



Statistical mass spectrometry-based proteomics

Olga Vitek

www.stat.purdue.edu

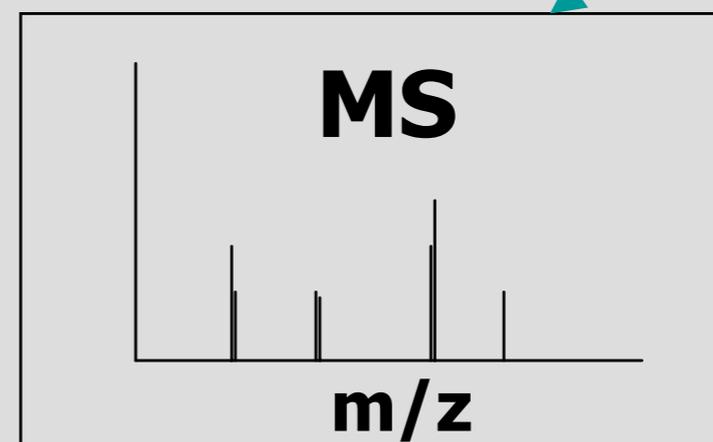
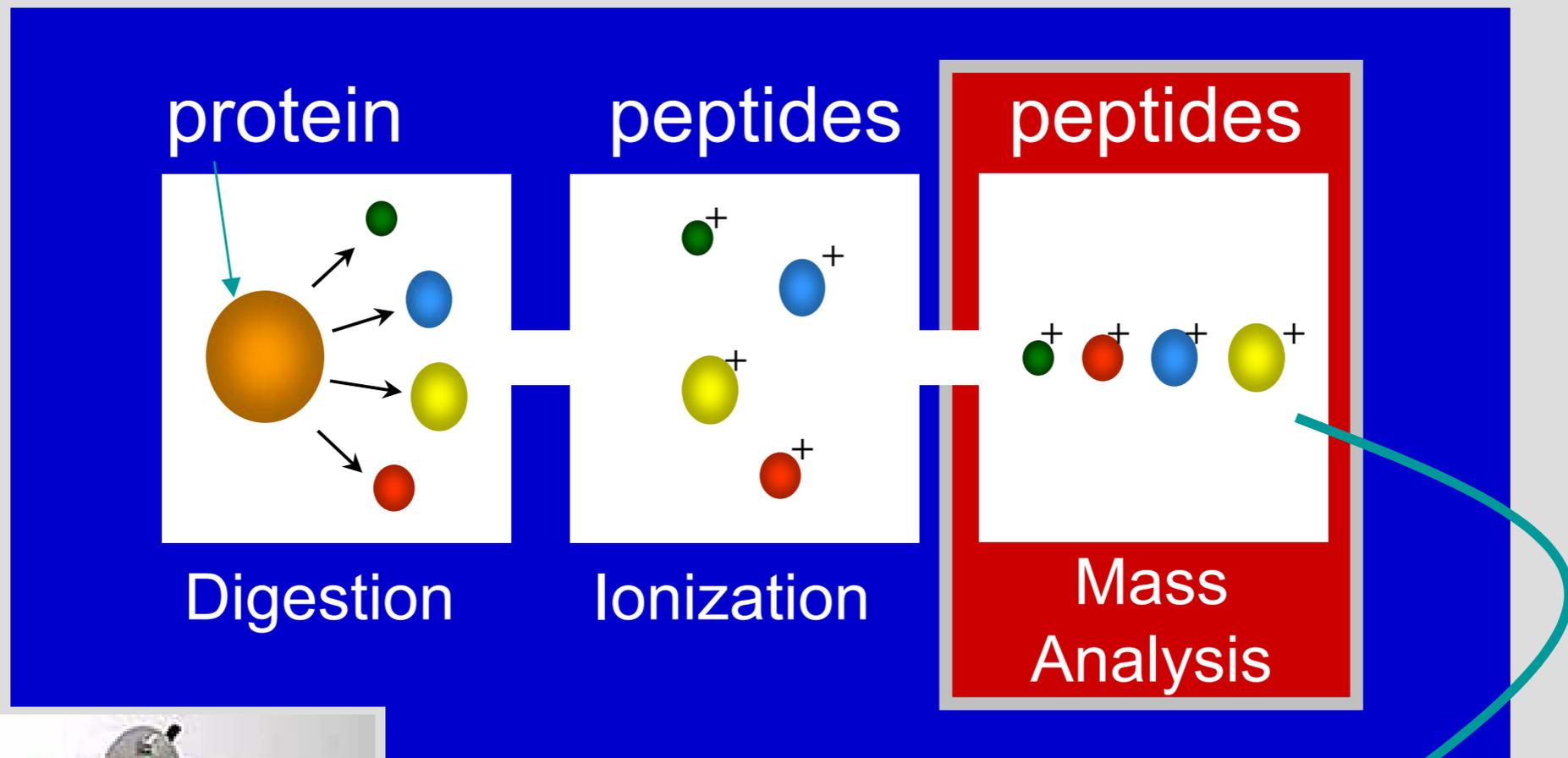
Outline

