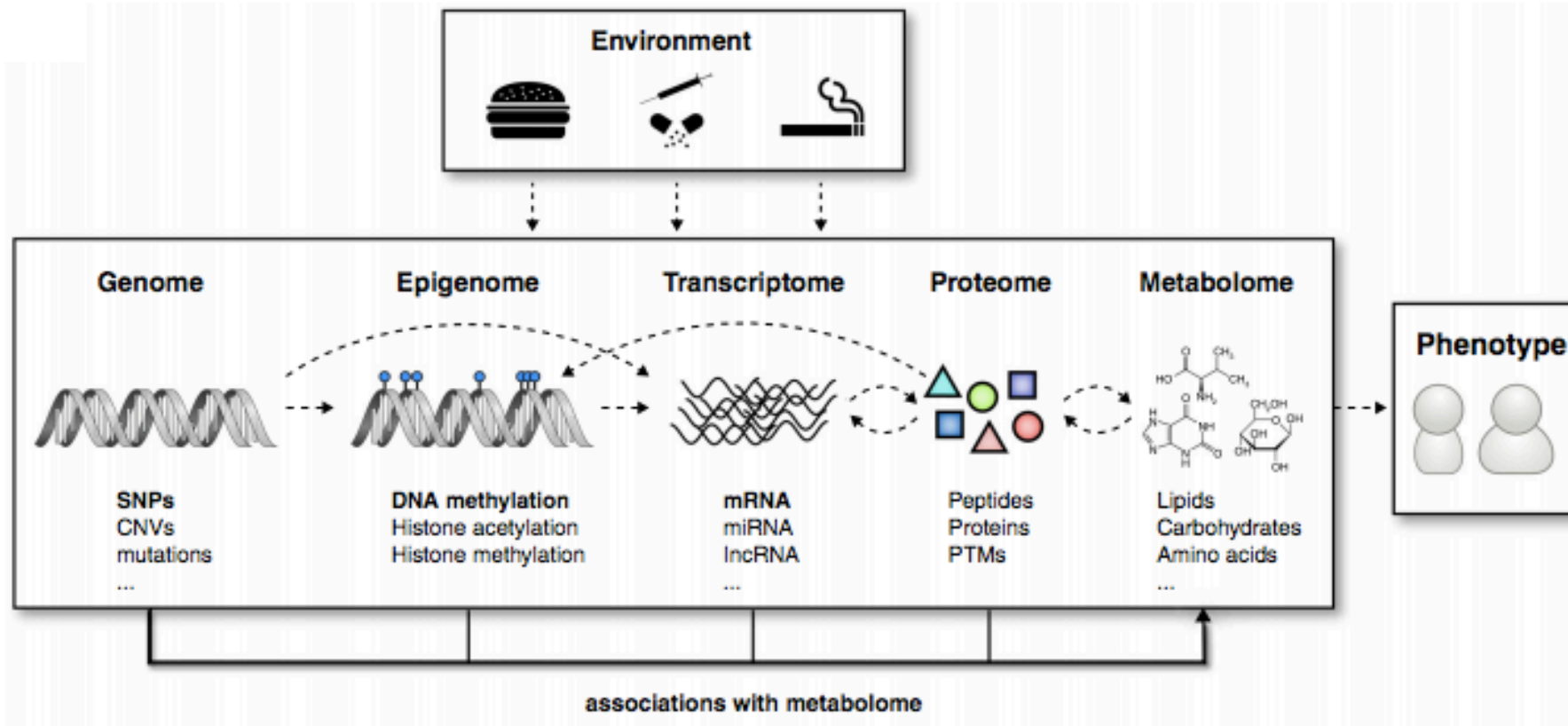# Intro to MS data processing for metabolomics and lipidomics

