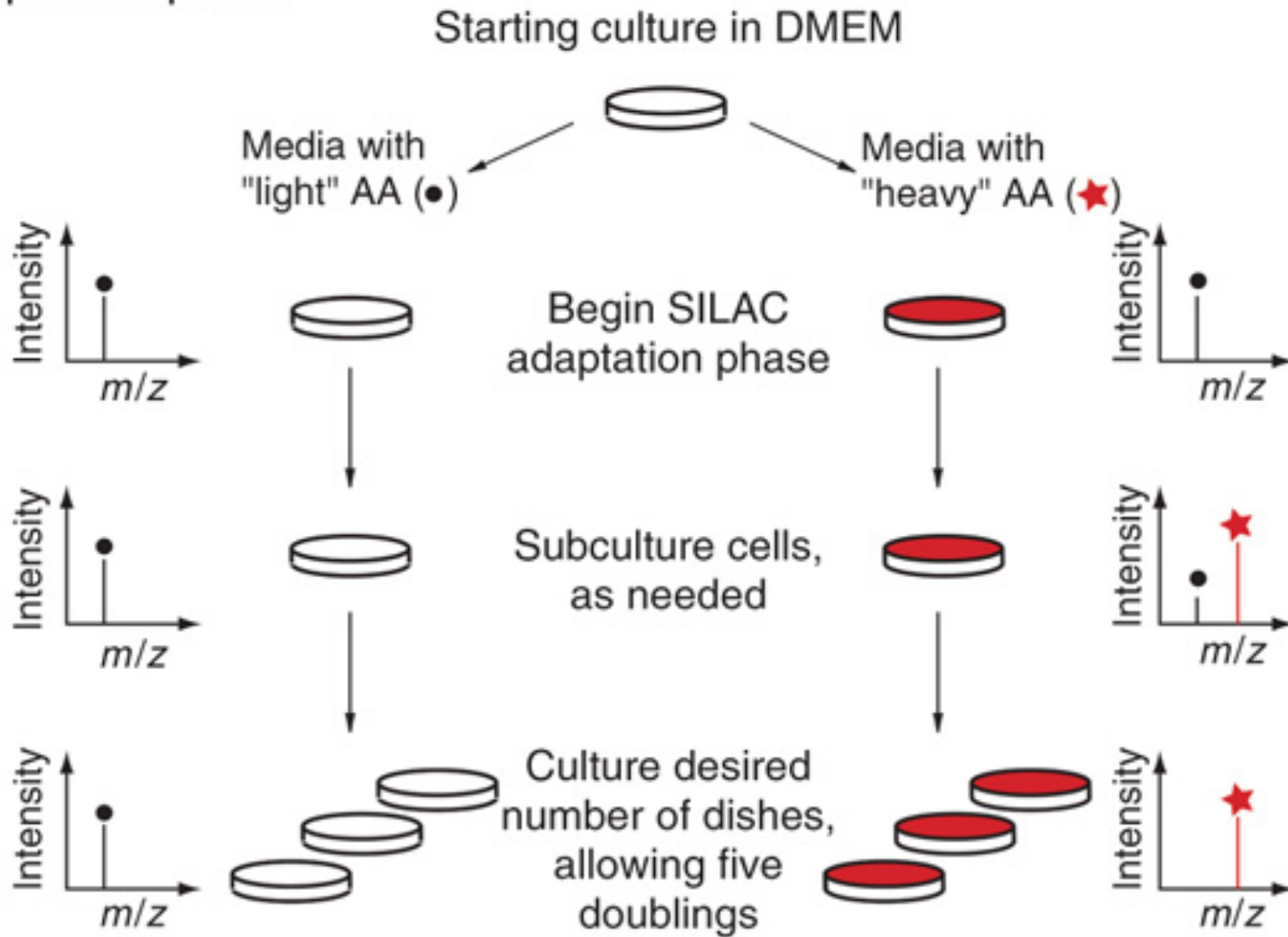
- Introduce stable labels by feeding labeled amino acids to the cell culture
- Labels will be integrated into all proteins after a reasonable amount of time
- Mix and compare with an unlabeled sample
- Tryptic digest ensures that each peptide contains at most one lysine!
- Peptides with heavy and light label are otherwise identical and coelute
- Spectra contain isotope patterns for both heavy and light peptides



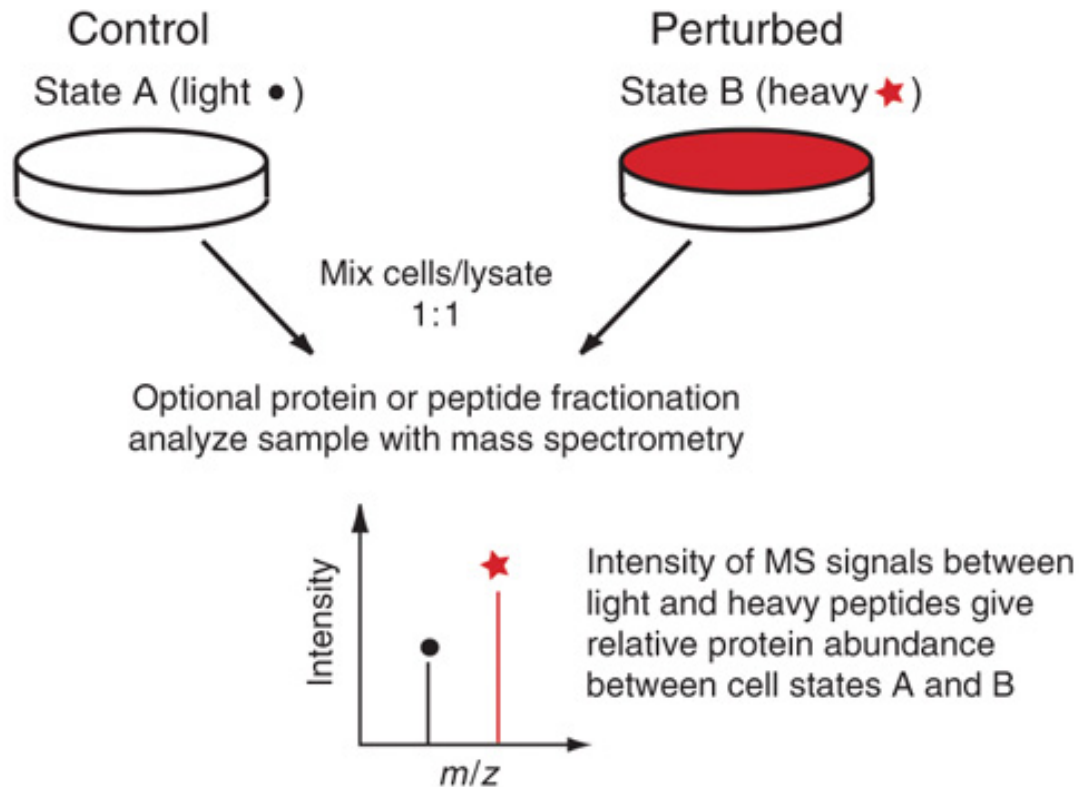
# SILAC



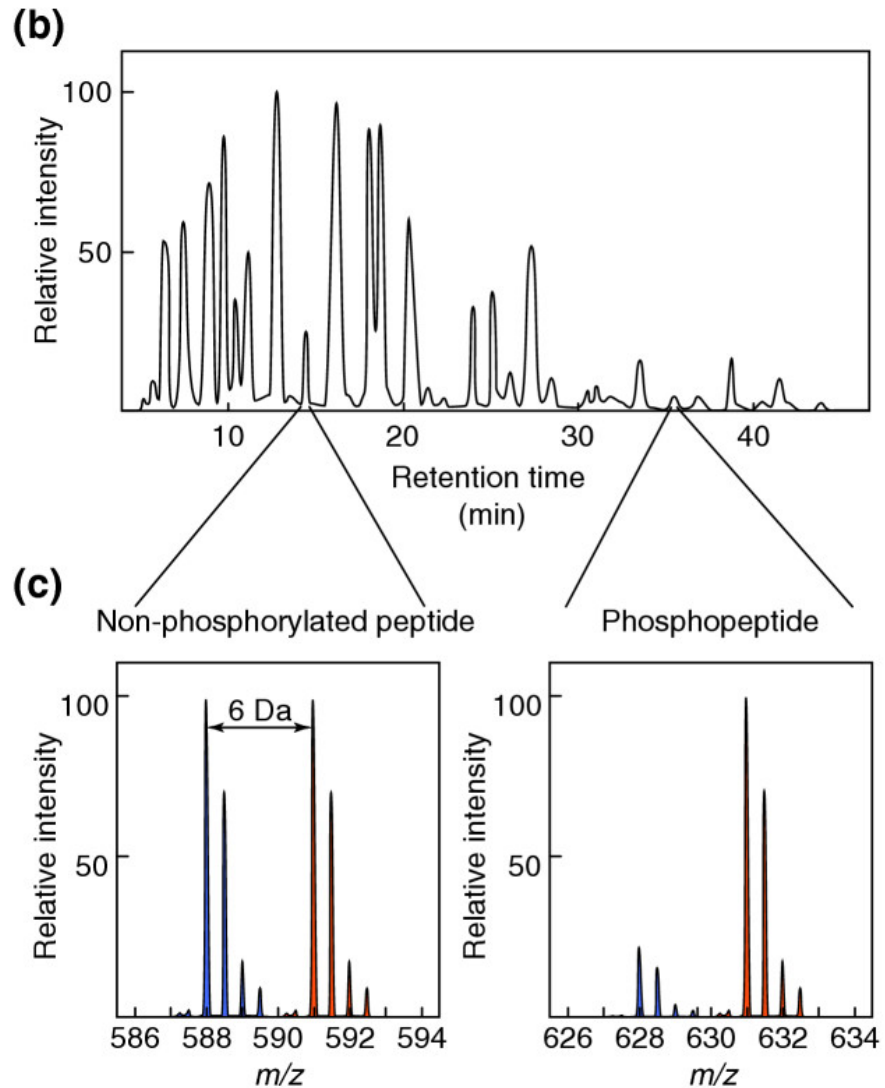
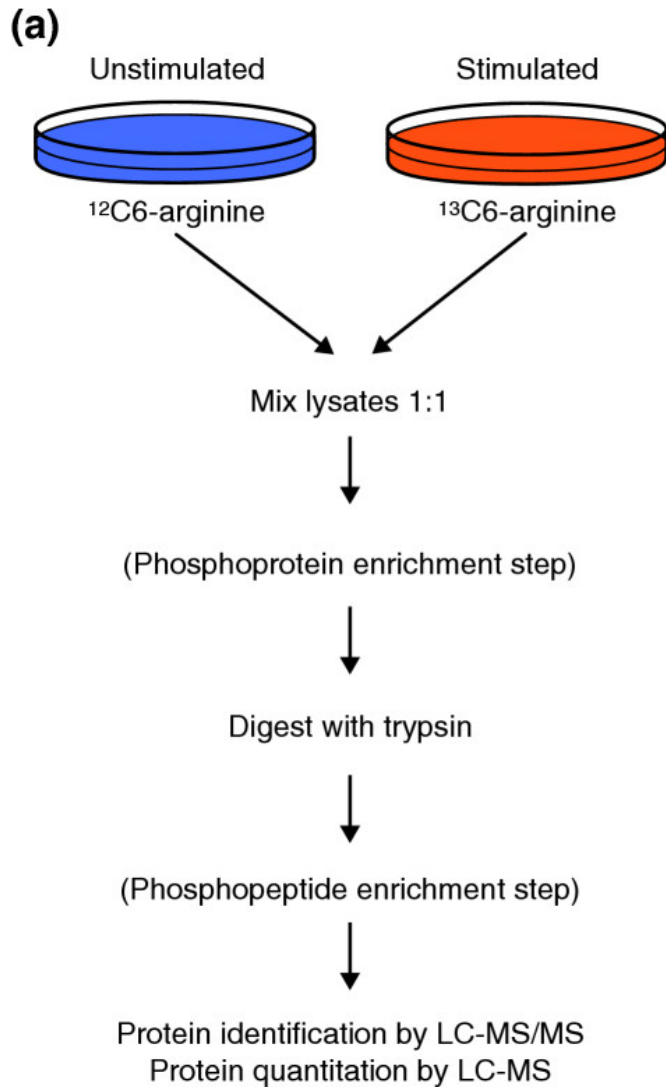
# SILAC



