

# Part I: Normalization & Summarization

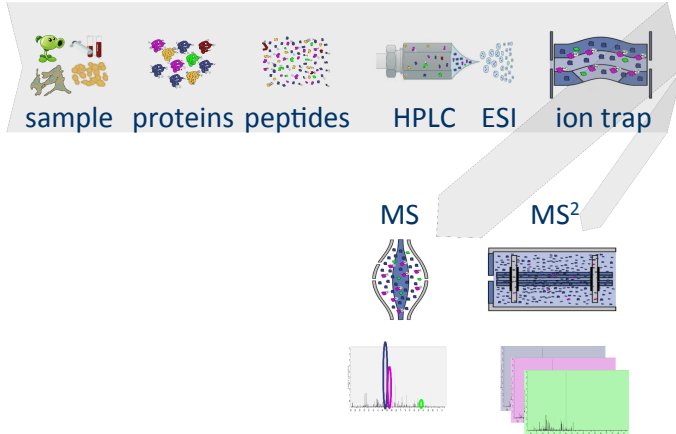
Lieven Clement

Proteomics Data Analysis 2018, Gulbenkian Institute, May 28 -June 1  
2018.

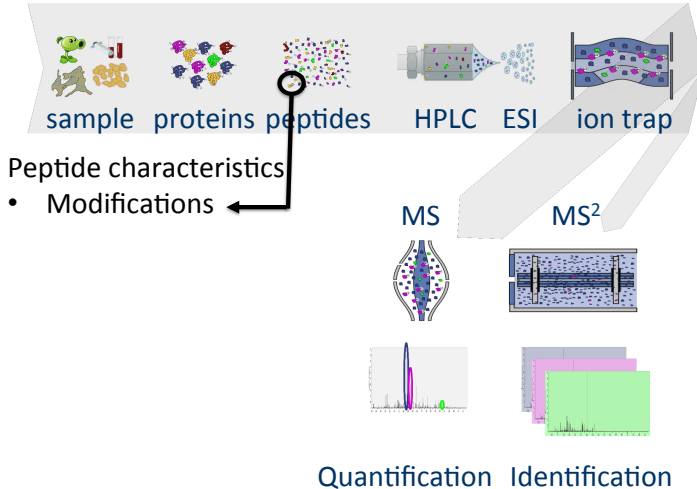
# Outline

- ① Introduction
  - ① Label free MS based Quantitative Proteomics Workflow and Challenges
- ② Preprocessing
  - ① Filtering
  - ② Log transformation
  - ③ Normalization
  - ④ Summarization

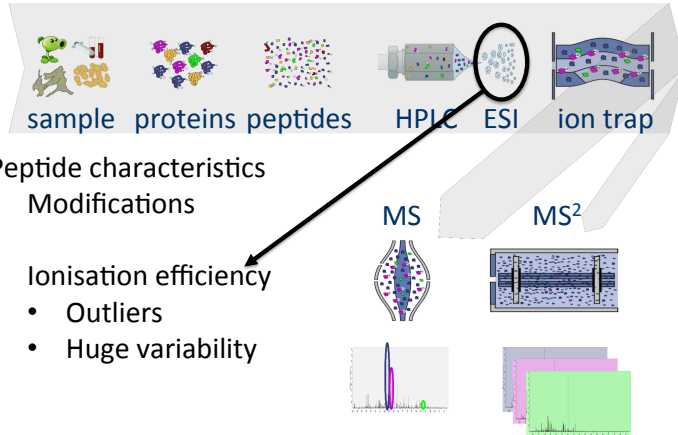
# Challenges in Label Free Quantitative Proteomics



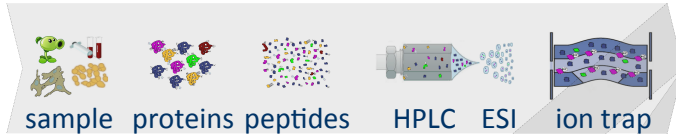
# Challenges in Label Free Quantitative Proteomics