Ilya Kurochkin

May 6, 2017

# Systems Biology
## -omics, environment and phenotype
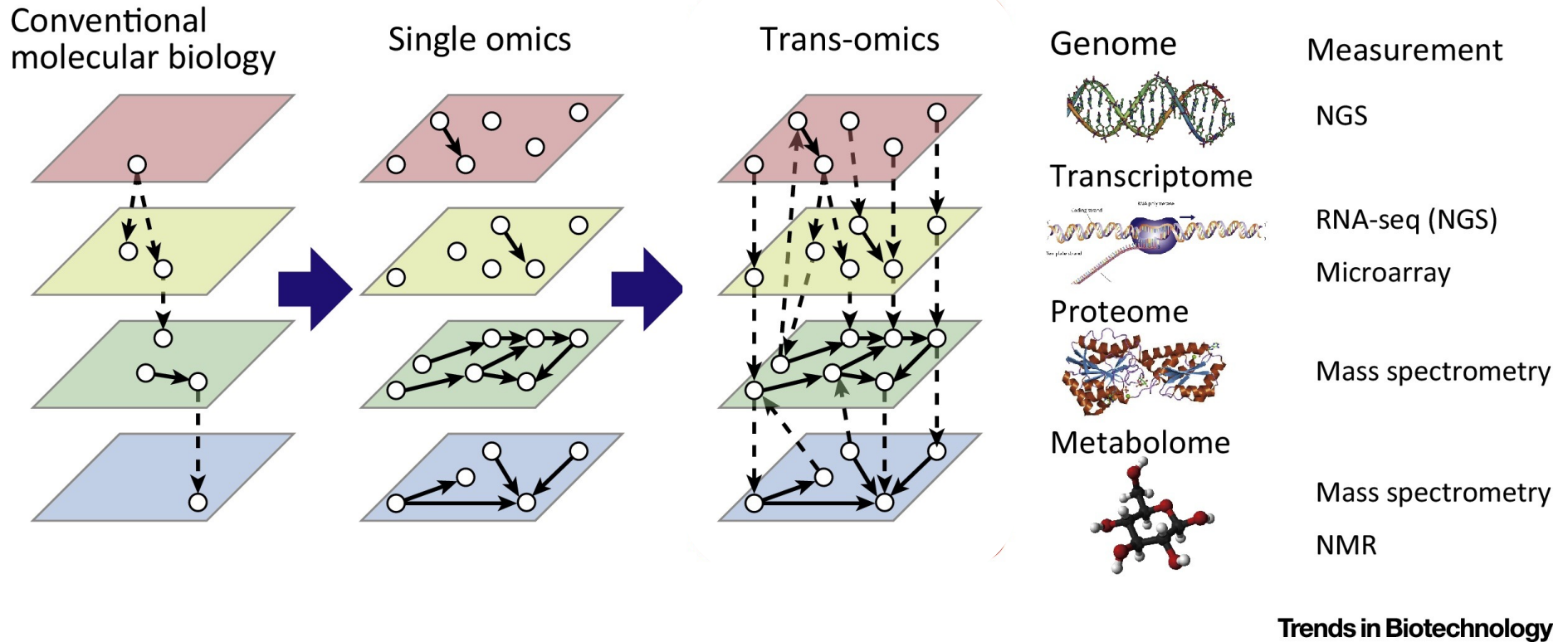


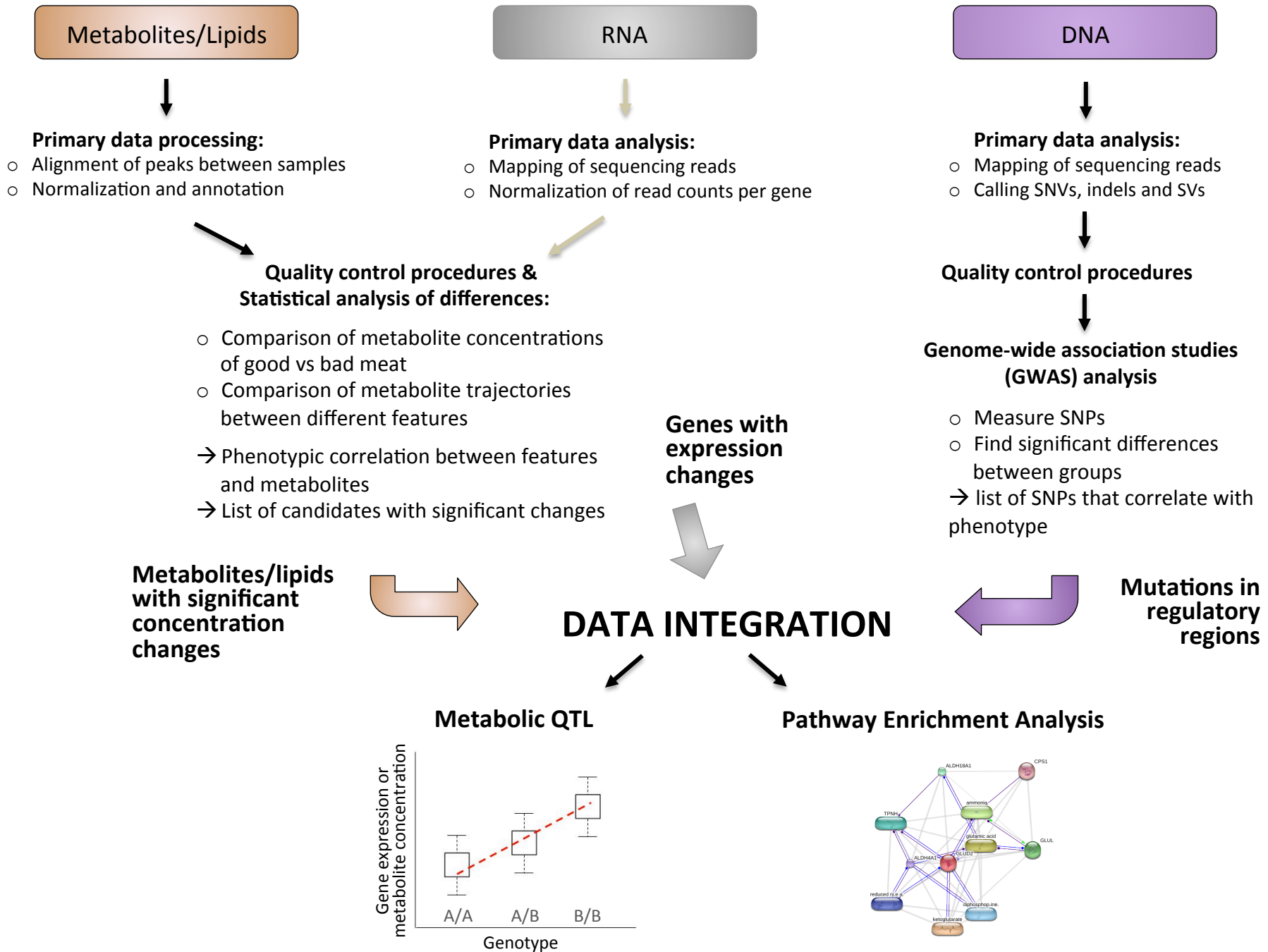DNA tells what is possibly…

…RNA what is probably….