# SILAC

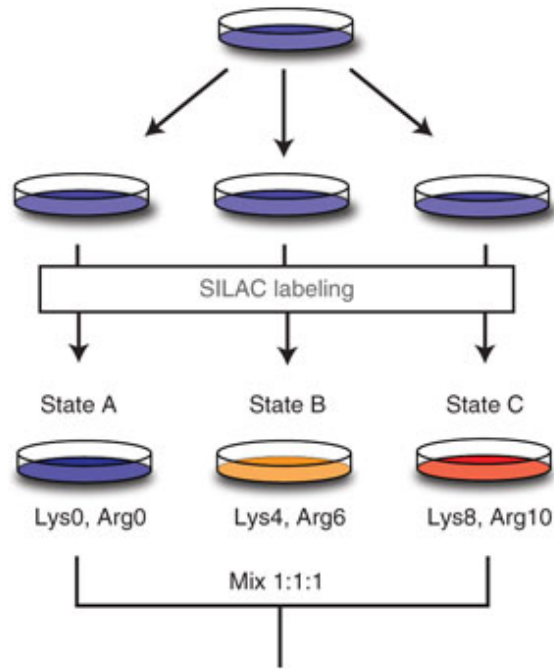


# SILAC

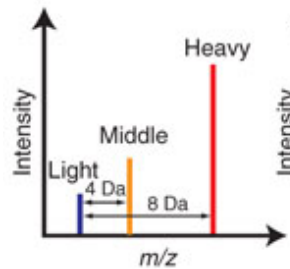


# Spike-In SILAC

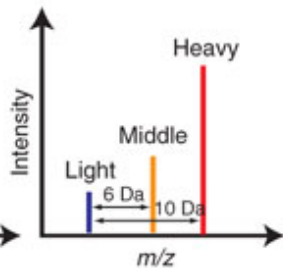
a



Lys-containing peptide



Arg-containing peptide



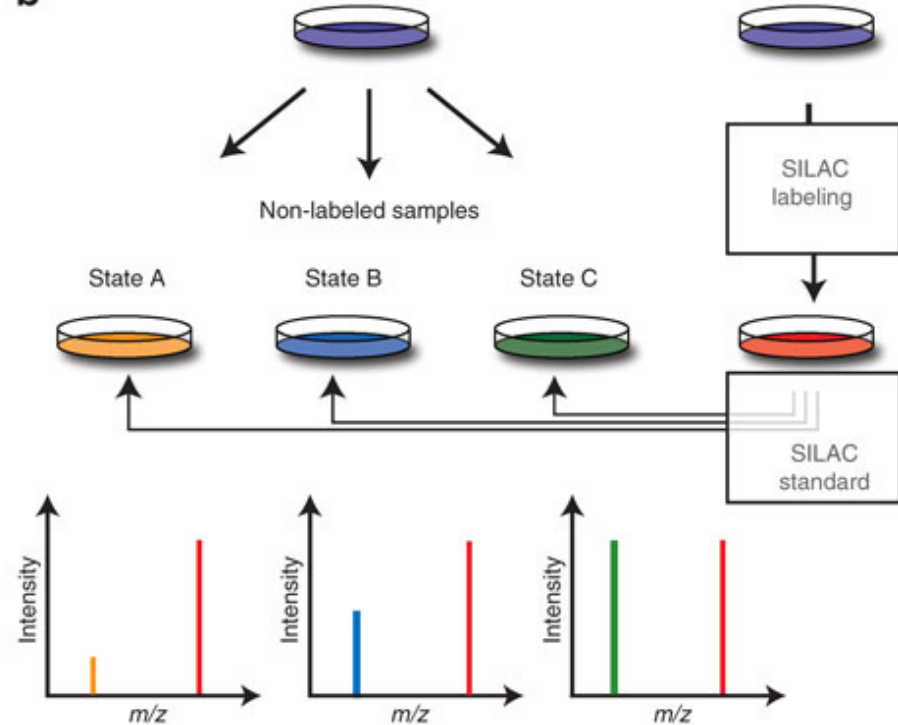
SILAC ratios

$$\text{Ratio}_1 = \frac{\text{Heavy (State C)}}{\text{Light (State A)}}$$

$$\text{Ratio}_2 = \frac{\text{Heavy (State C)}}{\text{Middle (State B)}}$$

$$\text{Ratio}_3 = \frac{\text{Middle (State B)}}{\text{Light (State A)}}$$

b



'Spike-in' SILAC ratios

$$\text{Ratio}_1 = \frac{\text{Heavy (SILAC standard)}}{\text{Light (State A)}}$$

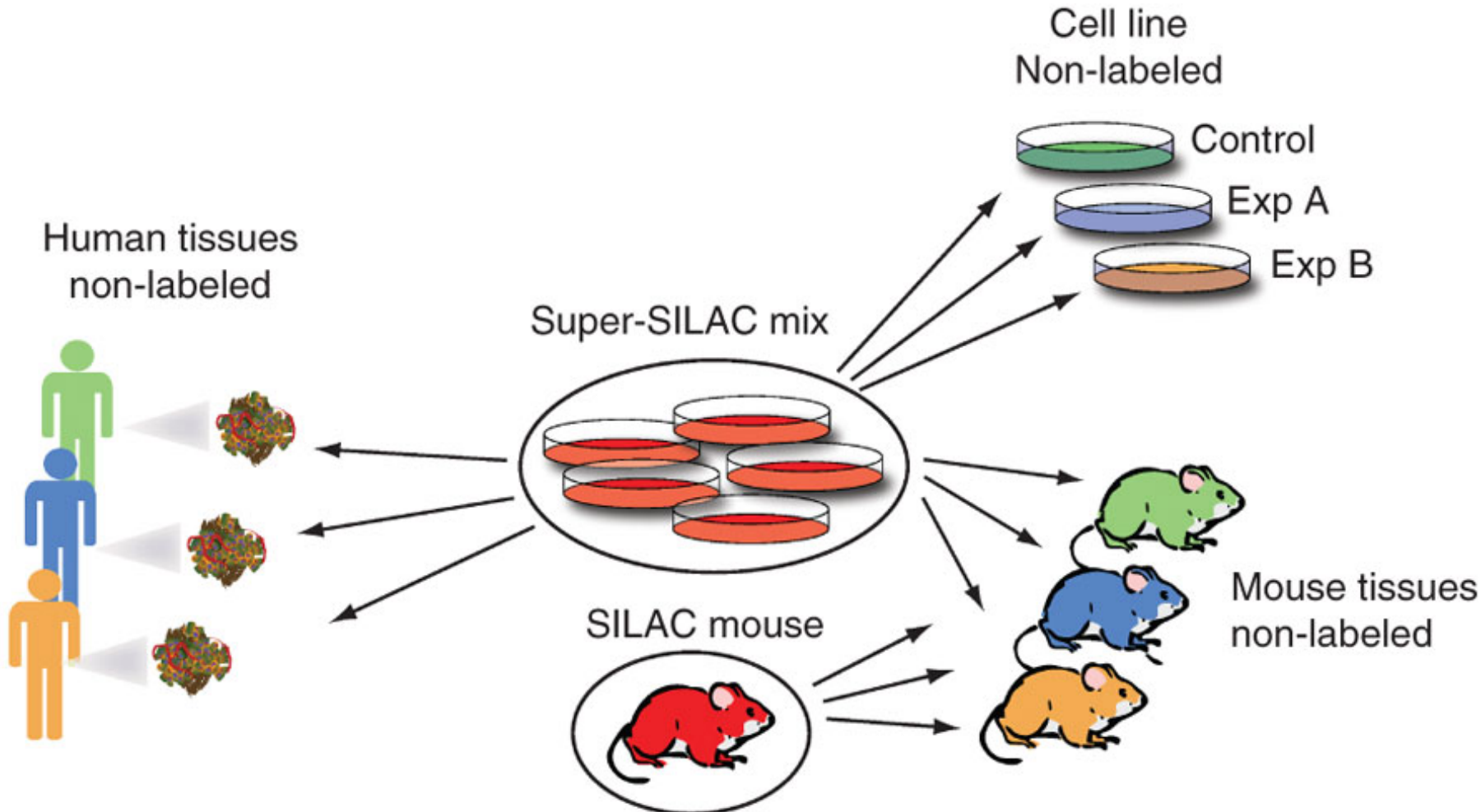
$$\text{Ratio}_2 = \frac{\text{Heavy (SILAC standard)}}{\text{Light (State B)}}$$

$$\text{Ratio}_3 = \frac{\text{Heavy (SILAC standard)}}{\text{Light (State C)}}$$

$$\frac{\text{Ratio}_1}{\text{Ratio}_2} = \frac{\text{Light (State B)}}{\text{Light (State A)}}$$