- What is proteomics?
 - ◆ Biological questions and technologies
- Protein quantification in label-free workflows
 - ◆ Joint analysis of multiple features and conditions
- Protein quantification in label-based workflows
 - ◆ Appropriately account for the labeling structure
- Mass spectrometry-based imaging
 - ◆ Account for the spacial heterogeneity of spectral data

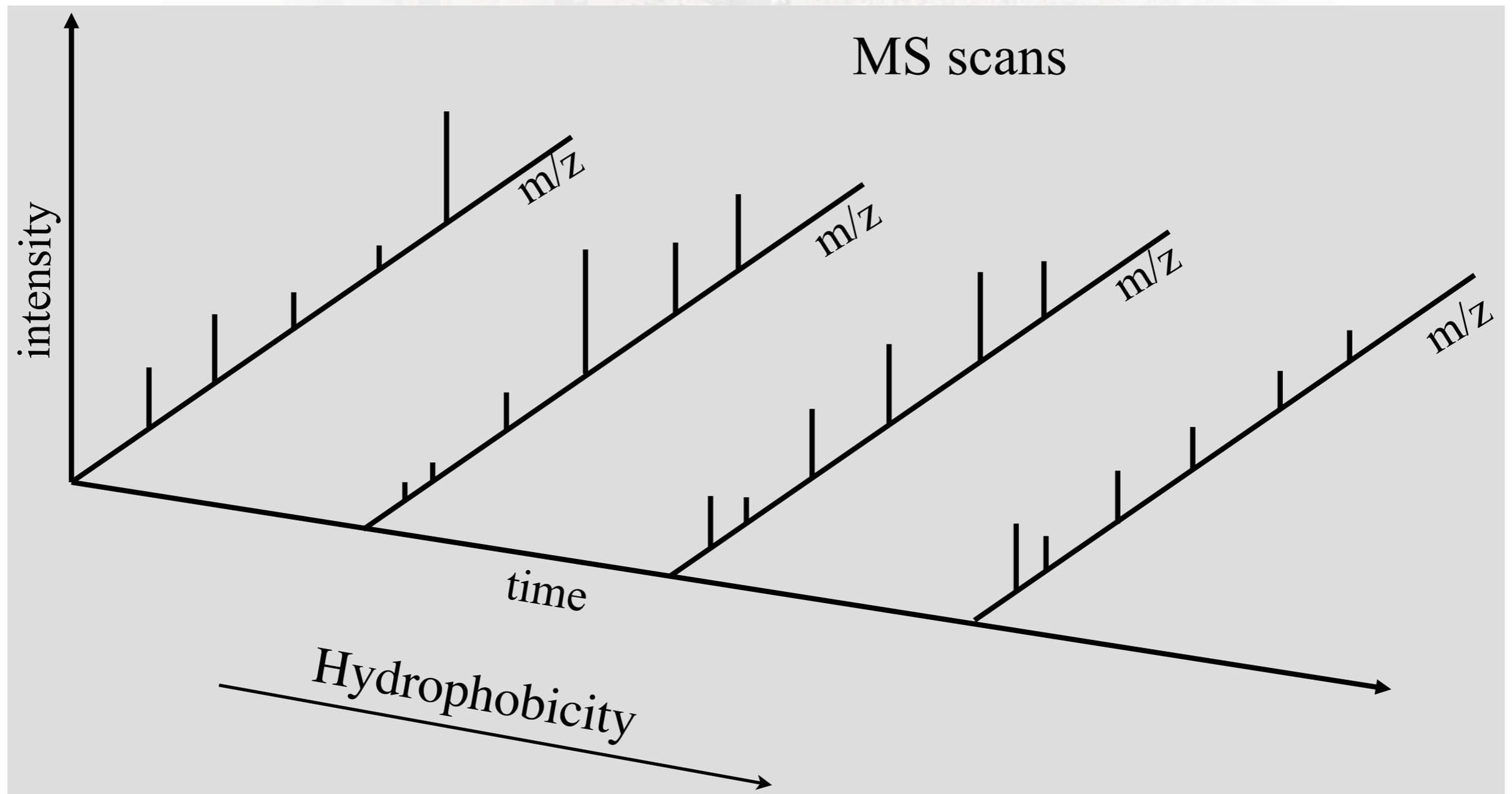
Goals of proteomics

- Proteomics: system-wide characterization of all proteins
 - ◆ Sequence, structure, localization, abundance, PTMs, interactions
- More challenging than gene expression
 - ◆ *Complexity*
 - Human genome: ~20,000 protein coding genes
 - Their translation+splicing+proteolysis: ~50,000–500,000 proteins
 - Somatic DNA rearrangements and PTM: ~10 million
 - ◆ *Dynamic range*
 - > 10 orders of magnitude in plasma
 - Unlike nucleotides, proteins cannot be amplified
- How to make progress?
 - ◆ Sample preparation, separation, sensitive instruments
 - ◆ Statistical experimental design and accurate analysis