…proteins and metabolites what actually happens

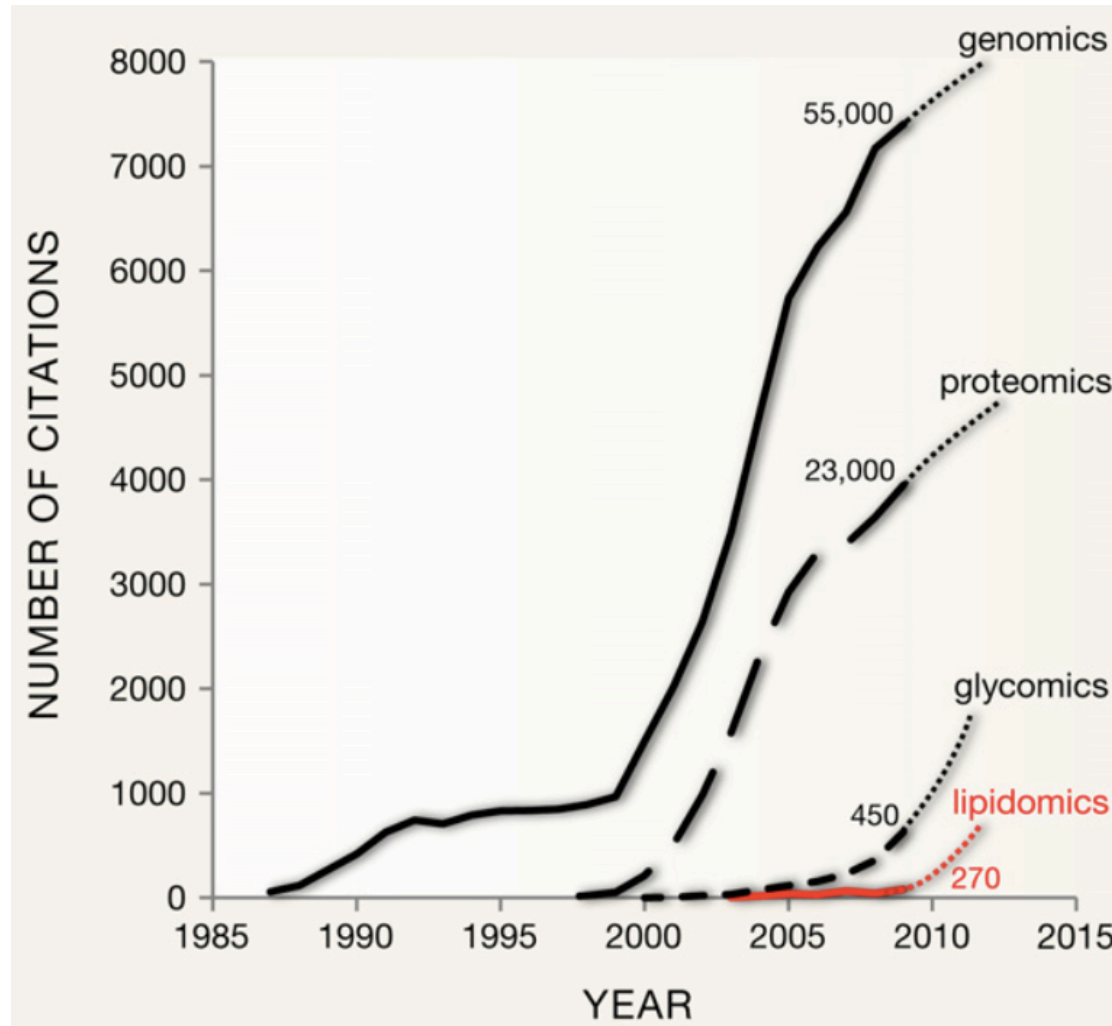Krumsiek 2016

# Systems Biology
## Technologies for different -omics layers



Yugi, Katsuyuki et al.
Trends in Biotechnology , Volume 34 ,
Issue 4 , 276 - 290

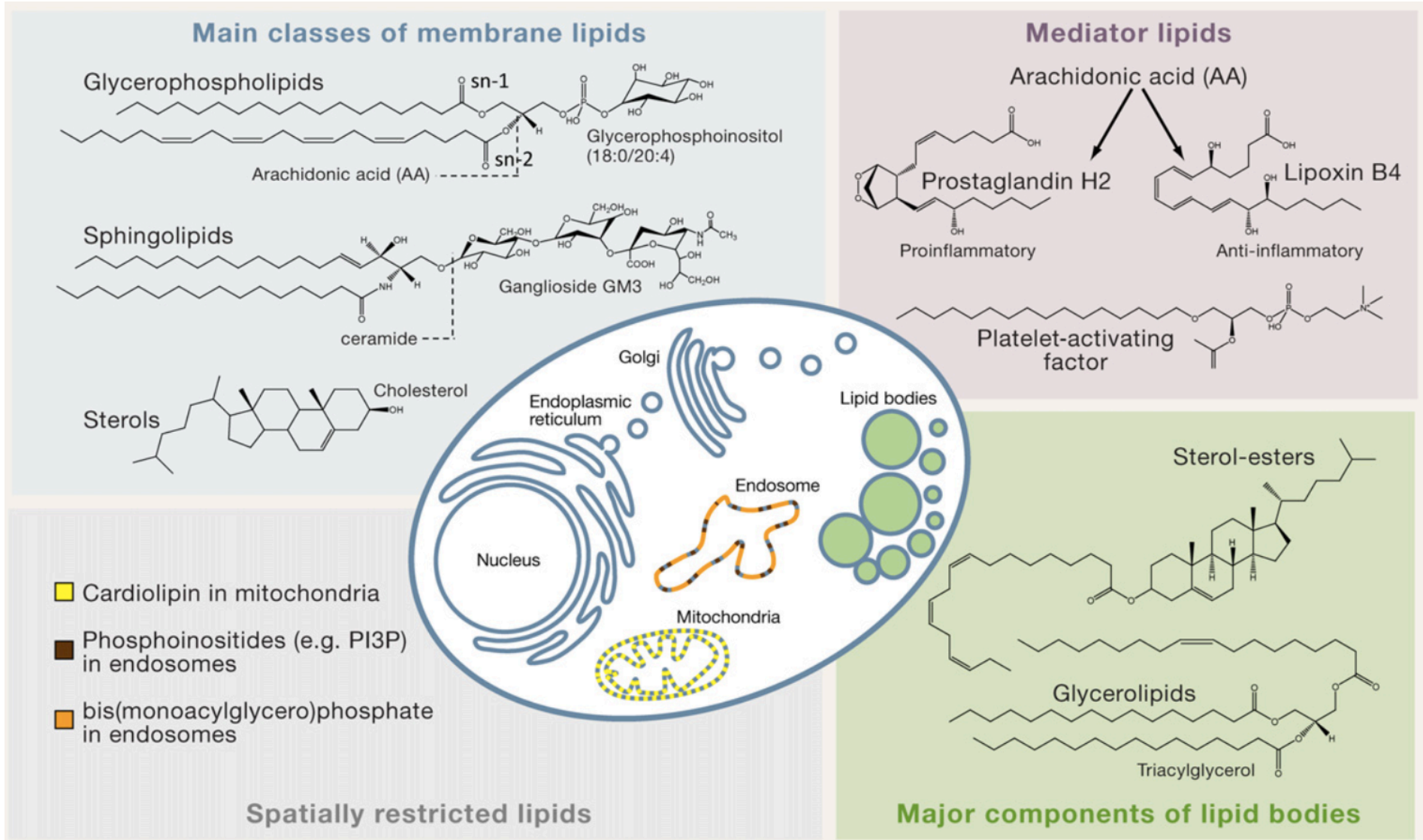# Lipidomics is a young, emerging field

Wenk 2010 – Cell 143(6):888-95
Lipidomics: New Tools and Applications

# Lipids
## Cellular Compartments of Common Biological Lipids

Wenk 2010 – Cell 143(6):888-95
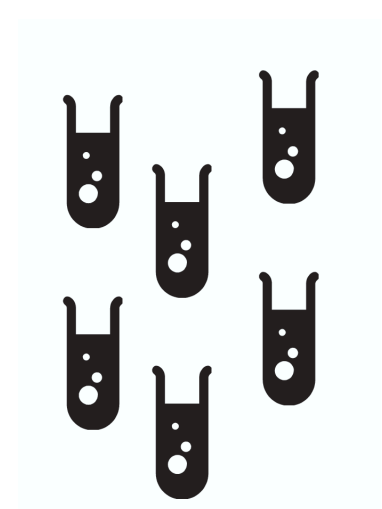Lipidomics: New Tools and Applications

# Study designs
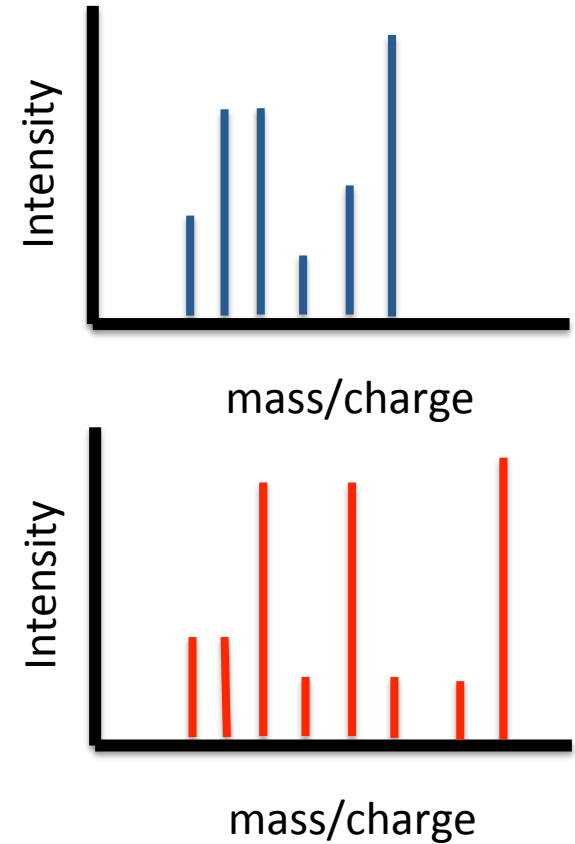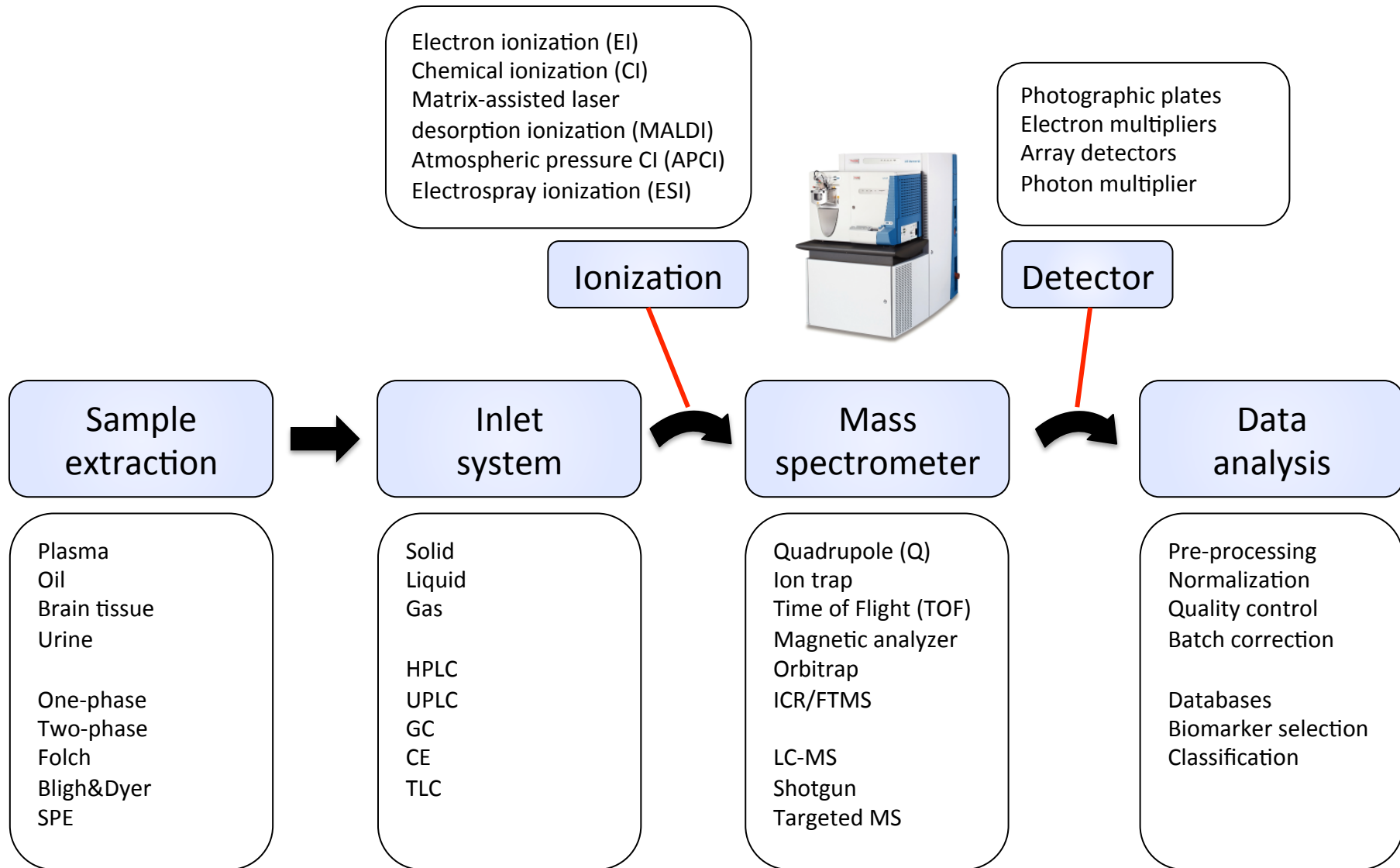
Tissue collection

Extraction of metabolites

Mass spectrometry analysis

# Mass spectrometry
## Workflow and variety

Electron ionization (EI)
Chemical ionization (CI)
Matrix-assisted laser
desorption ionization (MALDI)
Atmospheric pressure CI (APCI)
Electrospray ionization (ESI)