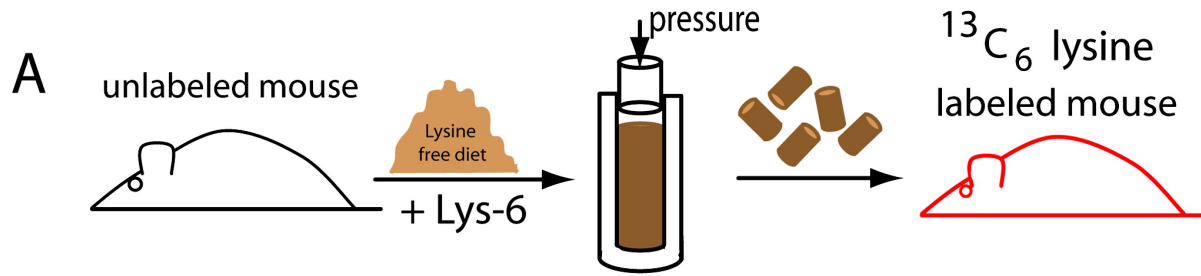
$$\frac{\text{Ratio}_2}{\text{Ratio}_3} = \frac{\text{Light (State C)}}{\text{Light (State B)}}$$

# Spike-In SILAC





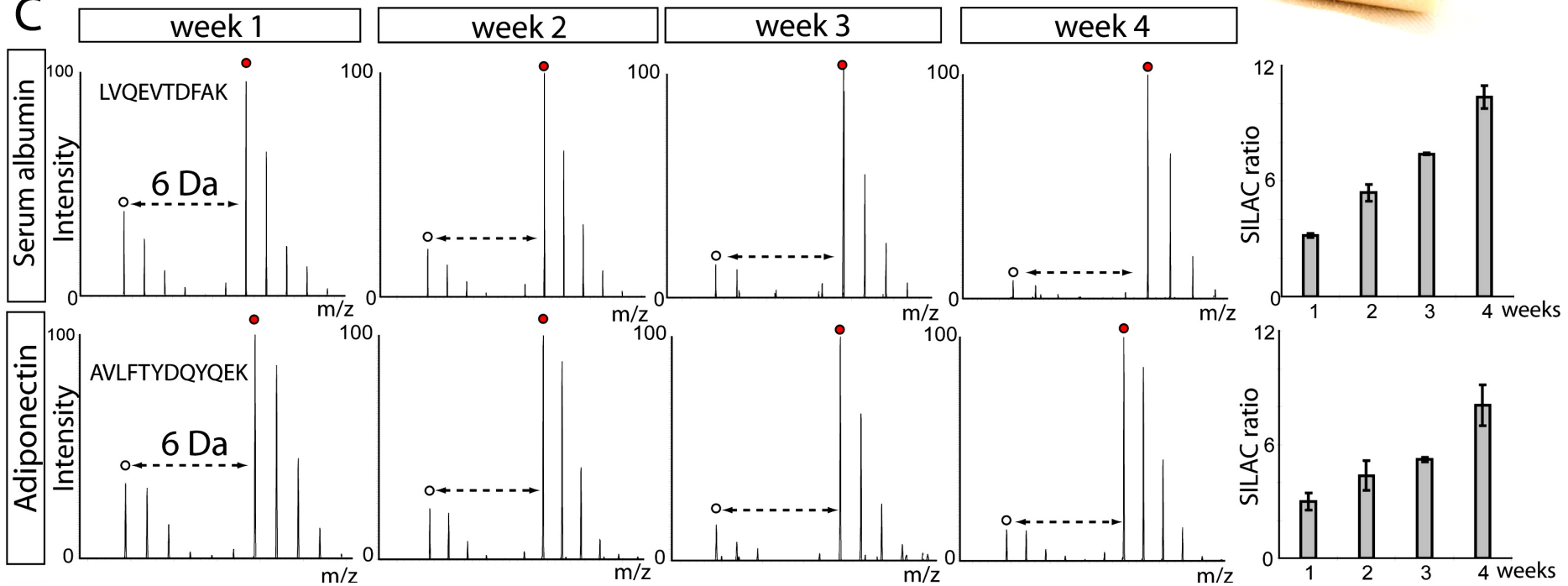
# SILAC Mouse



**B**

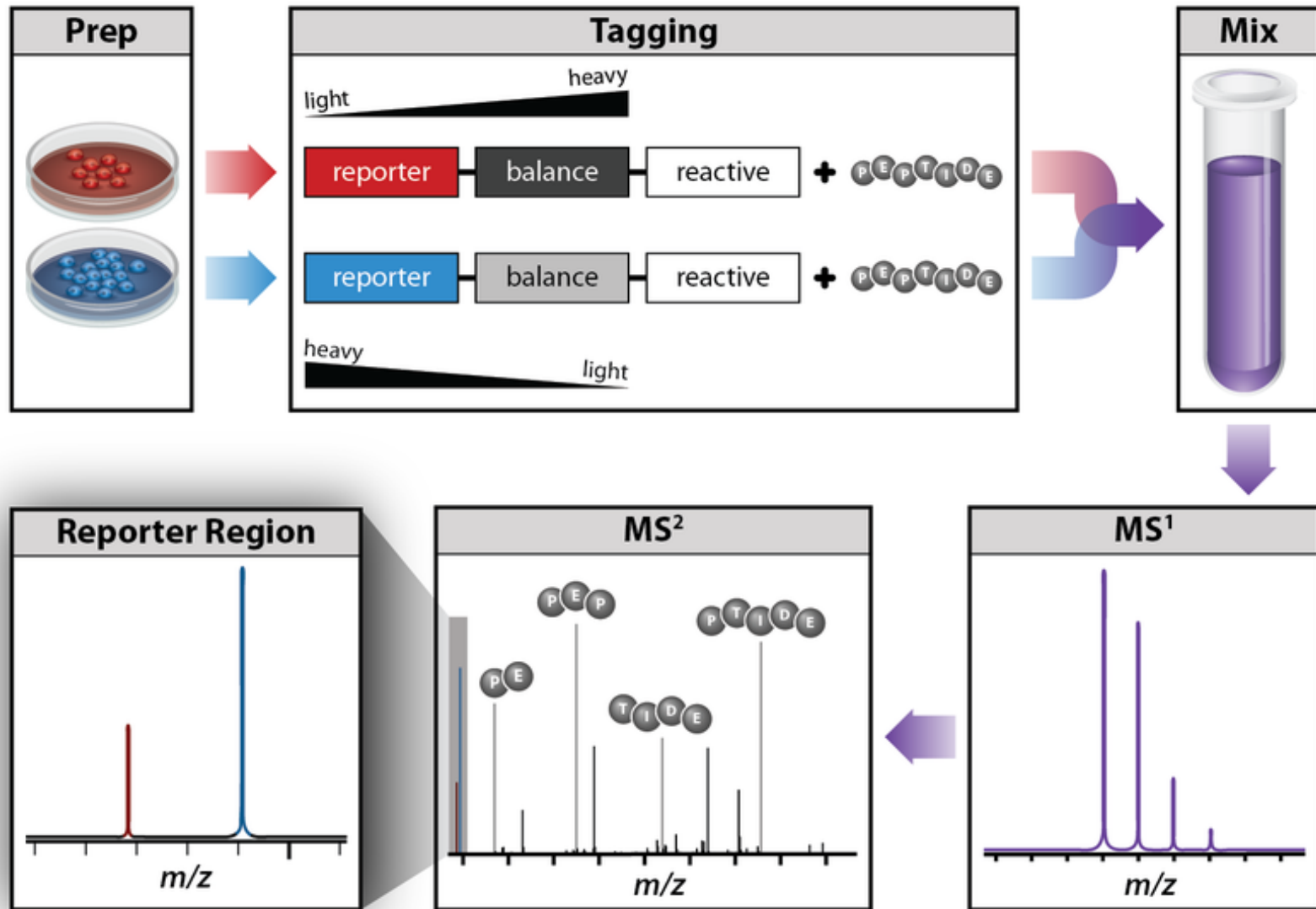


**C**





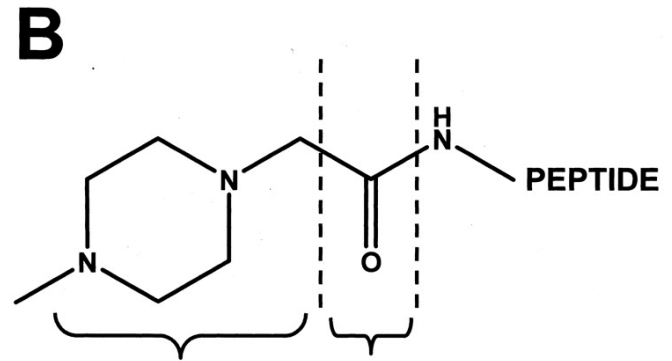
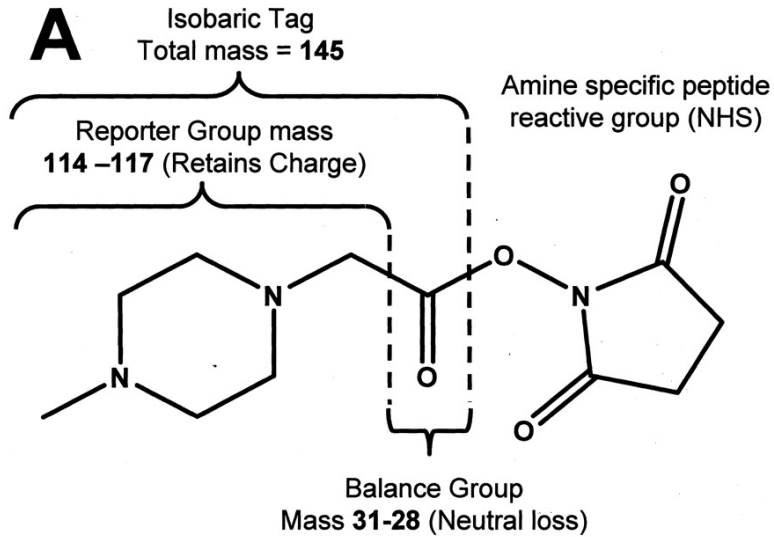
# Isobaric Labeling