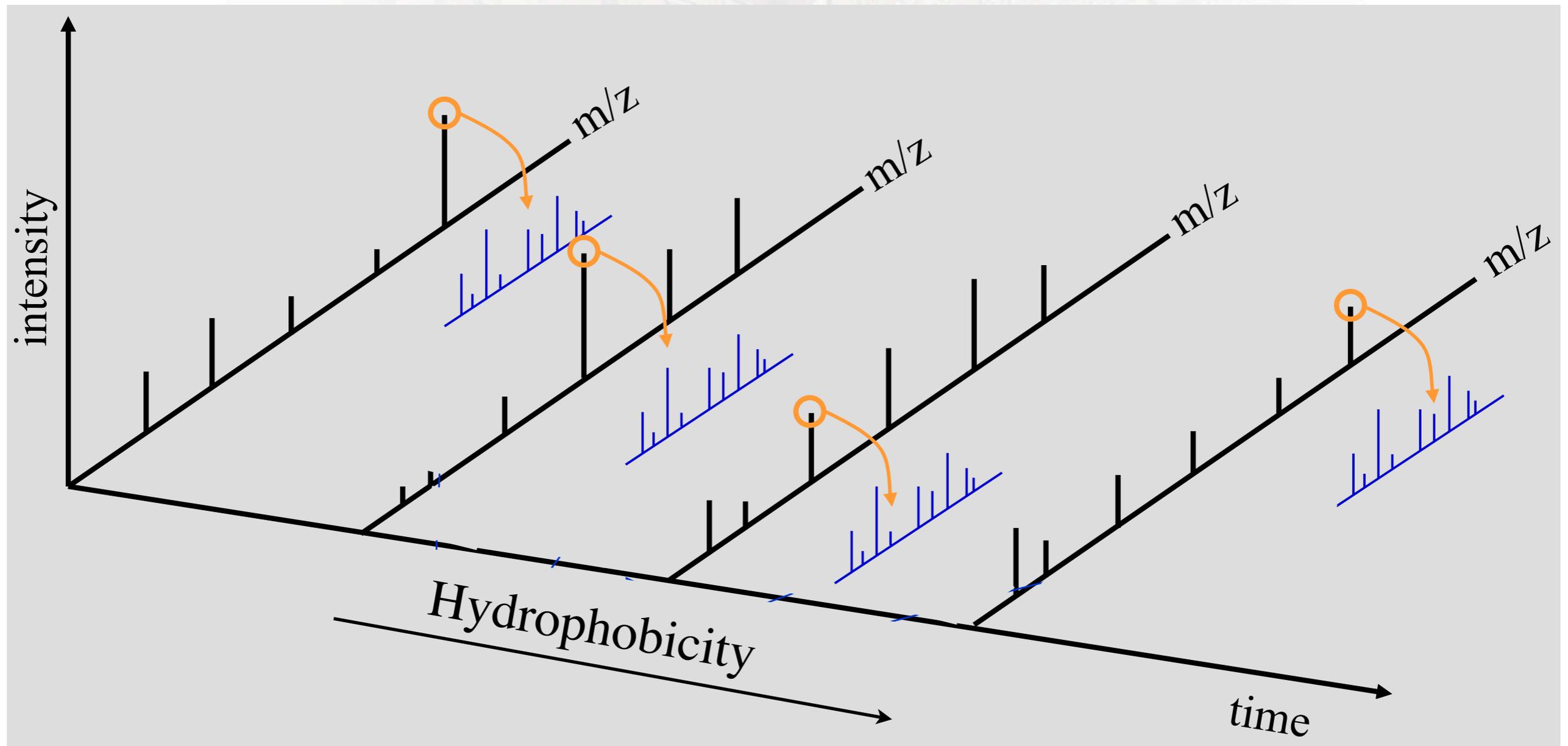
Acquisition of mass spectra