# Challenges in Label Free Quantitative Proteomics

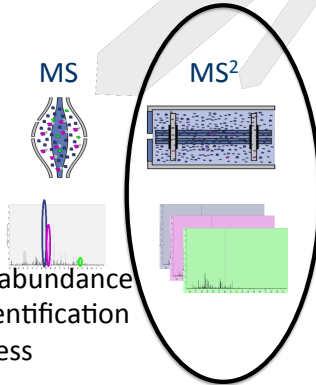


# Challenges in Label Free Quantitative Proteomics

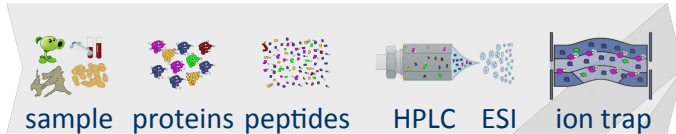


## Peptide characteristics

- Modifications
- Ionisation efficiency
  - Outliers
  - Huge variability
- MS<sup>2</sup> selection on peptide abundance
  - Context dependent Identification
  - Non-random missingness

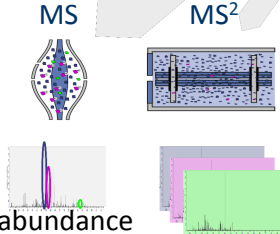


# Challenges in Label Free Quantitative Proteomics



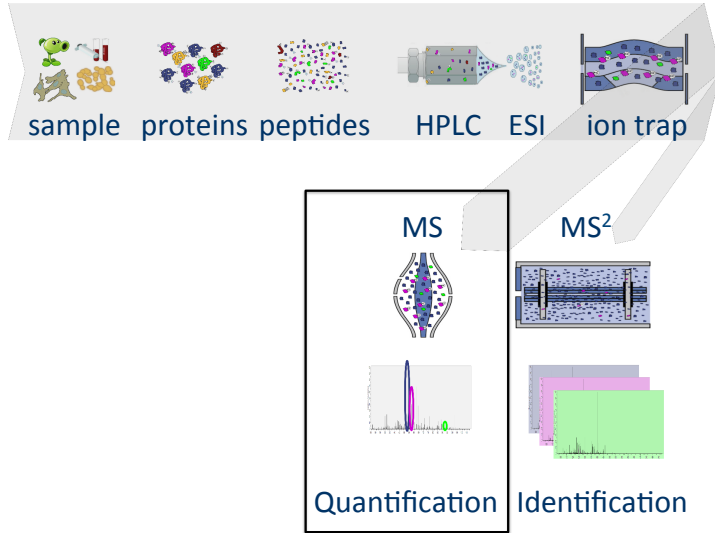
## Peptide characteristics

- Modifications
- Ionisation efficiency
  - Outliers
  - Huge variability
- MS<sup>2</sup> selection on peptide abundance
  - Context dependent Identification
  - Non-random missingness



**Unbalanced peptides identifications across samples and messy data**

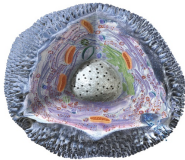
# Challenges in Label Free MS-based Quantitative proteomics