Photographic plates
Electron multipliers
Array detectors
Photon multiplier

**Ionization**

**Detector**

**Sample extraction** → **Inlet system** → **Mass spectrometer** → **Data analysis**

Plasma
Oil
Brain tissue
Urine

One-phase
Two-phase
Folch
Bligh&Dyer
SPE

Solid
Liquid
Gas

HPLC
UPLC
GC
CE
TLC

Quadrupole (Q)
Ion trap
Time of Flight (TOF)
Magnetic analyzer
Orbitrap
ICR/FTMS

LC-MS
Shotgun
Targeted MS

Pre-processing
Normalization
Quality control
Batch correction

Databases
Biomarker selection
Classification

# The data

# How does LC-MS data look like?
Chromatogram vs Spectrum

**Mass chromatogram**



Retention Time [min]

**Mass spectrum**



Mass [m/z]

# Data representation

# How does LC-MS data look like?
## Zooming in.....
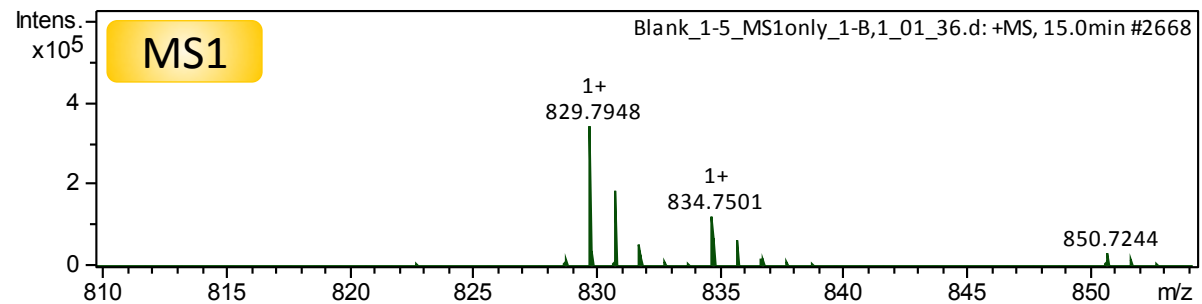
$$TAG_{(}15{:}0/18{:}1{-}d7/15{:}0)$$

$C_{51}H_{89}D_7O_6$
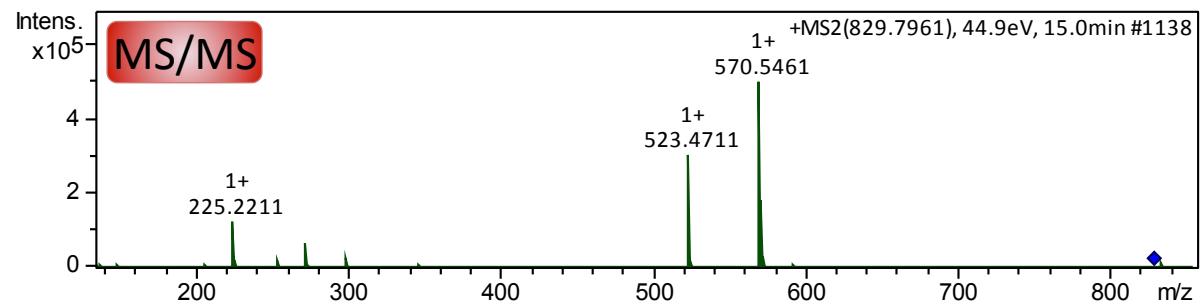
Neutral Mass: 811.76465