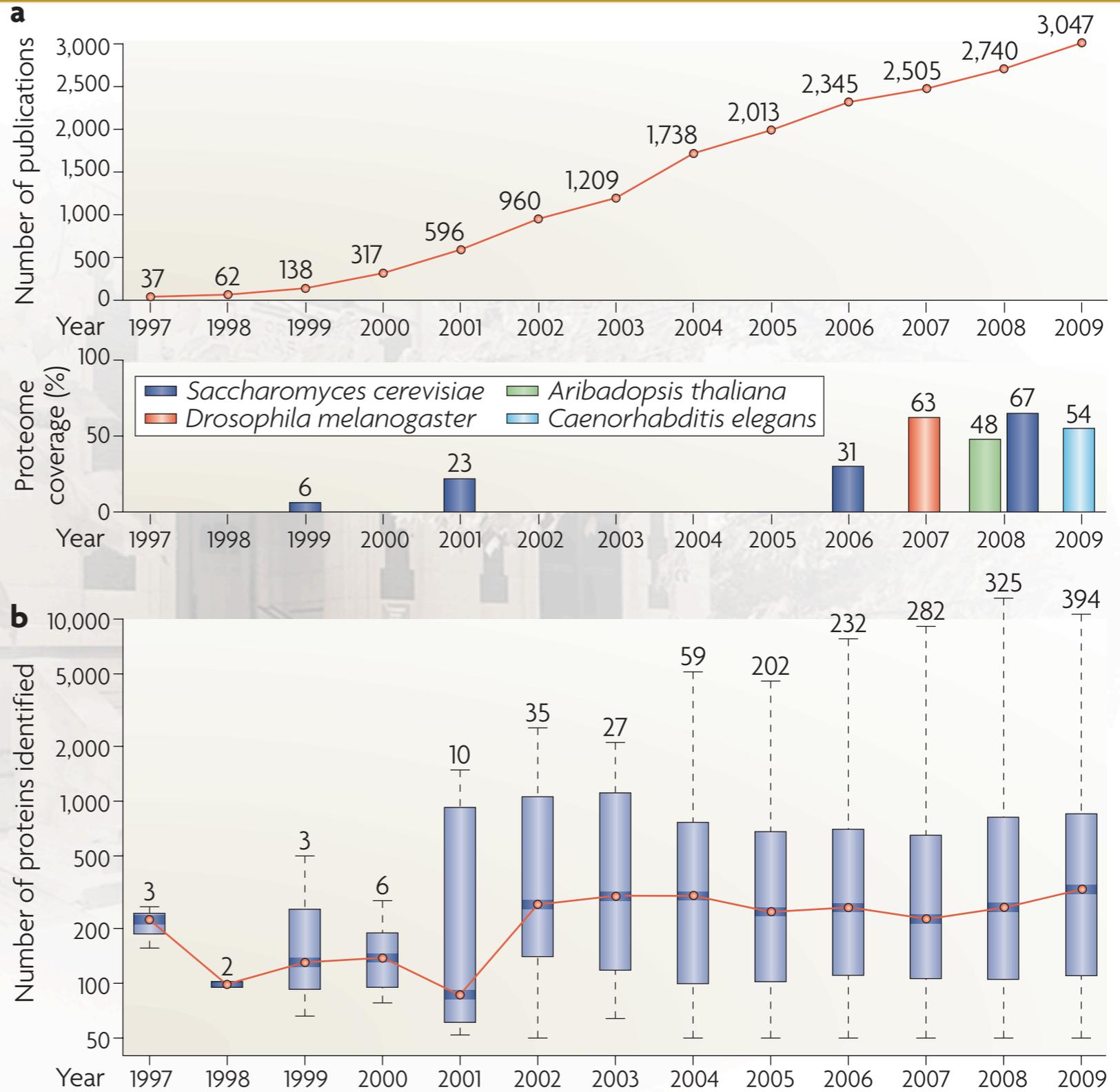
Liquid chromatography coupled with mass spectrometry (LC-MS)