# Isobaric Labeling

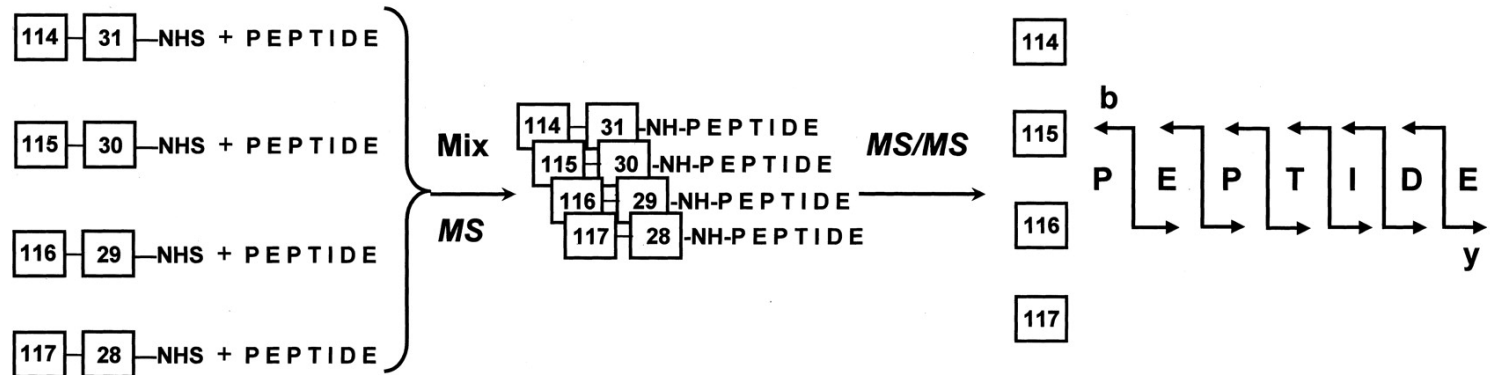
- Idea
  - Label the different samples with labels of the **same mass (isobaric)**
  - Design the label in a way that they fragment differently upon collision-induced dissociation
  - MS<sup>2</sup> spectra will then contain **reporter ions**
  - Quantification and identification are then both based on tandem spectra only
- Key method: **iTRAQ – isobaric tags for relative and absolute quantification**
  - Based on covalent modification of N-terminus of peptides
  - Labeling performed after digestion (also applicable to clinical samples)
  - Kits available for 4 or 8 distinct labels ('quadropex', 'octopex')

# iTRAQ

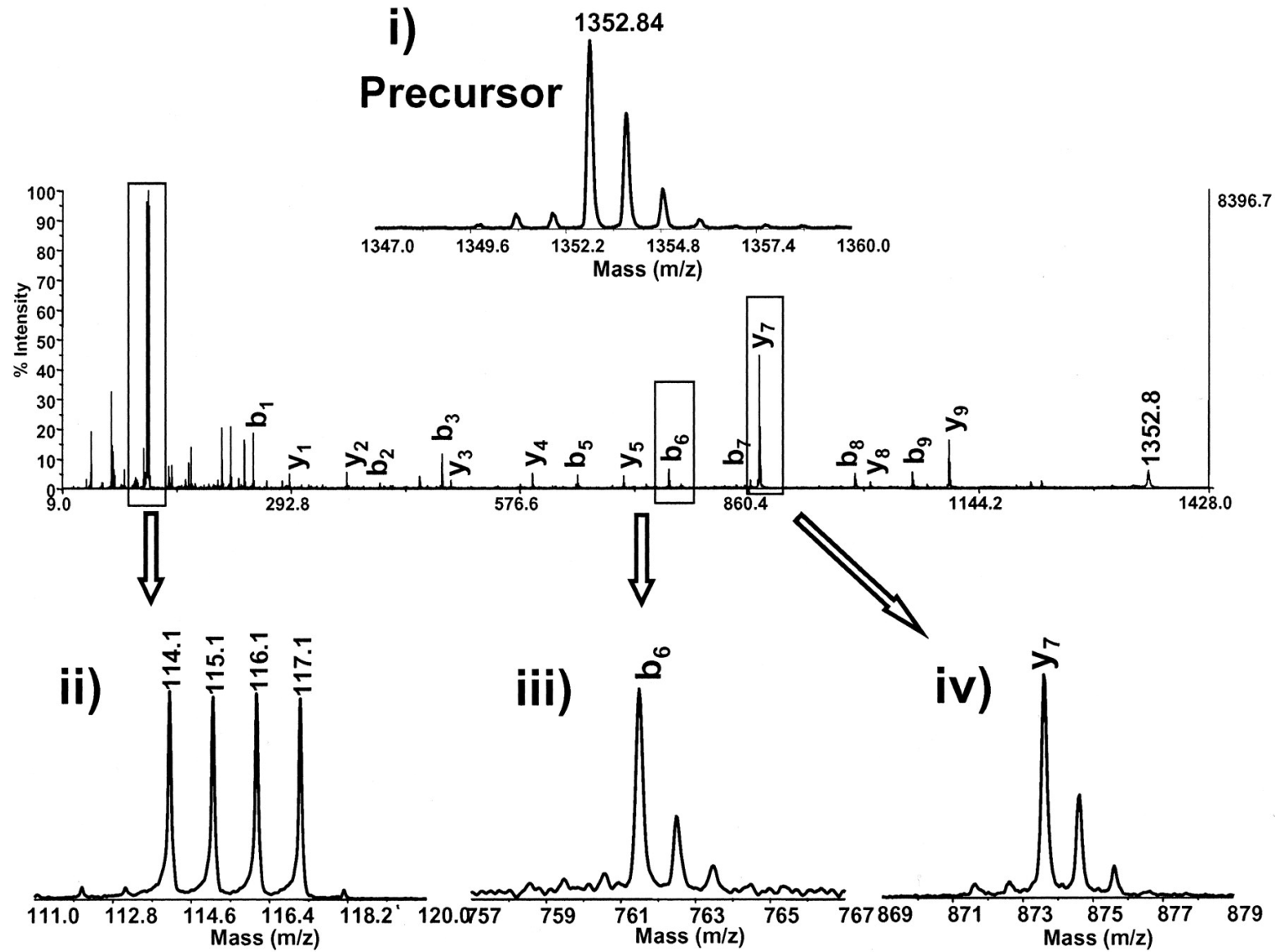


m/z 114 (+1)	<sup>13</sup> C	<sup>13</sup> C <sup>18</sup> O	(+3)
m/z 115 (+2)	<sup>13</sup> C <sub>2</sub>	<sup>18</sup> O	(+2)
m/z 116 (+3)	<sup>13</sup> C <sub>2</sub> <sup>15</sup> N	<sup>13</sup> C	(+1)
m/z 117 (+4)	<sup>13</sup> C <sub>3</sub> <sup>15</sup> N		(+0)