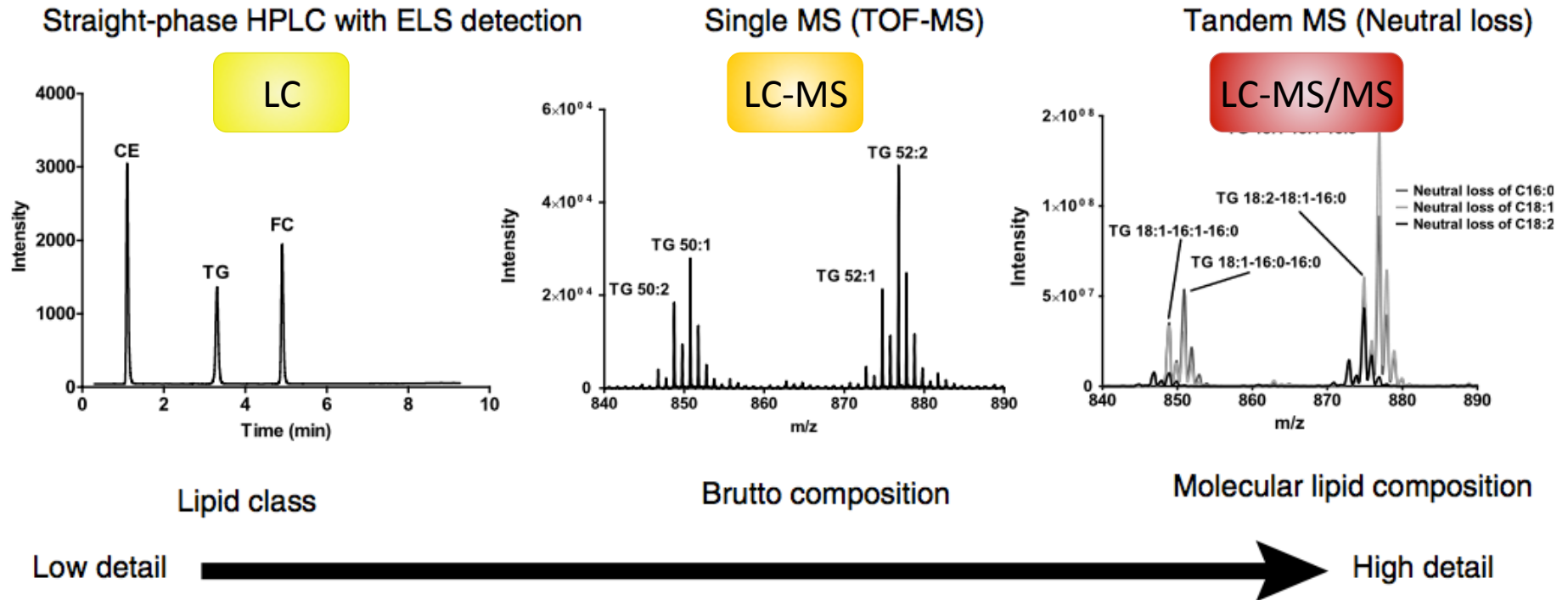**Chromatogram**



**Mass spectrum**



**Fragment mass spectrum**

# How do you get from data to compound?
## Why tandem MS (MS/MS)?



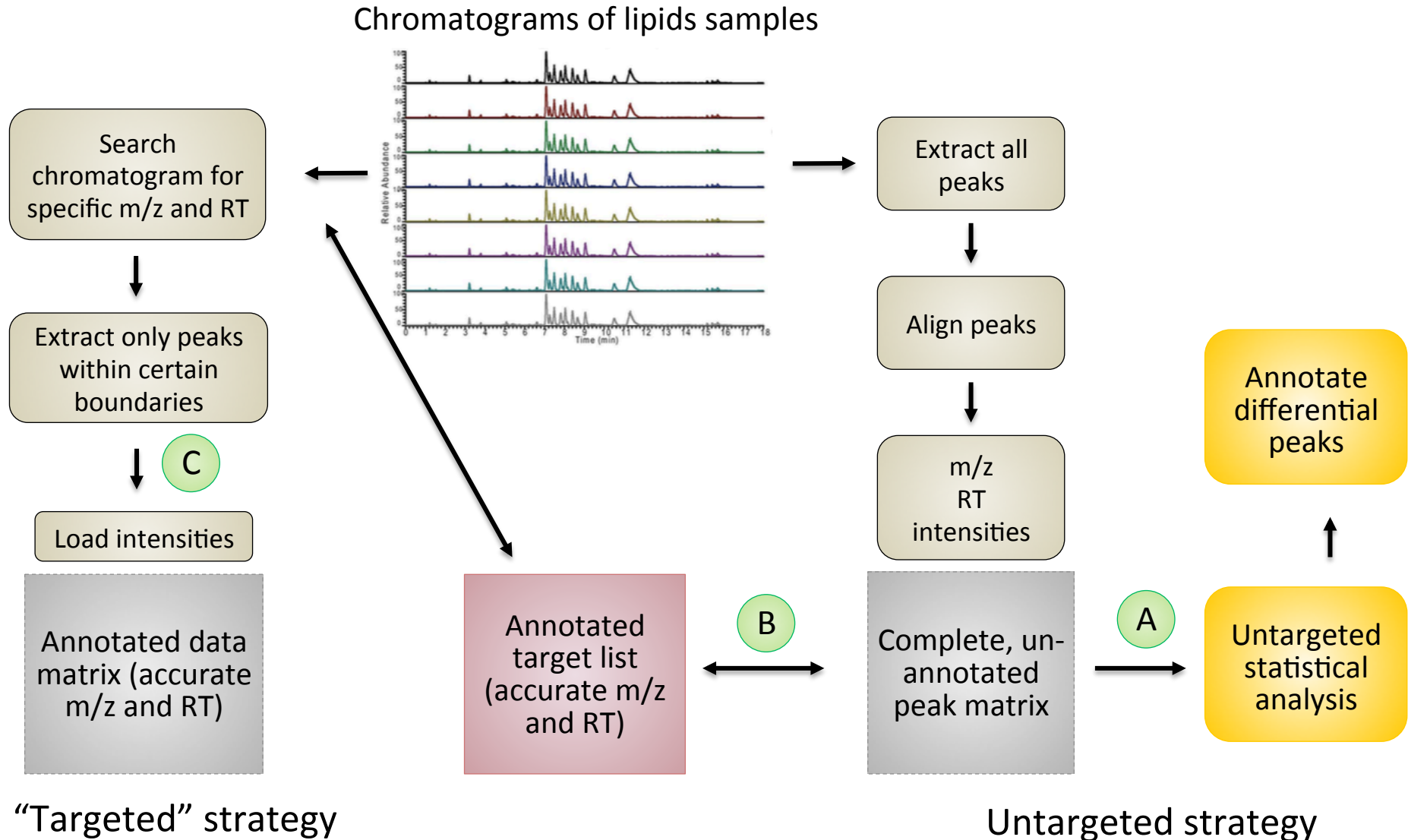- **LC-MS** allows for elucidation of molecular mass and most of the times brutto composition, but little structural information

- **LC-MS/MS** allows for the 1) detection of structurally informative fragment ions, and 2) the confirmation of ambiguous annotation of lipid species
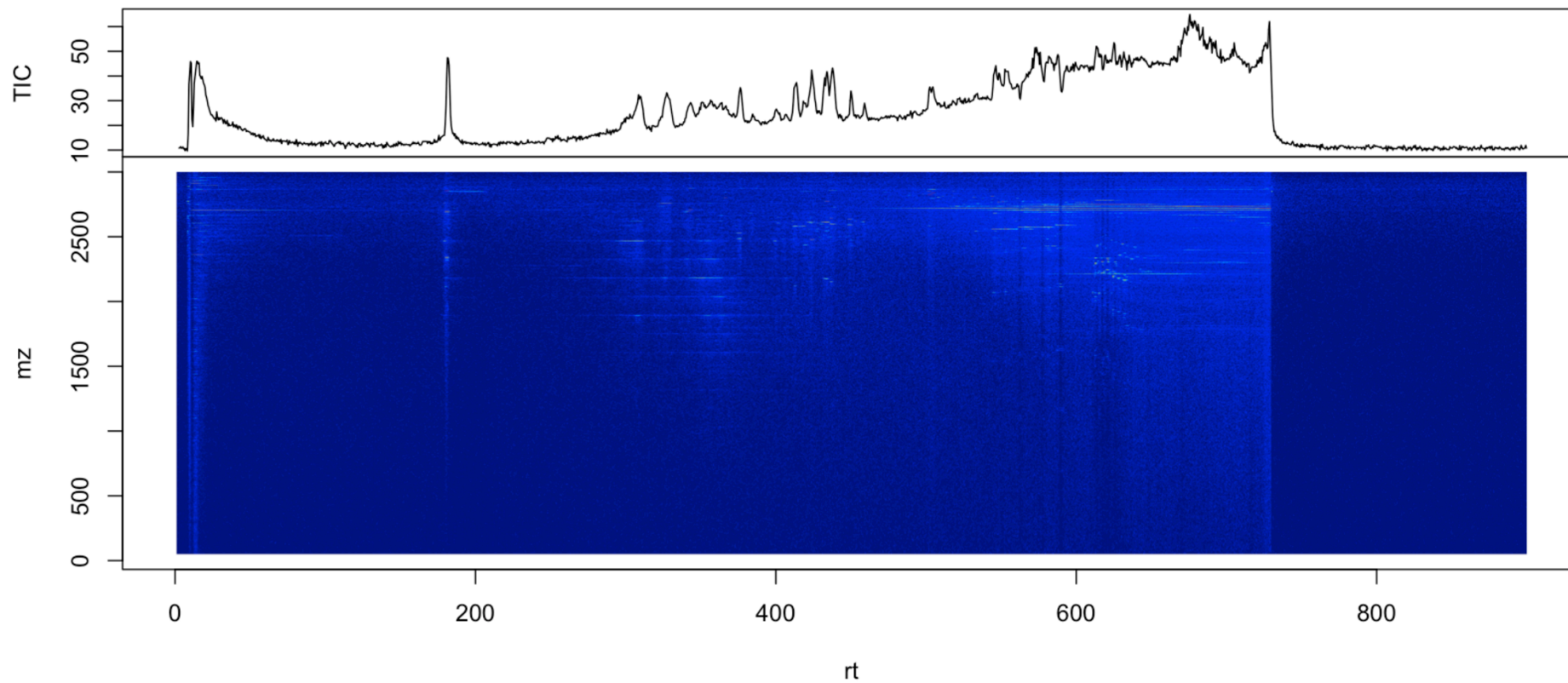
Ekroos, K. (2013). Lipidomics - Technologies and Applications. (K. Ekroos, Ed.)
(pp. 1–345). Wiley-VCH Verlag GmbH & Co. KGaA.

# How do you get from data to compound?
## LC-MS strategies (MS1)

Chromatograms of lipids samples



Search chromatogram for specific m/z and RT

Extract all peaks

Extract only peaks within certain boundaries

Align peaks

**C**

Load intensities

m/z RT intensities

Annotate differential peaks

Annotated data matrix (accurate m/z and RT)

Annotated target list (accurate m/z and RT)