Liquid chromatography coupled with mass spectrometry (LC-MS)



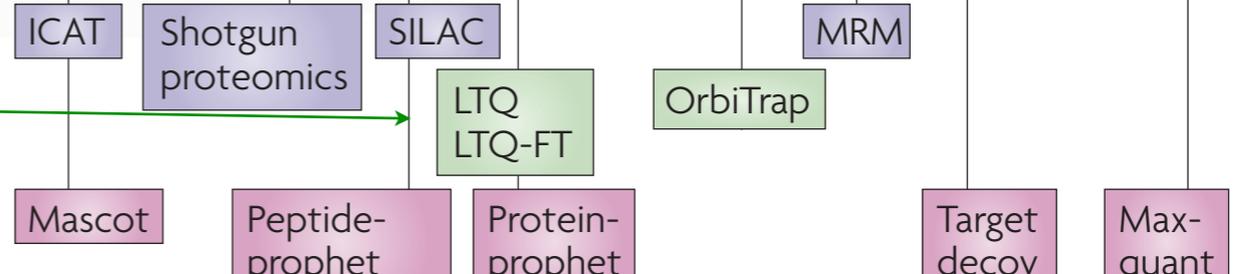
Proteomics is increasingly comprehensive



Assay technology

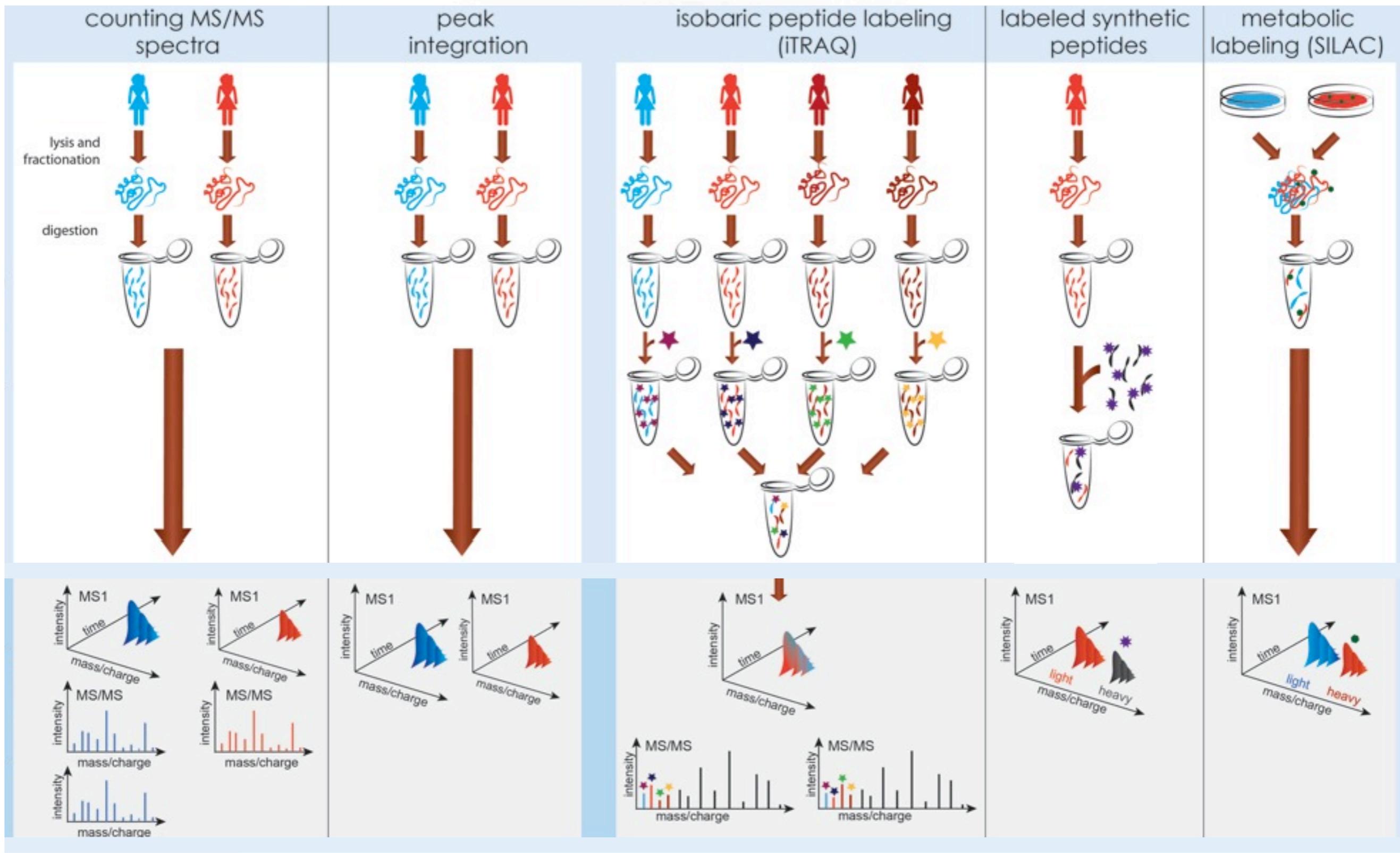
Instrumentation

Computational solutions



Global proteomic workflows

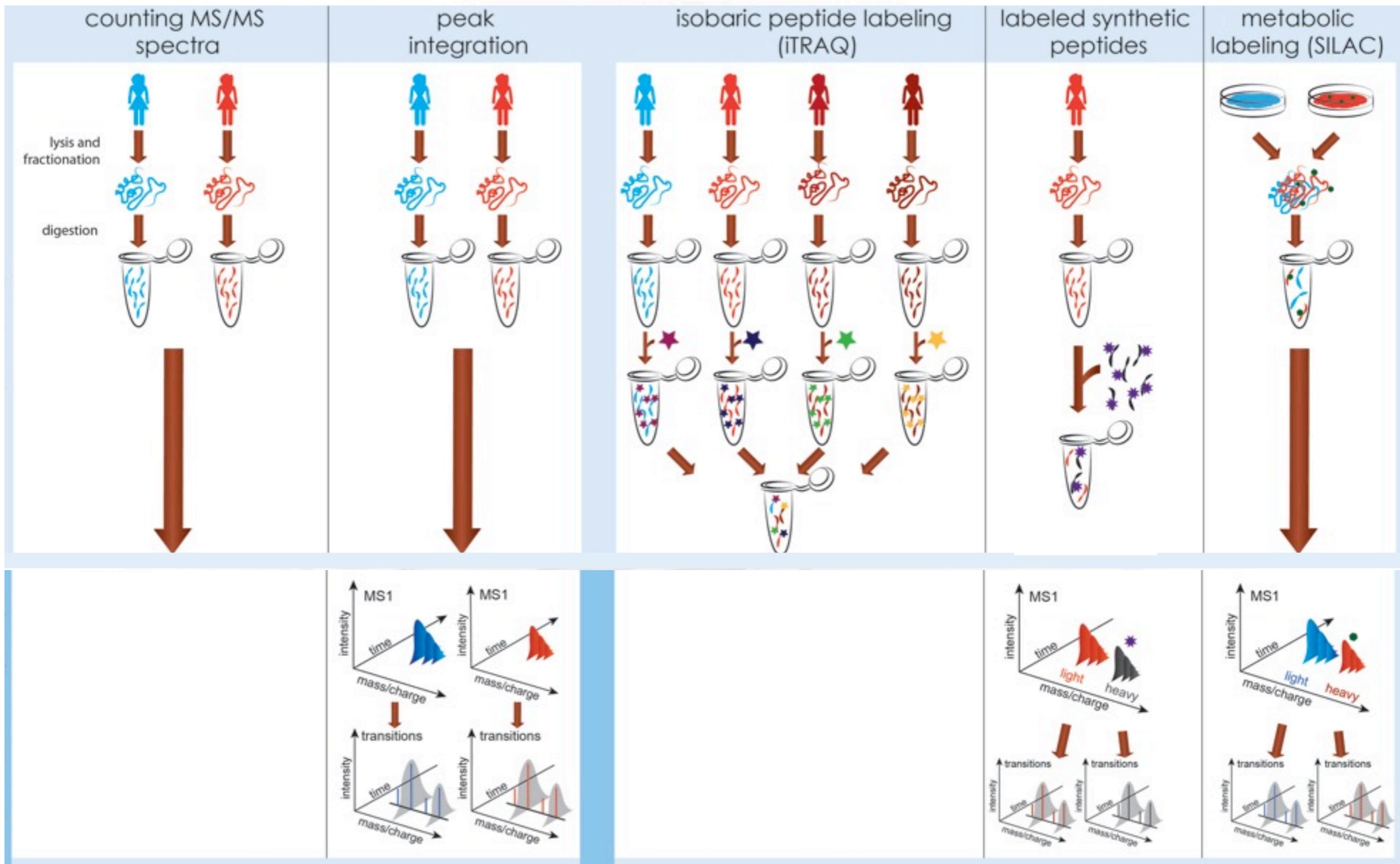