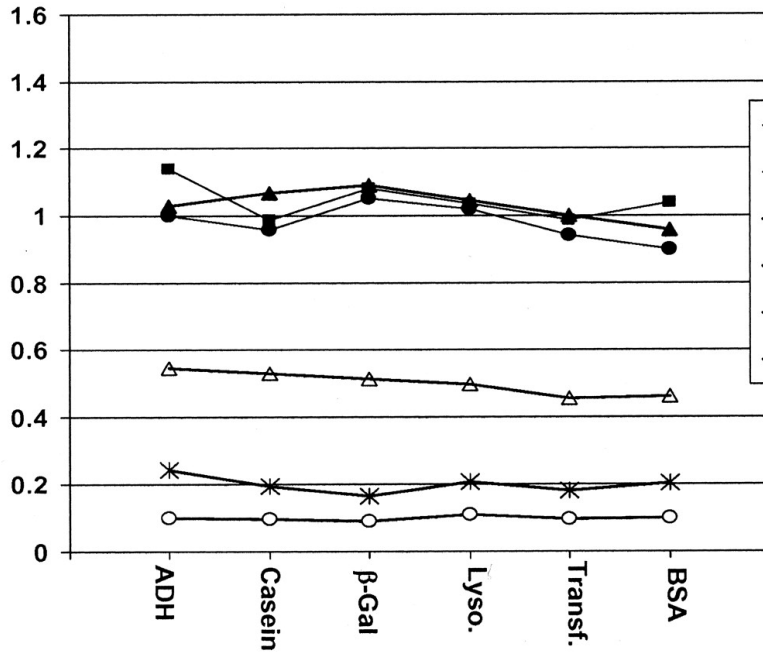
**C**



# iTRAQ

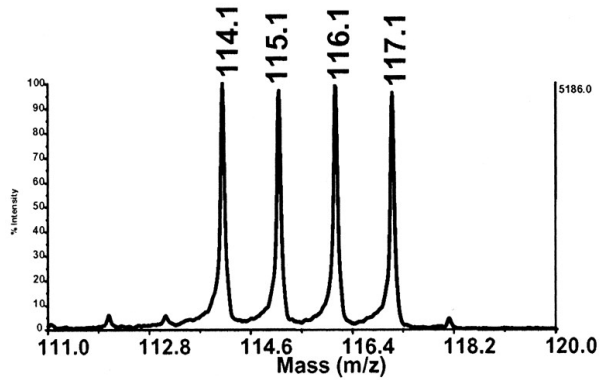


# iTRAQ

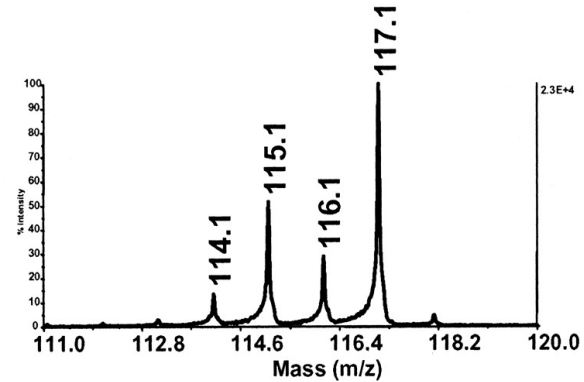


Ratio	Mean	SD
1:1	1.03	0.16
1:2	0.514	0.12
1:5	0.204	.045
1:10	0.097	0.023

**1:1:1:1 Mixture**

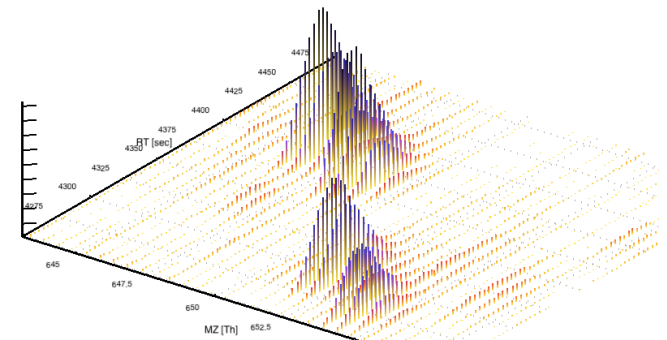
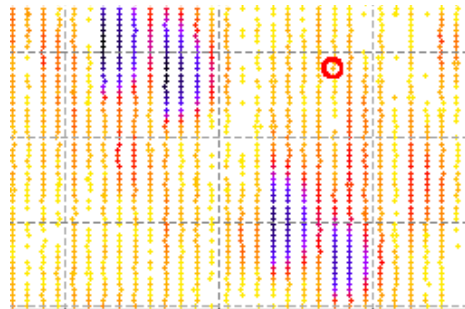
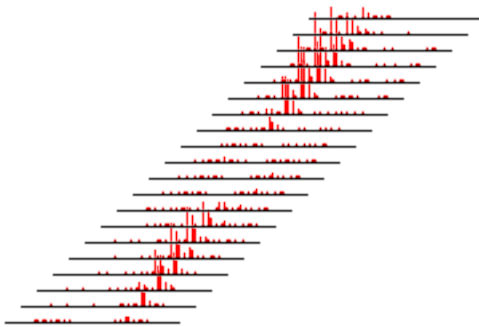


**1:5:2:10 Mixture**



# Quantitative Data – LC-MS Maps

- Spectra are acquired with rates up to dozens per second
- Stacking the spectra yields **maps**
- Resolution:
  - Up to millions of points per spectrum
  - Tens of thousands of spectra per LC run
- Huge 2D datasets of up to hundreds of GB per sample
- MS intensity follows the chromatographic concentration

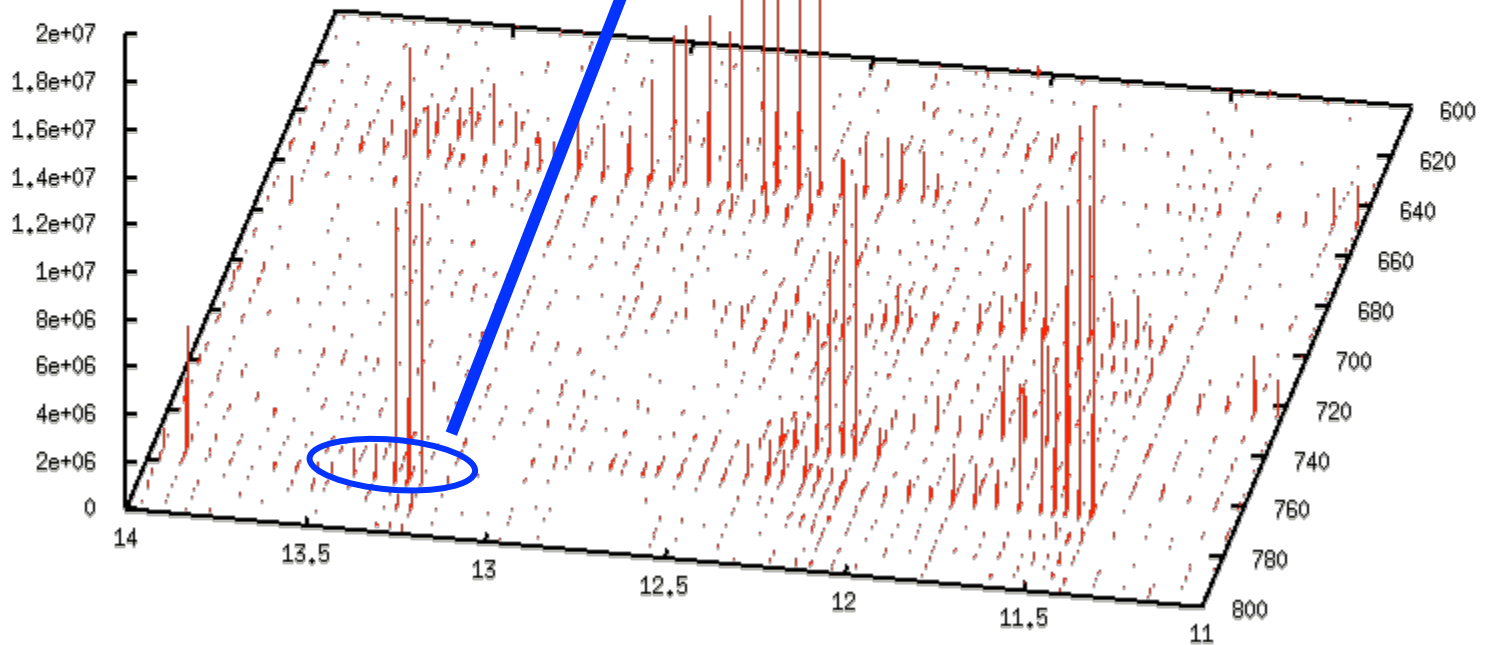


# LC-MS Data (Map)

-----

## Quantification

(15 nmol/ $\mu$ l, 3x over-expressed, ...)



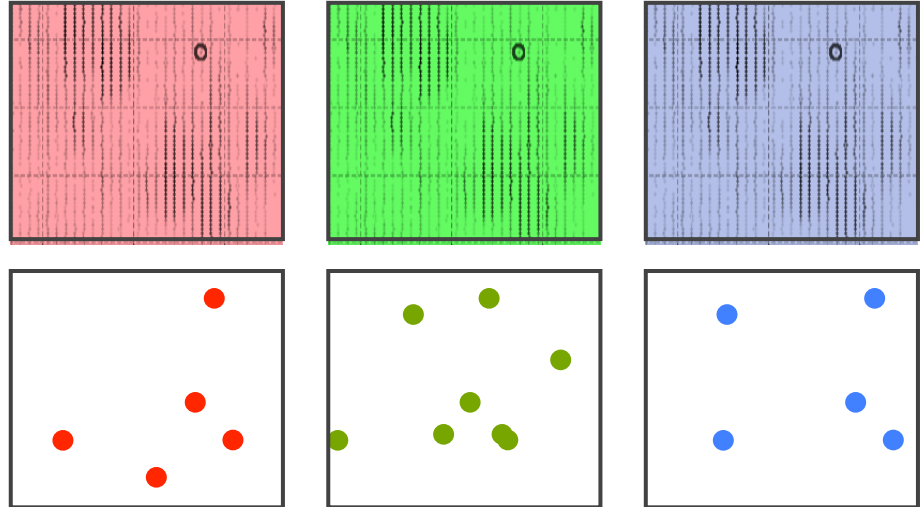
# Label-Free Quantification (LFQ)

- Label-free quantification is probably the most natural way of quantifying
  - **No labeling required**, removing further sources of error, no restriction on sample generation, cheap
  - Data on different samples acquired in different measurements – **higher reproducibility needed**
  - **Manual analysis difficult**
  - **Scales very well** with the number of samples, basically no limit, no difference in the analysis between 2 or 100 samples