**B**

Complete, un-annotated peak matrix

**A**

Untargeted statistical analysis

"Targeted" strategy

Untargeted strategy

# Data Processing

# Data representation

# LC/MS: Extracted Ion Chromatogram



**Peak:**
- **m/z 607.2925**
- **retention time 11.49 min**
- **intensity = max or area?**

Expected peak shape is the same for a certain compound => max ~ area. So no difference (hopefully) in case if you don't compare compounds between each other, just the same compound between samples.
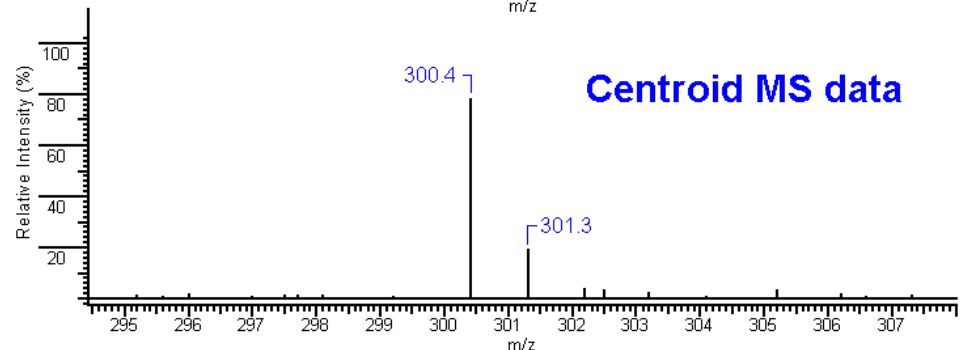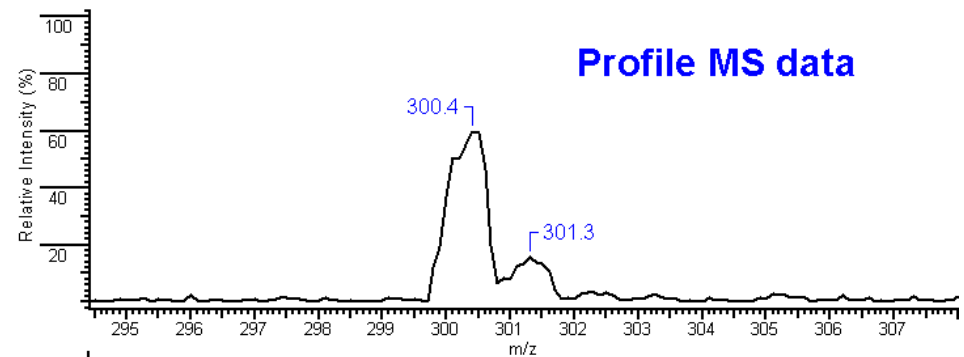
# Data types
## Profile vs Centroid

***Profile data*** (aka continuous) – intensity records for all the range of mz and retention time (RT).
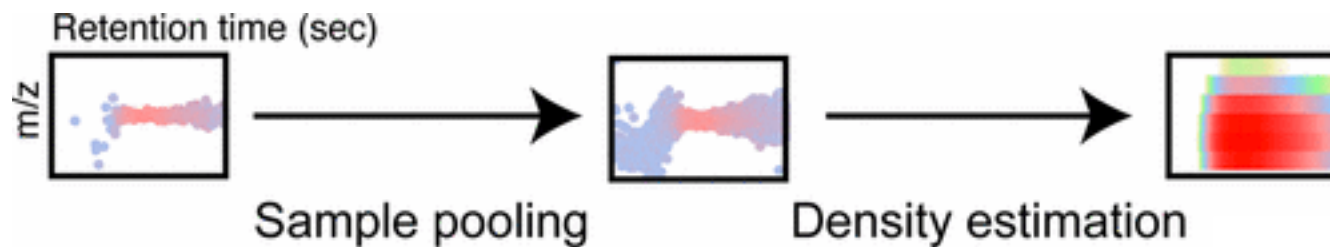
***Centroided data*** – only local maximums are detected and saved.

- Pros profile:
    - More options for peak detection, better detection
    - Less ambiguous => less false positive values

- Cons profile:
    - Big data volume
    - Slow conversion and analysis

# Two paradigms

1.  Peak picking then alignment (Do peak picking for each sample separately)

2.  Alignment and peak picking  (Do peak picking on each sample simultaneously)
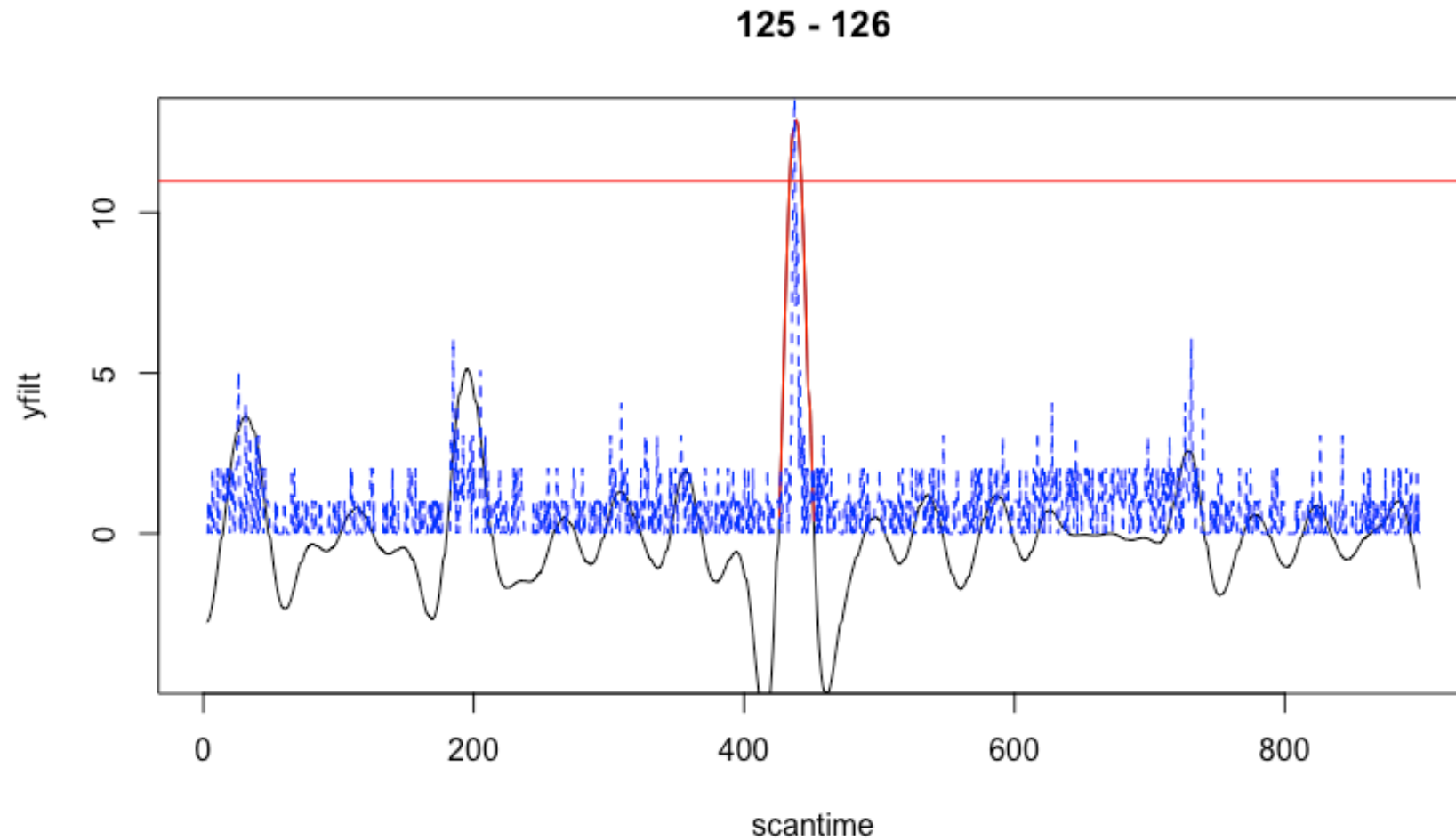
**Peak Picking / Peak Detection**
Methods

- In literature there are a lot of different peak picking algorithm. But no best solution, only better solutions.