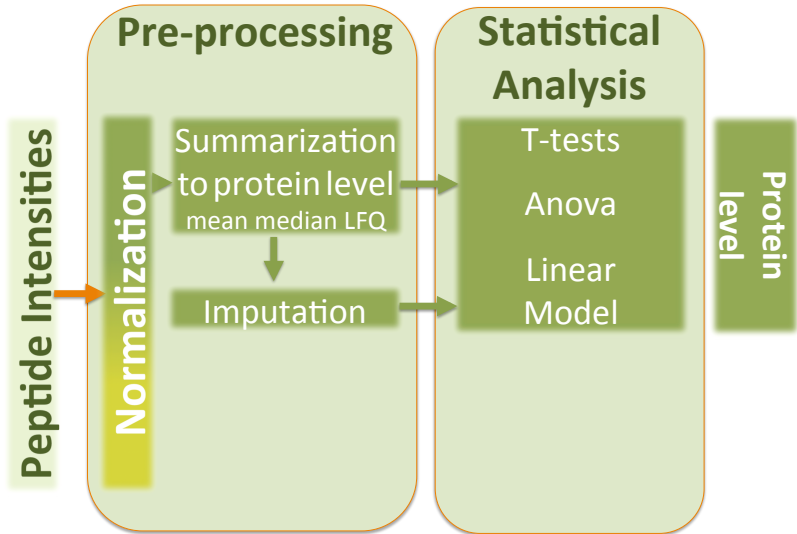


# Challenges in Label Free MS-based Quatitative proteomics

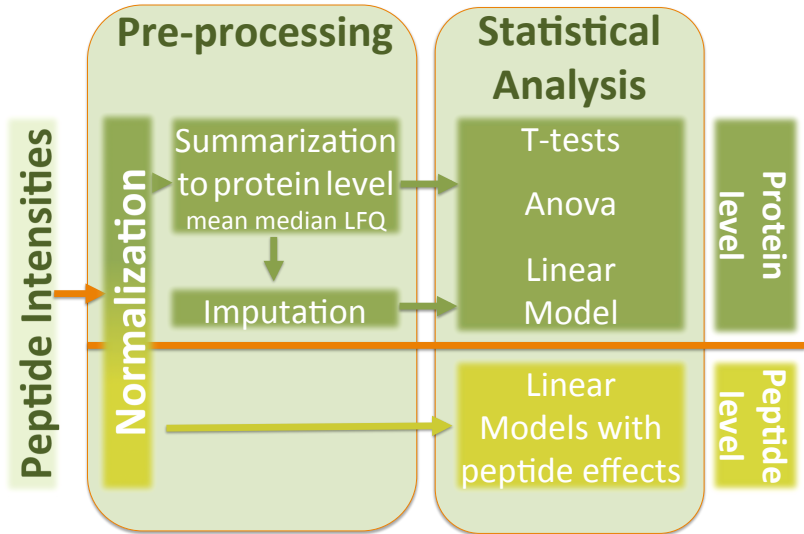
We need information on protein level!



# Label-free Quantitative Proteomics Data Analysis Pipelines



# Label-free Quantitative Proteomics Data Analysis Pipelines



# CPTAC Spike-in Study

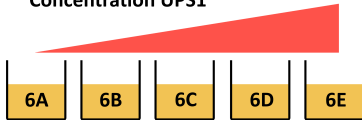
Digested  
UPS1 protein mix



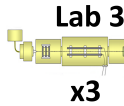
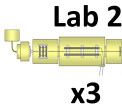
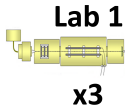
Digested  
yeast proteins



Concentration UPS1



5 spike-in concentrations: 6A to 6E



- Same trypsin-digested yeast proteome background in each sample
  - Trypsin-digested Sigma UPS1 standard: 48 different human proteins spiked in at 5 different concentrations (treatment A-E)
  - Samples repeatedly run on different instruments in different labs
  - After MaxQuant search with match between runs option
    - 41% of all proteins are quantified in all samples
    - 6.6% of all peptides are quantified in all samples
- **vast amount of missingness**

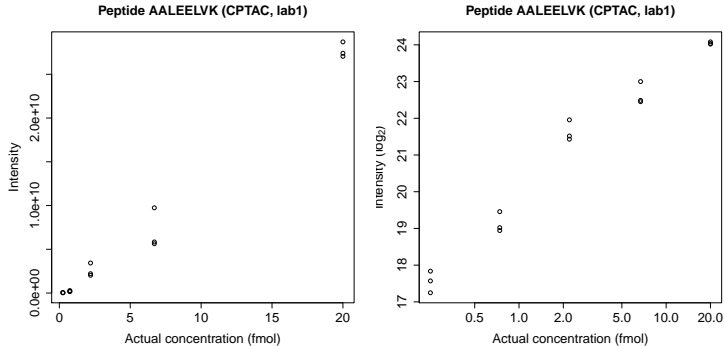
# Preprocessing

- Typical preprocessing steps
  - ① Filtering
  - ② Log-transformation
  - ③ Normalization
  - ④ (Summarization)
- Many methods exist