Sample preparation
Global LC-MS/MS



Targeted proteomic workflows

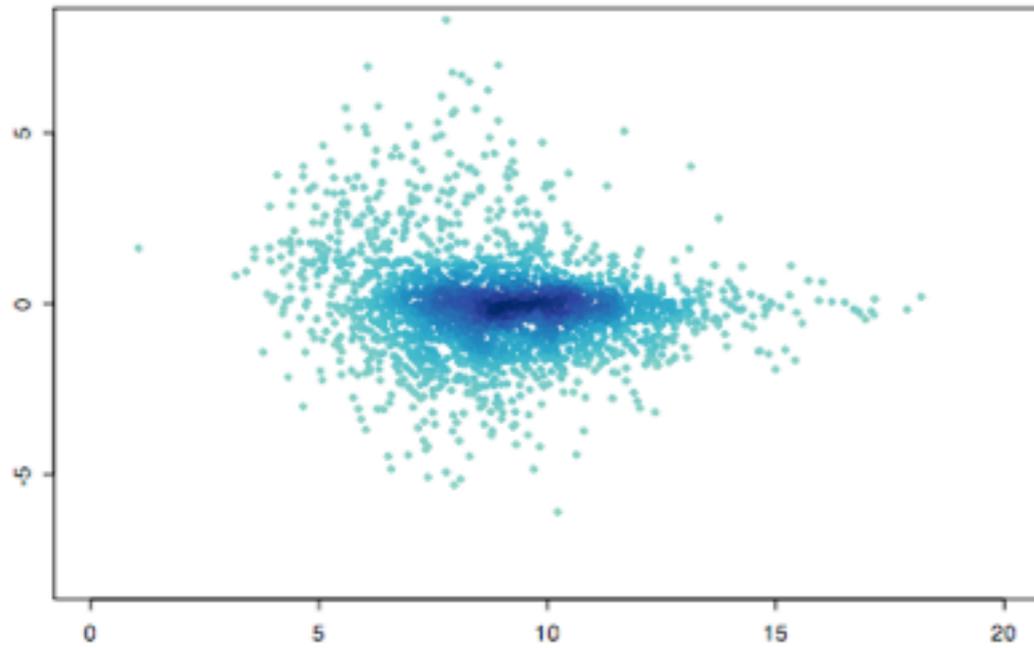
Sample preparation