- *Know your data!*

- Gaussian model peak width – standard

Smoothing, baseline correction may be applied, not for all methods.
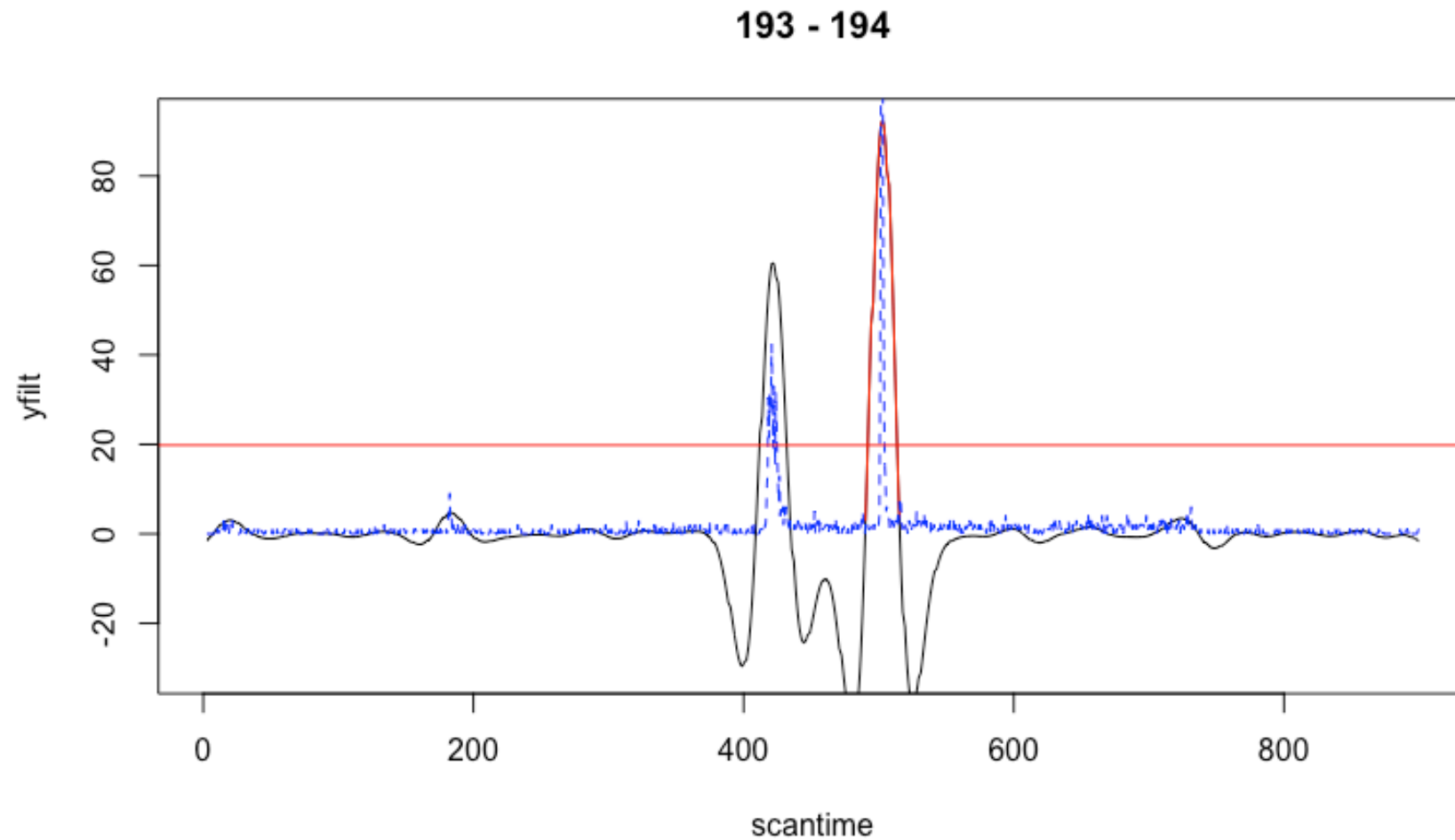
Peak picking is a crucial step of analysis. The main question: how to choose method and parameters?

- Tryout => tradition

- Repeating for others

- Attempt to define objective metrics of peak-picking quality and build a parameter selection based on their maximization: Brodsky L. et al. (2010) Evaluation of Peak Picking Quality in LC–MS Metabolomics Data. *Anal. Chem.*
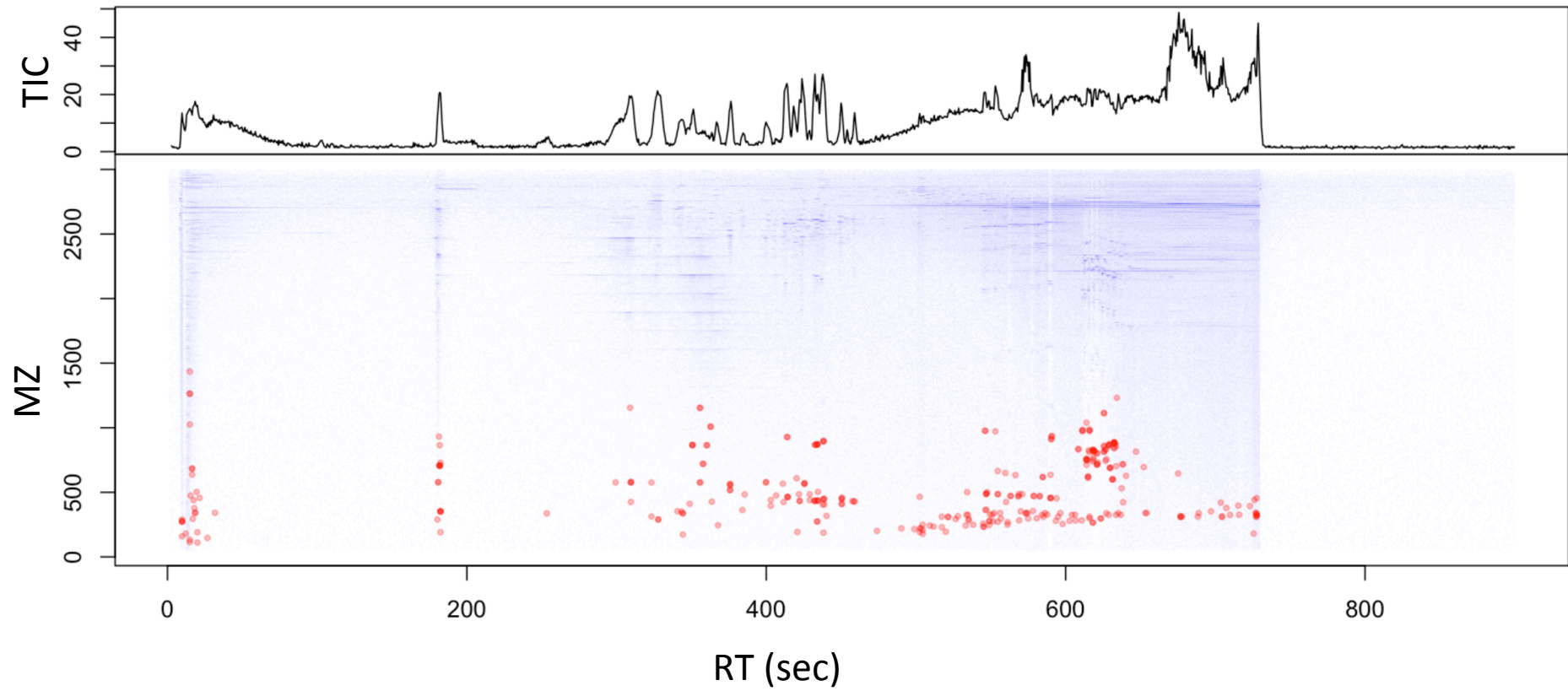
# Matched filter (gaussian model) – noisy example