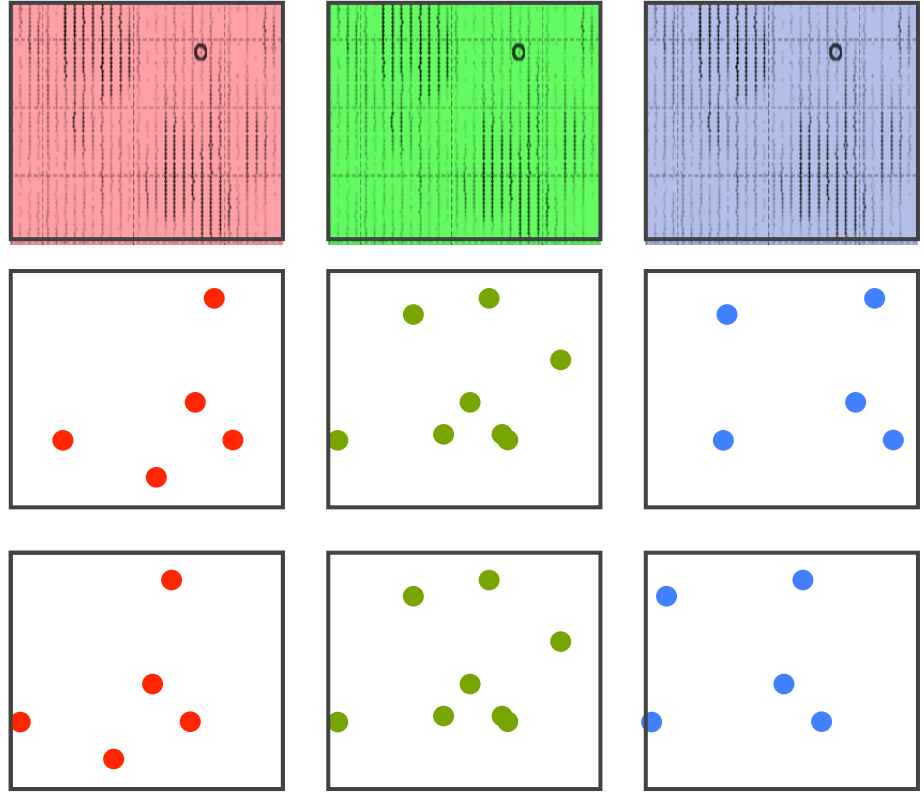
# LFQ – Analysis Strategy

1. Find features in all maps



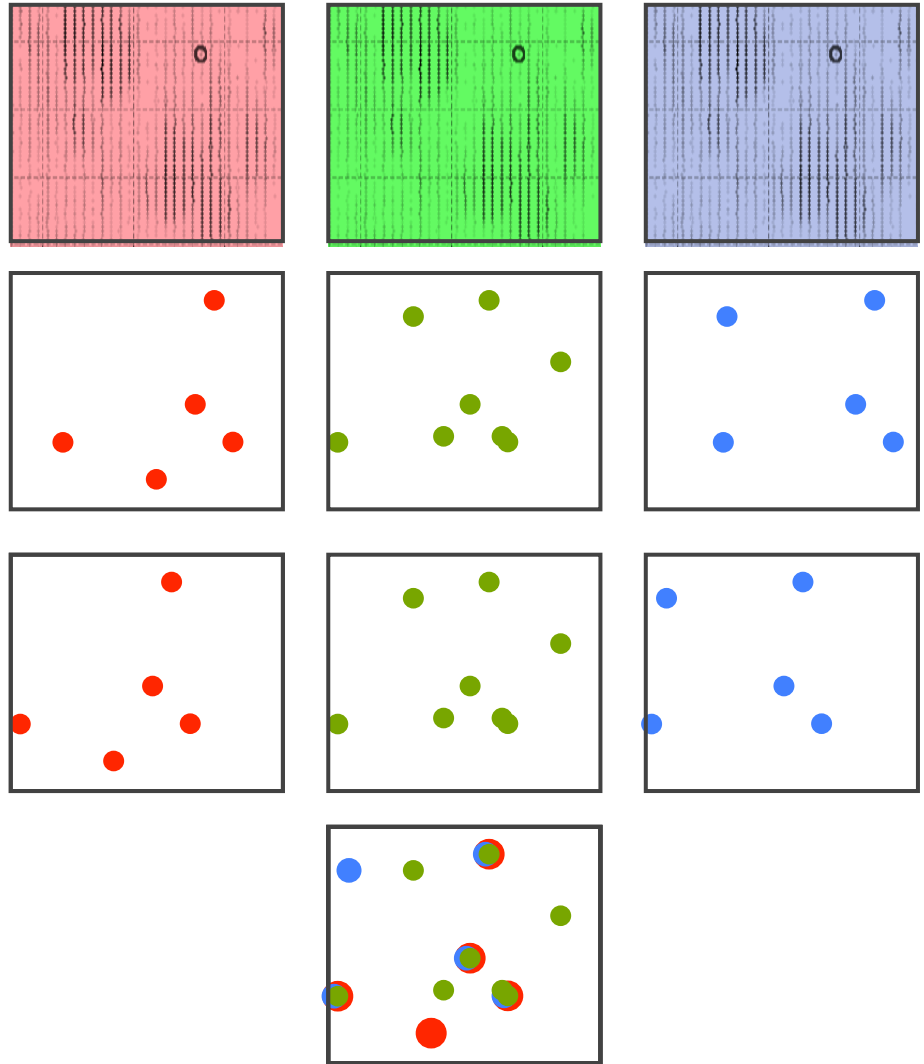
# LFQ – Analysis Strategy

1. **Find** features in all maps
2. **Align** maps



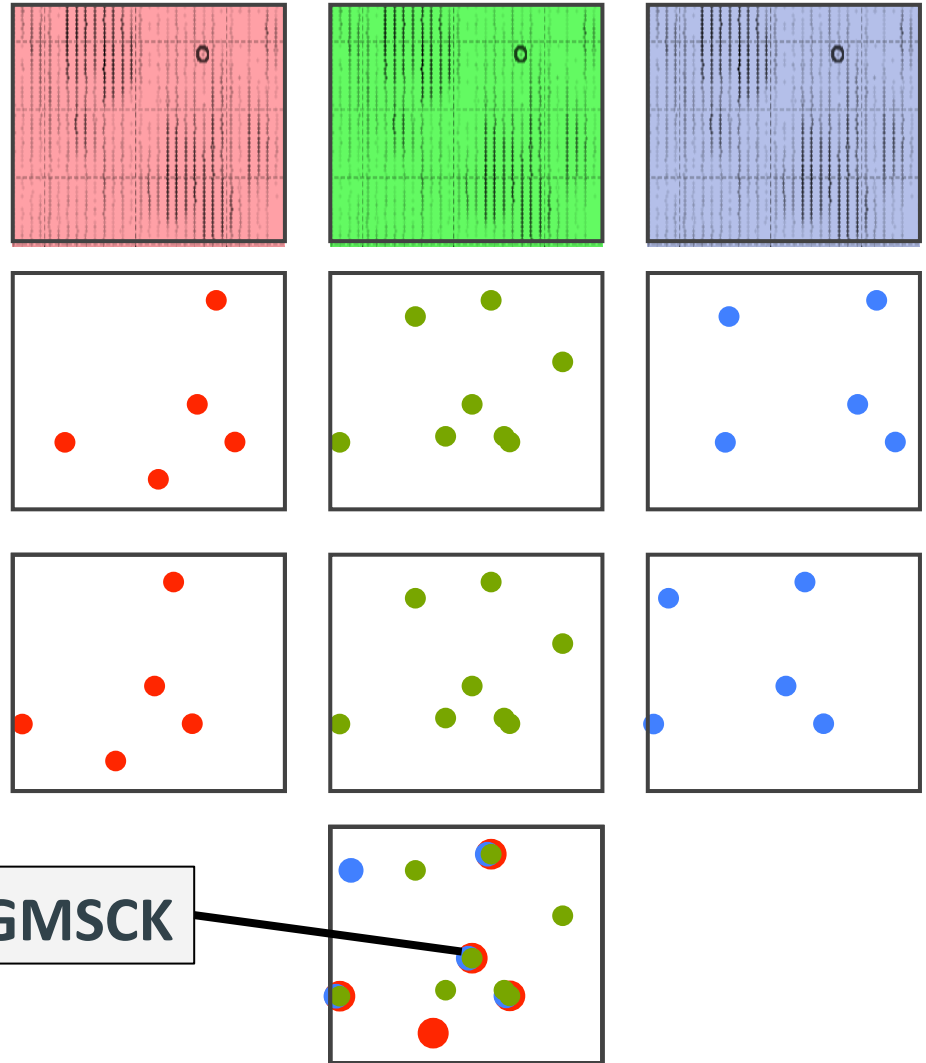
# LFQ – Analysis Strategy

1. **Find** features in all maps
2. **Align** maps
3. **Link** corresponding features



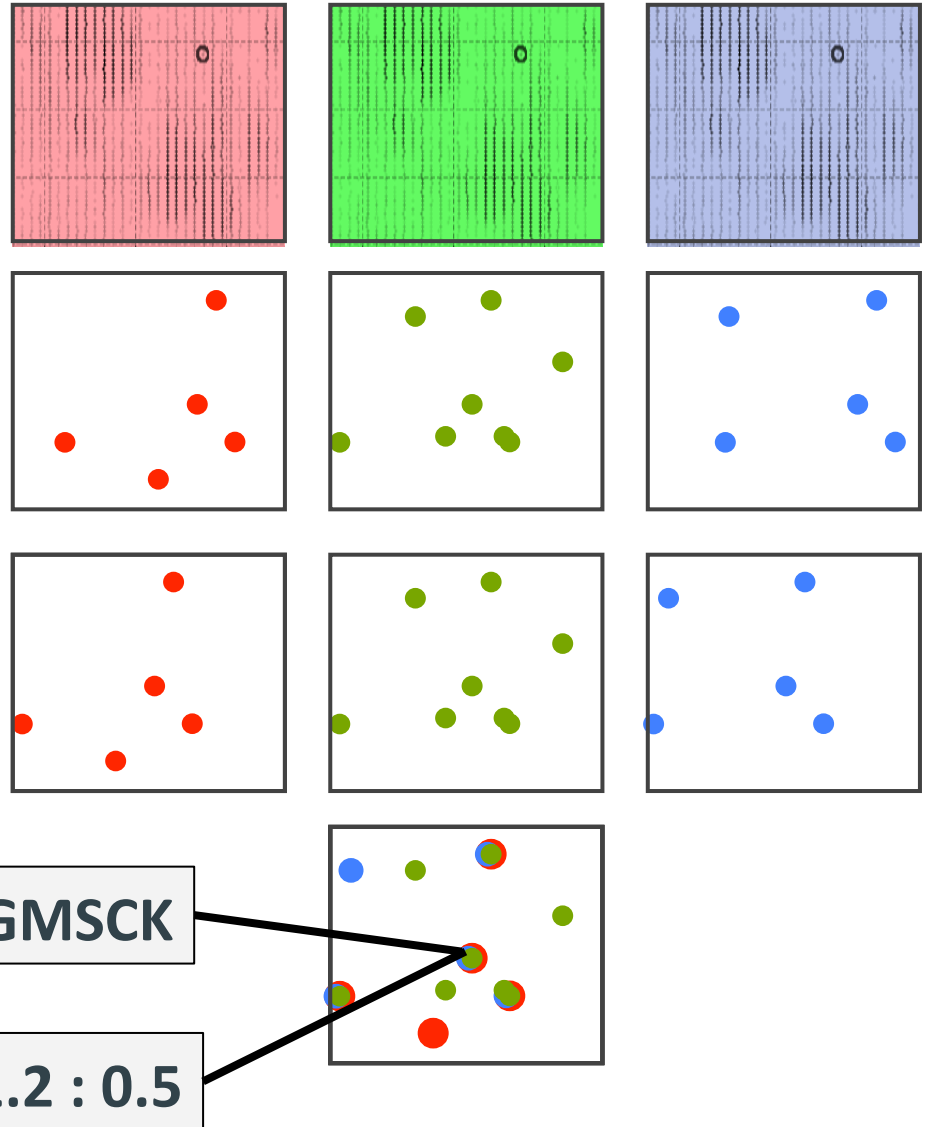
# LFQ – Analysis Strategy

1. **Find** features in all maps
2. **Align** maps
3. **Link** corresponding features
4. **Identify** features