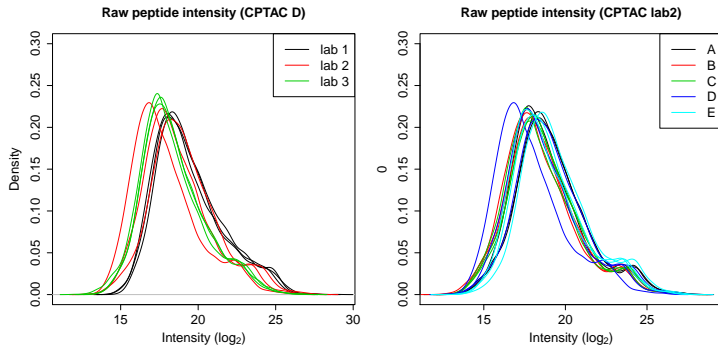
# Filtering

- Reverse sequences
  - Only identified by modification site (only modified peptides detected)
  - Razor peptides: non-unique peptides assigned to the protein group with the most other peptides
  - Contaminants
  - Peptides few identifications
  - Proteins that are only identified with one or a few peptides
- 
- Filtering does not induce bias if the criterion is independent from the downstream data analysis!

# Log-transformation



Variability more equal upon log transformation: often multiplicative error structure of intensity-based read-outs



Even in very clean synthetic dataset (same background, only 48 UPS proteins can be different) the marginal peptide intensity distribution across samples can be quite distinct

- Considerable effects between and within labs for replicate samples
  - Considerable effects between samples with different spike-in concentration
- Normalization is needed

# Mean or median?

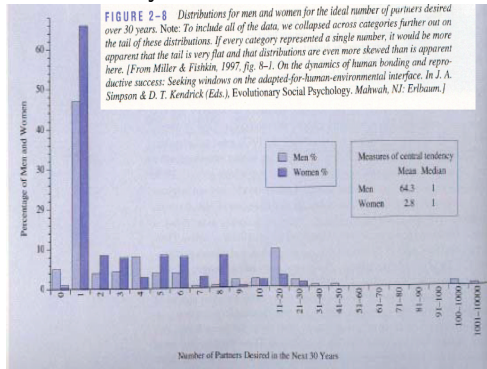
- Over a period of 30 years males desire to have on average 64.3 partners and females 2.8. (Miller and Fishkin, 1997)

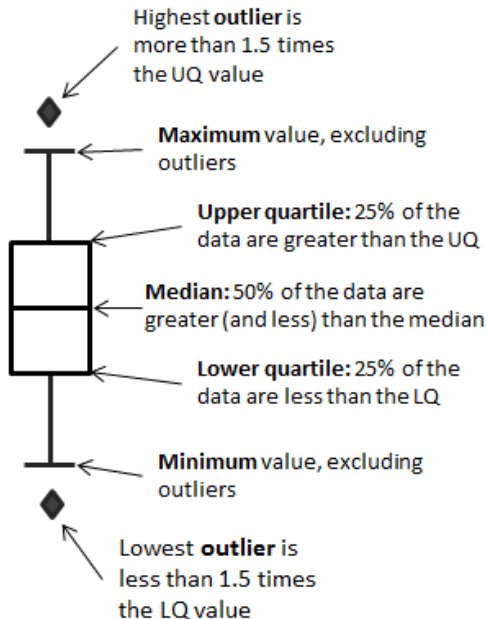
# Mean or median?

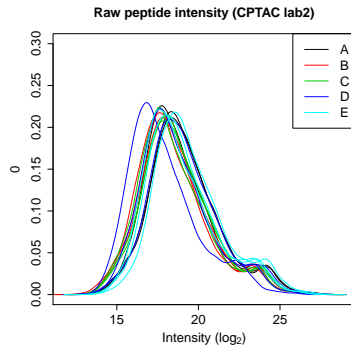
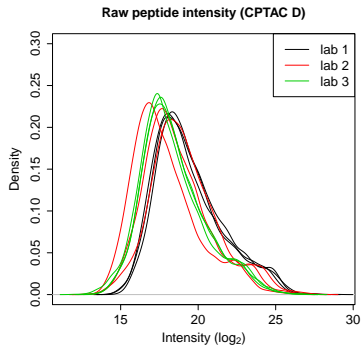
- Over a period of 30 years males desire to have on average 64.3 partners and females 2.8. (Miller and Fishkin, 1997)
- Over a period of 30 years males, is the median of the number of desired partners is 1 for both males and females. (Miller and Fishkin, 1997)