125 - 126

# Matched filter (gaussian model) – good example



**193 - 194**

# Peak Picking / Peak Detection

# Align different samples

- Construct the data matrix

- Combine the single samples

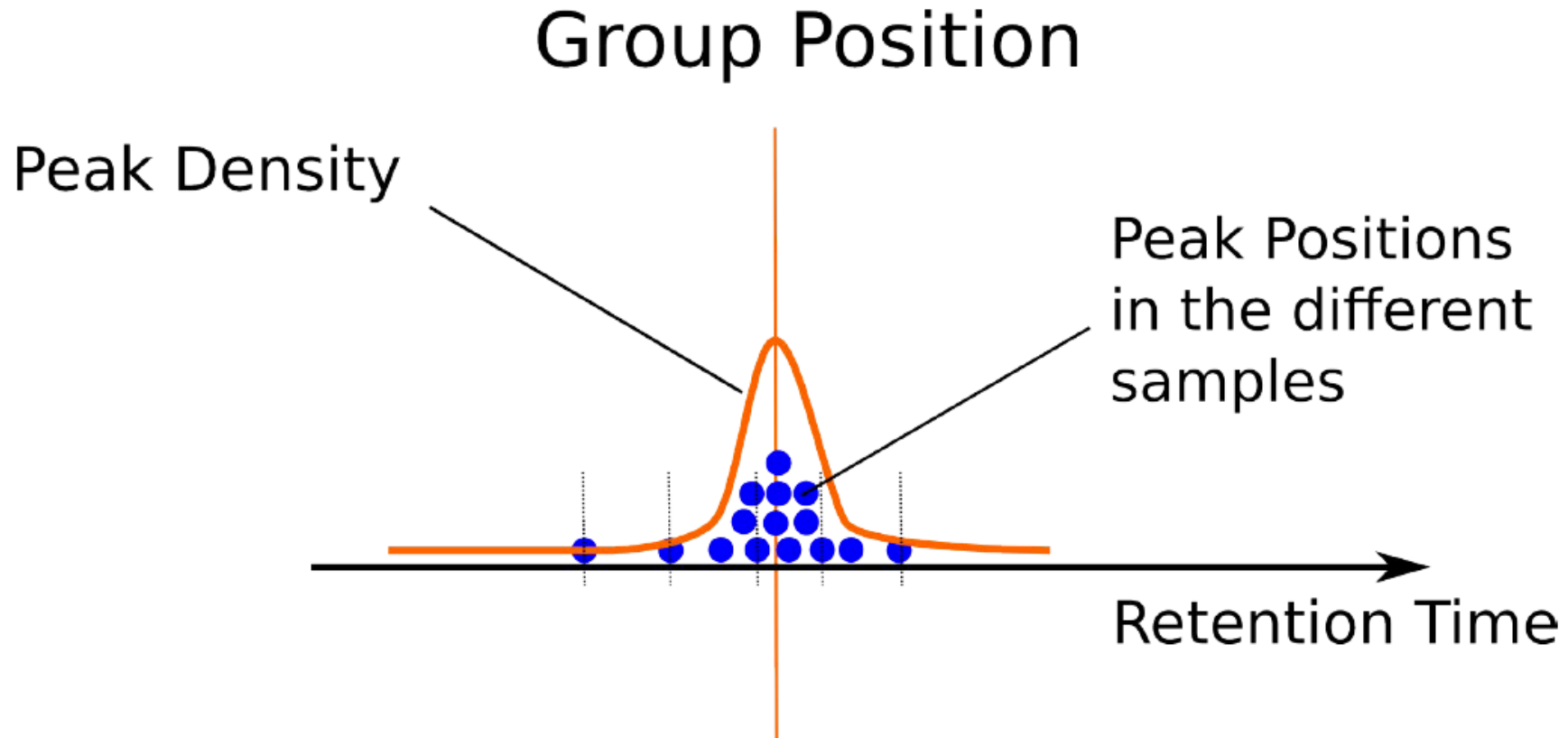- With the ultimate goal of correcting for retention time shifts



**We HAVE to compare the right variable across the samples**

# Density Based Grouping

# Peak alignment



Franceschi et al., J. Chemom. (2012)

# Rt correction

- Need "hook" groups.
- Ideally each sample is represented by one feature in a "hook" group.
- Correct the hooks and interpolate elsewhere
- Unfortunately You can have more or fewer features per group.

# Rt correction



Retention Time Deviation vs. Retention Time

- Can be performed before peak picking (chromatogram alignment)

- Linear or polynomial or whatever correction

- May afford to exclude any ambiguous peaks

- You could run it iteratively till RT deviation is less than your window for peak grouping

# Feature Detection Evaluation

- Compare mzMine, XCMS
- Gold standard via
  - Technical replicates
  - Democracy
- Evaluation via
  - Dilution series
  - Mix of complex samples
- F-Measure:  sum of
  - Precision (TP/(TP+FP))
  - Recall or sensitivity (TP/P)



Tautenhahn, Böttcher, Neumann. High Sensitive Feature Detection For High Resolution LC-MS. BMC Bioinformatics (2008)

# Peak Filtering

Remove peaks from data table based on:

- Number of missing values for a peak

- Max/mean/median intensity (total or within groups of replicates)

- Variability in intensity – coefficient of variance, standard deviation, interquartile range, etc. (total or within groups of replicates)

- …

# Missing values
## Discrimination

**NAs could be:**

- real zero/low concentration

- mispicked/misaligned peaks (in general feature is detected correctly)

- incorrectly detected feature

**Considerations:**

- Total % of NAs for a feature

- presence in replication groups

- amplitude, variability

- ...

**Missing values**
Treatment

- **Unreliable features:**
  - Remove

- **True zeros:**
  - Look at raw specters data
  - Generate random baseline-level noise

- **False zeros:**
  - Replace by mean/median/etc. for this feature
  - Replace by mean/median/etc. for this feature & replication group
  - PCA-based (BPCA, PPCA, …), KNN-based imputation methods

See:
Stacklies, W. et al. (2007). pcaMethods — a bioconductor package providing PCA methods for incomplete data. *Bioinformatics.*

**Normalization**
Methods

- Not changing intensity distribution – all intensities in one sample have the same normalization factor:
  - by biomass
  - by a single internal standard
  - by mean/median/sum intensity of features in this sample
  - probabilistic quotient normalization (PQN)
  - ...
- Changing intensity distribution – each feature in each sample has it's own norm factor, i.e.:
  - by multiple internal standards (i.e. NOMIS)
  - quantile normalization – "stretching" distributions of all samples to make them similar
  - ...
- General assumption for normalization is that most of the compounds are not affected. Is that true? For different treatment? For different species? For different tissues? Does it matter if we have no choice? :/

# Normalization
## Probabilistic Quotient Normalization (PQN)



- Reference spectrum could be a single "golden" spectra or an average/median spectrum of control group

- Divide each spectrum by a reference spectrum (feature by feature)

- Plot distribution of ratios (quotients)

- Find median

- This is your scaling factor

Dieterle, F. et al. (2006). **Probabilistic Quotient Normalization as Robust Method to Account for Dilution of Complex Biological Mixtures. Application in 1H NMR Metabonomics**. *Anal Chem.*

# Centering and Scaling
Applied to features across the samples

**Nature of MS data:**
- Features are extremely different in amplitude
- Heteroscedasticity – biological (induced and uninduced) and technical variance are higher for features with high intensity

**Scaling:**
- Equalizes contributions of features to separation in multivariate space
- Makes features comparable (i.e. for looking at time profile)

**Types of scaling:**
- Range scaling – by [max – min] – sensitive to outliers; undesirable
- Auto-scaling – by standard deviation (SD) – data loose dimensionality
- Pareto-scaling – by root of SD – features with higher intensity decrease more

**Centering** is subtracting mean/median from all the values:
- Necessary for some methods like PCA and makes no sense for others like fold change

# Transformation

Certain function applies to all the values in a data table.

- Log-transformation

- General logarithmic transformation (glog) – approximately log for high values and linear close to zero

- Cube-root transformation

Why?

- Transformation has a scaling-like effect making features more comparable.

- Log/glog-transformation helps to reveal multiplicative relations between features.

**Annotation**

Retention indexing / Retention projection

RT is extremely variable. Idea of ***retention indexing:***

save an exemplary LC as a "scale" for the future and then align all the times by this database.

**Limitations:**

*   limited number of tested compounds – extrapolate several compounds to a class?

*   interactions between compounds => RT could depend on a sample composition – databases of complex mixtures?

*   only certain LC system/conditions – retention projection? (see the next slide)

All additional experiments =>  time, money

**Annotation**

Databases

Annotation could be manual or with more or less automatic tools coupled with databases:

- Commercial – really?

- Open source

- In-house:
  - works for you, specified for your needs, possible to include retention indexing
  - but costs additional work, money, time

Fragment MS/MS (or GC/MS) databases:

- Experimental
  - specific: instrument, ionization parameters, etc.
- In-silico (e.g. LipidBlast)
  - theoretical, but wide coverage

# Problems

## Experimental

1. Batch effect (48 per run)
2. Platform-based effect
3. Poor correspondence between experiments
4. Concentration estimation

## Data Analysis

1. Annotation (low percent of annotated compounds ~20-40%)
2. No golden software standard
3. Technical effects
4. Poor alignment of samples

# Acknowledgments

Questions?