# LFQ – Analysis Strategy

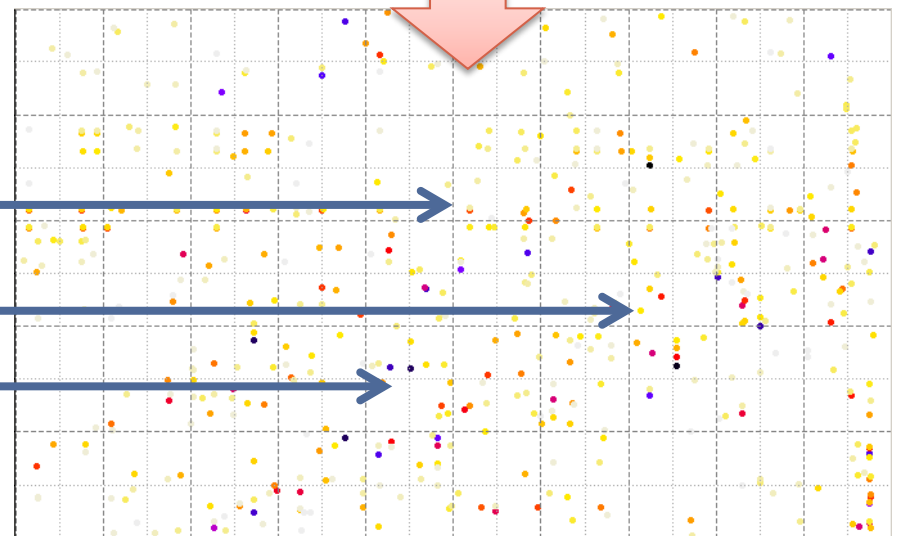
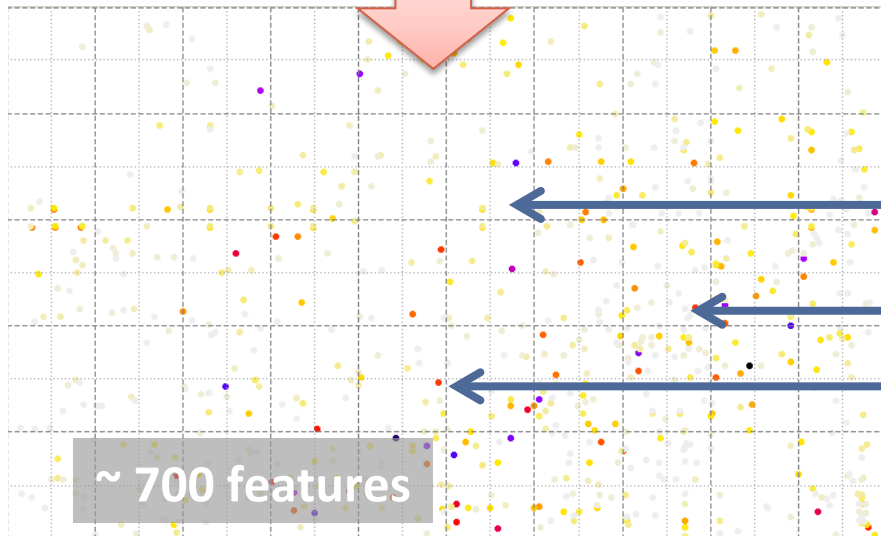
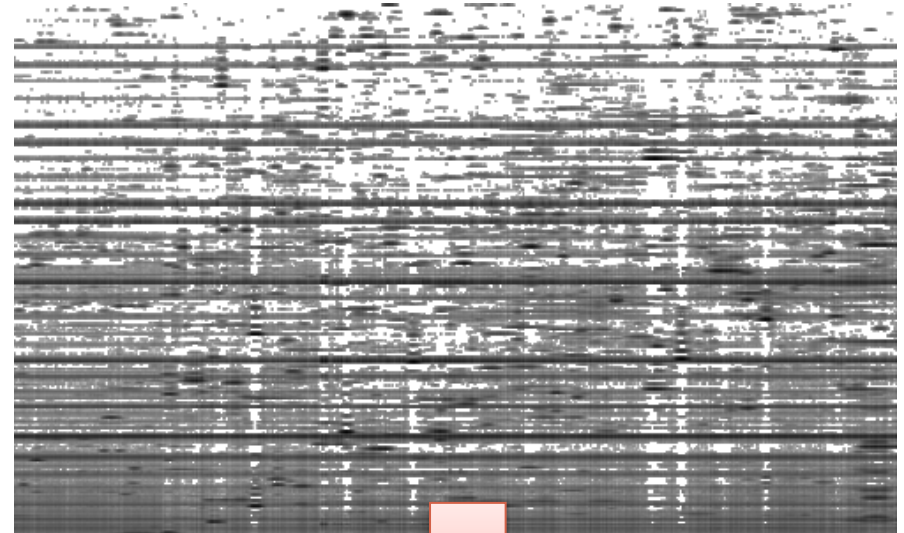
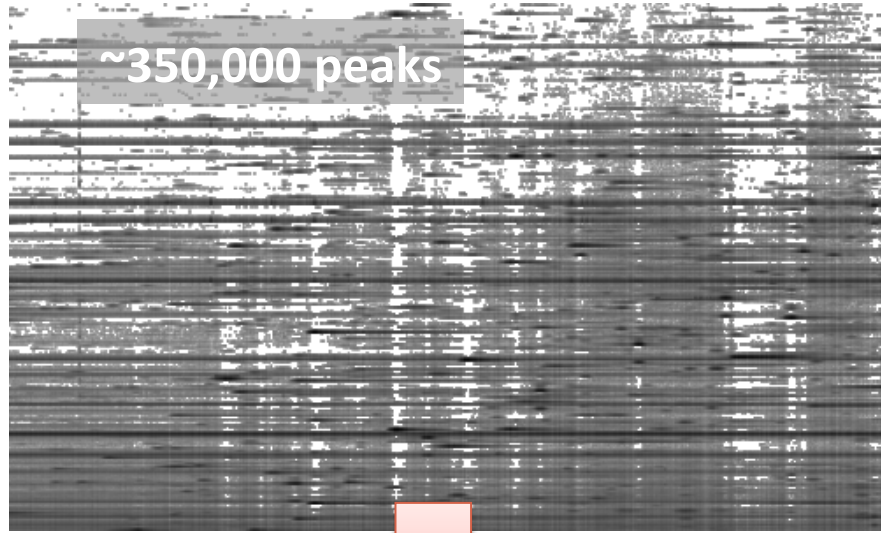
1. **Find** features in all maps
2. **Align** maps
3. **Link** corresponding features
4. **Identify** features
5. **Quantify**



# Feature-Based Alignment

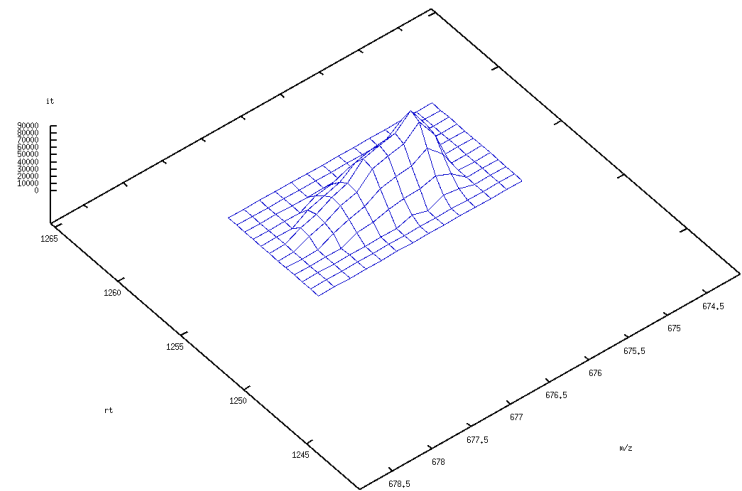
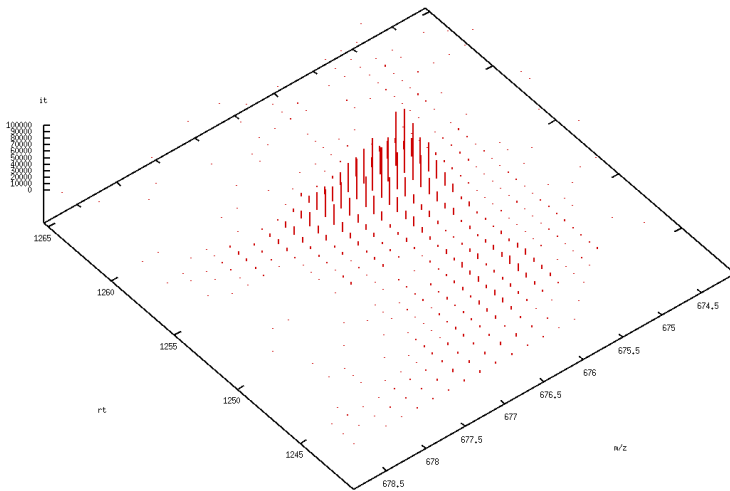
- LC-MS maps can contain millions of peaks
- Retention time of peptides and metabolites can shift between experiments
- In label-free quantification, maps thus need to be aligned in order to identify corresponding features
- Alignment can be done on the raw maps (where it is usually called 'dewarping') or on already identified features
- The latter is simpler, as it does not require the alignment of millions of peaks, but just of tens of thousands of features
- Disadvantage: it relies on an accurate feature finding

# Feature-Based Alignment



# Feature Finding

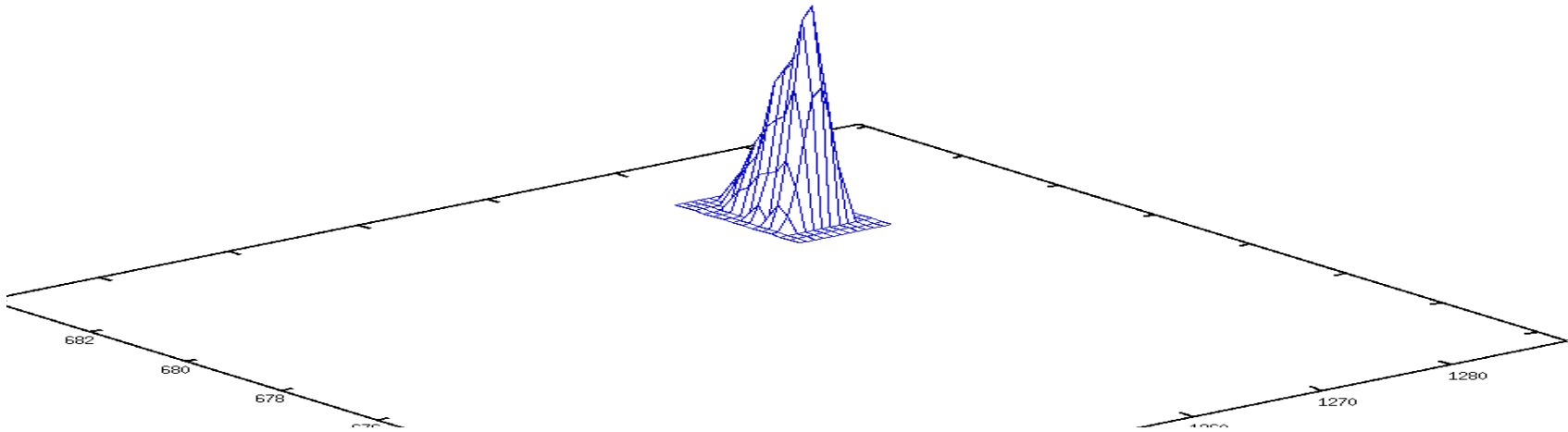
- Identify all peaks belonging to one peptide
- Key idea:
  - Identify suspicious regions
  - Fit a **model** to that region and identify peaks explained by it