# Mean or median?

## Mean is very sensitive to outliers!

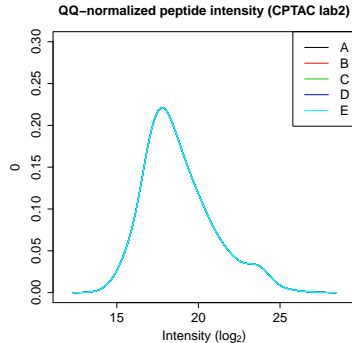
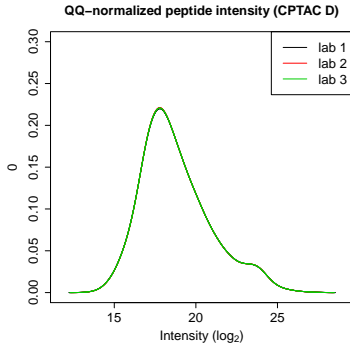






Even in very clean synthetic dataset (same background, only 48 UPS proteins can be different) the marginal peptide intensity distribution across samples can be quite distinct

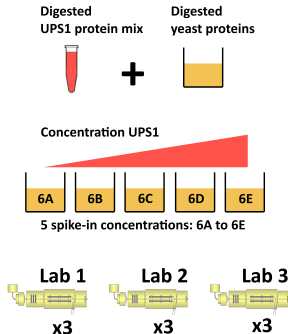
- Considerable effects between and within labs for replicate samples
  - Considerable effects between samples with different spike-in concentration
- Normalization is needed



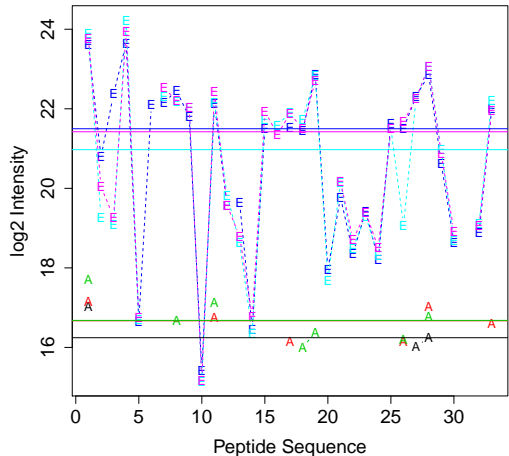
Even in very clean synthetic dataset (same background, only 48 UPS proteins can be different) the marginal peptide intensity distribution across samples can be quite distinct

- Considerable effects between and within labs for replicate samples
  - Considerable effects between samples with different spike-in concentration
- Normalization is needed, e.g. **quantile normalization**

# Summarization

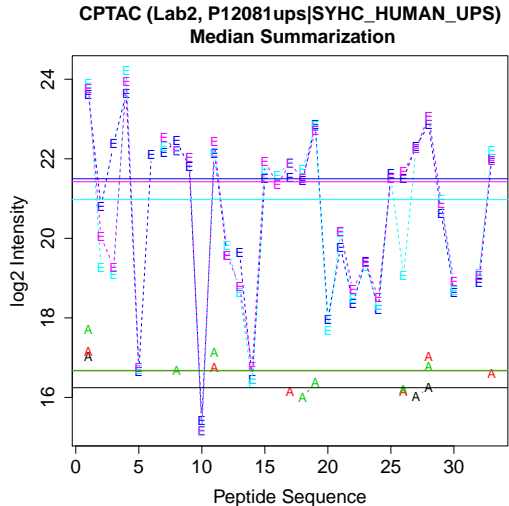


**CPTAC (Lab2, P12081ups|SYHC\_HUMAN\_UPS)**  
Median Summarization



# Summarization

- Strong peptide effect
- Unbalanced peptide identification
- Summarization bias
- Different precision of protein level summaries