# Feature Finding

- **Extension:** collect all data points close to the seed
- **Refinement:** remove peaks that are not consistent with the model
- **Fit an optimal model** for the reduced set of peaks
- Iterate this until no further improvement can be achieved

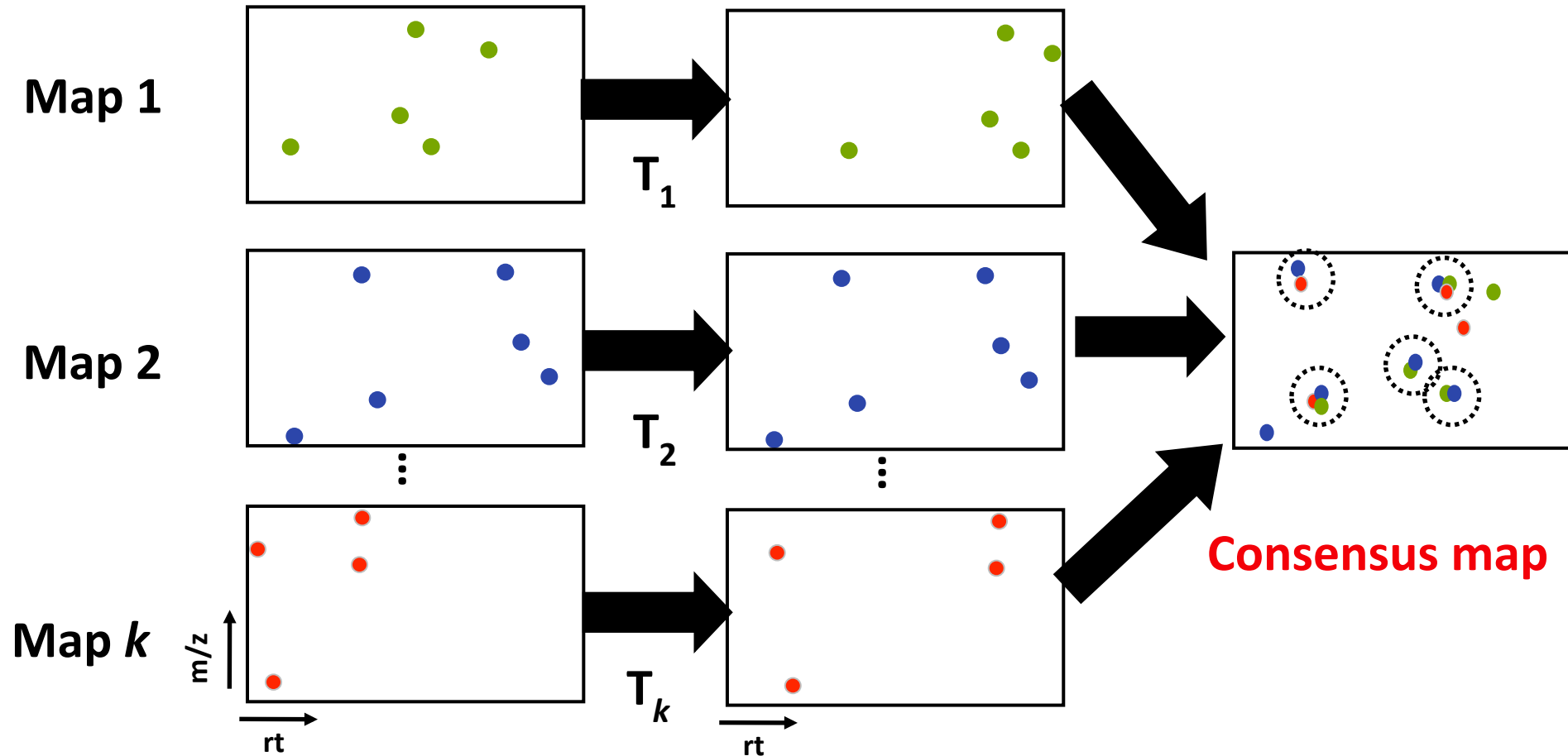


# Linear Alignment

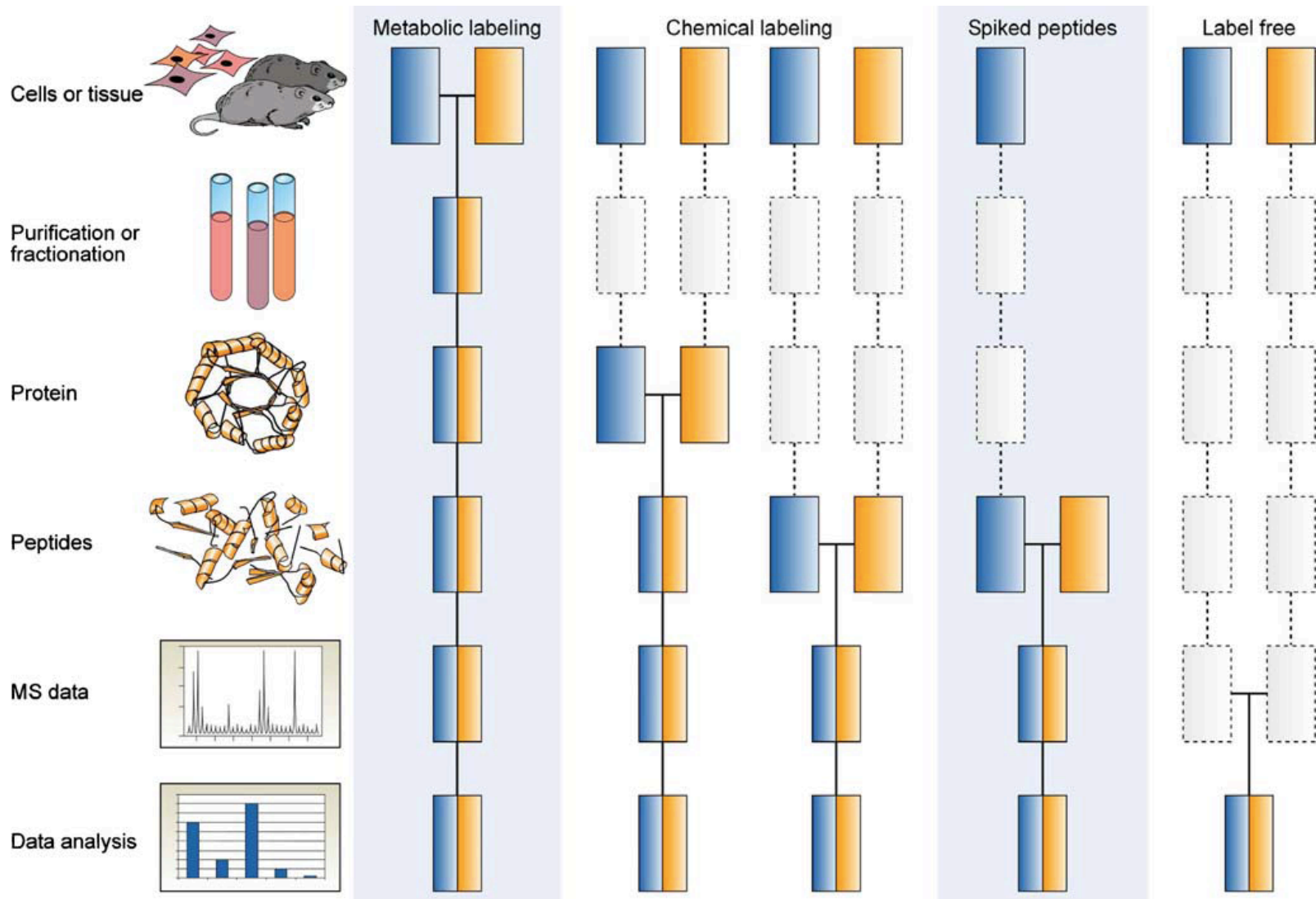
- Lange et al. proposed an efficient feature-based alignment of maps based on pose clustering
- The algorithm takes a pair of maps and computes an optimal linear alignment
- It can be applied for multiple alignment of an arbitrary amount of maps by applying it multiply and align the maps in a star-like fashion onto one reference map ( $k-1$  alignments for  $k$  maps)
- The algorithm relies on accurate feature detection but is rather runtime efficient

# Multiple Alignment

- Dewarp  $k$  maps onto a comparable coordinate system
- Choose one map (usually the one with the largest number of features) as reference map (here: map 2  $\rightarrow T_2 = \mathbf{1}$ )



# Quantification Strategies



Common quantitative mass spectrometry workflows. Boxes in blue and yellow represent two experimental conditions. Horizontal lines indicate when samples are combined. Dashed lines indicate points at which experimental variation and thus quantification errors can occur.

# Materials

- Quantification in general:
  - Bantscheff *et al.*, Quantitative mass spectrometry in proteomics: a critical review, *Anal Bioanal Chem* (2005), 389, 1017-1031 [PMID: 17668192]
- Experimental methods
  - SILAC: Ong, Mann, *Nat Prot* 1 (2007), 2650-2660.
  - iTRAQ: Ross *et al.*, *Mol Cell Prot* (2004), 3, 1154-1169.
- Pose clustering algorithm
  - Lange *et al.*, A geometric approach for the alignment of liquid-chromatography—mass spectrometry data, *Bioinformatics* (2007), 23:i273-i281 [PMID: 17646306]
- Nonlinear alignment
  - Podwojski *et al.*, Retention time alignment algorithms for LC/MS data must consider non-linear shifts, *Bioinformatics* (2009), 25 (6): 758-764. [PMID: 19176558]

# Materials

- Online Materials
  - Learning Unit 4[A,B,C]
- Background
  - Chromatography: Learning Unit 2A
  - Statistical concepts: Learning Unit 3A