# MaxLFQ summarization

a

&gt;P63208

MPSIKLQSSDGEIFEVDVEIAKQSVTIKTMLEDLGMDDEGDD  
 DPVPLPNVNAAILKKVIQWCTHHKDDPPPPEDDENKEKRTDD  
 IPVVDQEFLEKVDQGTFLFELILAANYLDIKGLLDVTCKTVANM  
 IKGTPEEIRKTFNIKNDFTEEEEAQVRKENQWCEEK

b

Peptide species	Sequence	Charge	Mod.
P <sub>1</sub>	LQSSDGEIFEVDVEIAK	2	–
P <sub>2</sub>	LQSSDGEIFEVDVEIAK	3	–
P <sub>3</sub>	RTDDIPVVDQEFLEK	2	–
P <sub>4</sub>	TVANMIK	2	–
P <sub>5</sub>	TVANMIK	2	Oxid.
P <sub>6</sub>	TPEEIRK	3	–
P <sub>7</sub>	NDFTEEEEAQVR	2	–

c

Sample	P <sub>1</sub>	P <sub>2</sub>	P <sub>3</sub>	P <sub>4</sub>	P <sub>5</sub>	P <sub>6</sub>	P <sub>7</sub>
A		+				+	
B		+	+			+	
C	+	+	+	+		+	+
D	+	+		+		+	+
E		+		+			+
F		+			+		

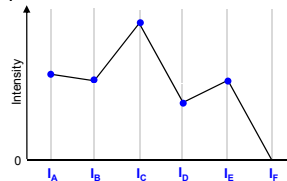
d

A						
B	$r_{BA}$					
C	$r_{CA}$	$r_{CB}$				
D	$r_{DA}$	$r_{DB}$	$r_{DC}$			
E	$r_{EA}$	$r_{EB}$	$r_{EC}$	$r_{ED}$		
F	$r_{FA}$	$r_{FB}$	$r_{FC}$	$r_{FD}$	$r_{FE}$	
	A	B	C	D	E	F

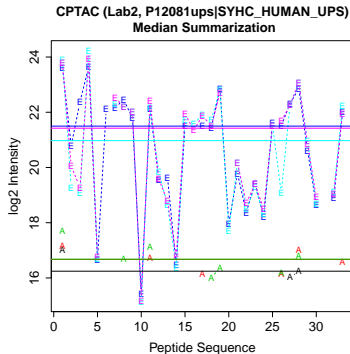
e

$r_{BA} = I_B / I_A$	$r_{CA} = I_C / I_A$	$r_{CB} = I_C / I_B$
$r_{DA} = I_D / I_A$	$r_{DB} = I_D / I_B$	$r_{DC} = I_D / I_C$
$r_{EC} = I_E / I_C$	$r_{ED} = I_E / I_D$	$I_F = 0$

f

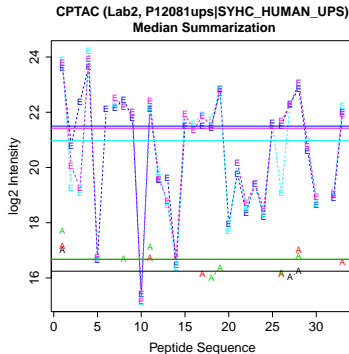


# Peptide based model



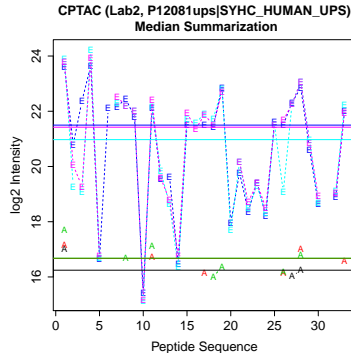
- 1 y: log2 transformed peptide intensities

# Peptide based model



- 1 y: log2 transformed peptide intensities
- 2 Protein by protein analysis of peptide level data with linear model

# Peptide based model



- ①  $y$ : log2 transformed peptide intensities
- ② Protein by protein analysis of peptide level data with linear model
 

peptide level

$y_{pept} \sim peptide$

+

protein level

$sample$