Targeted SRM

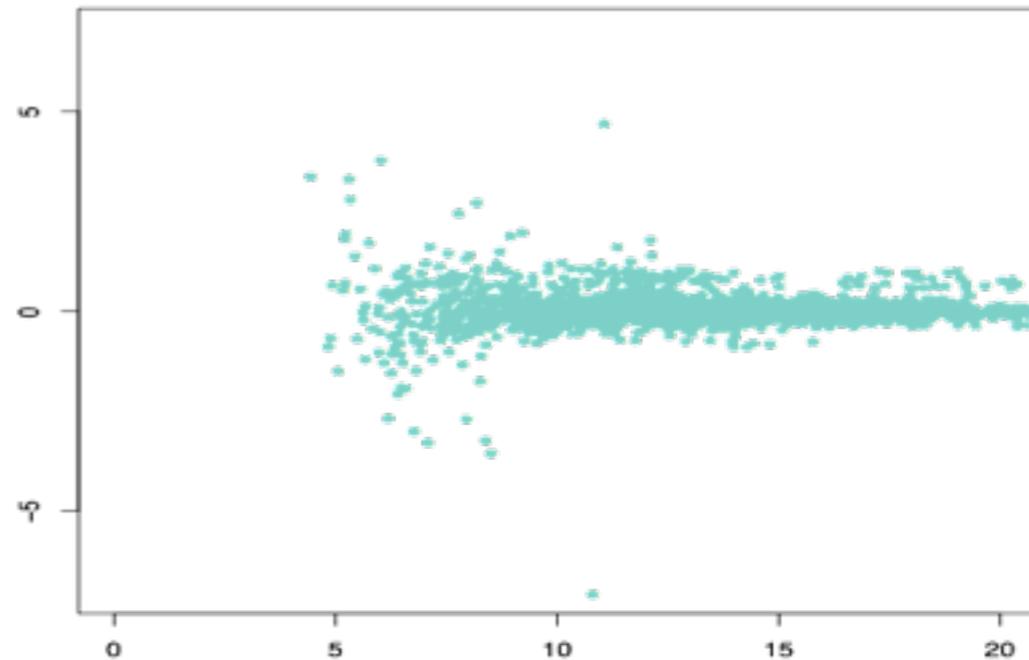


Variation is experiment-specific

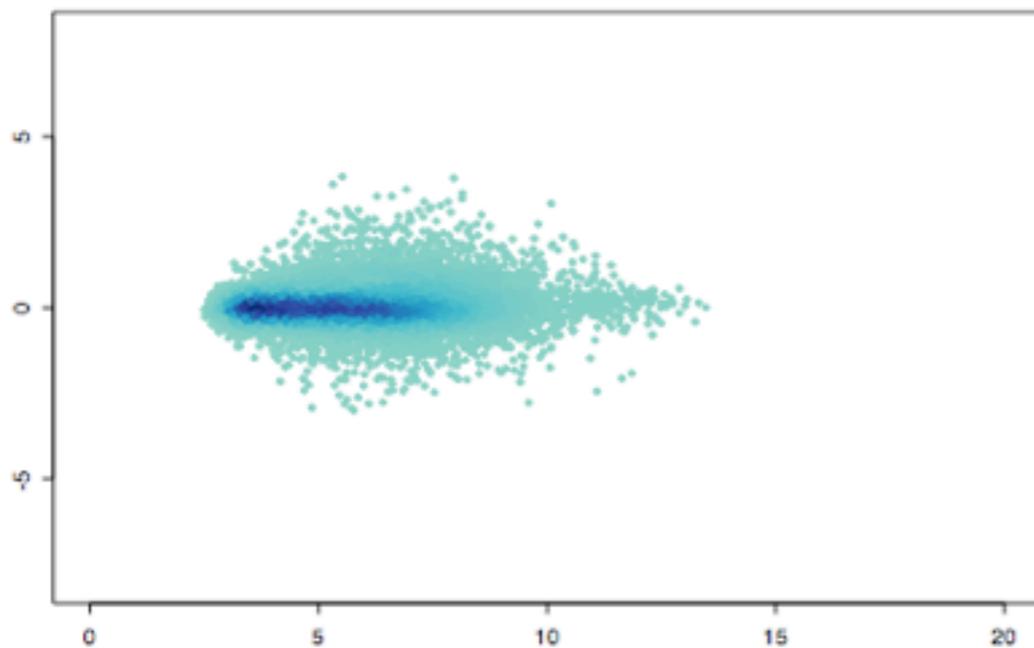
LC-MS: Human diabetes



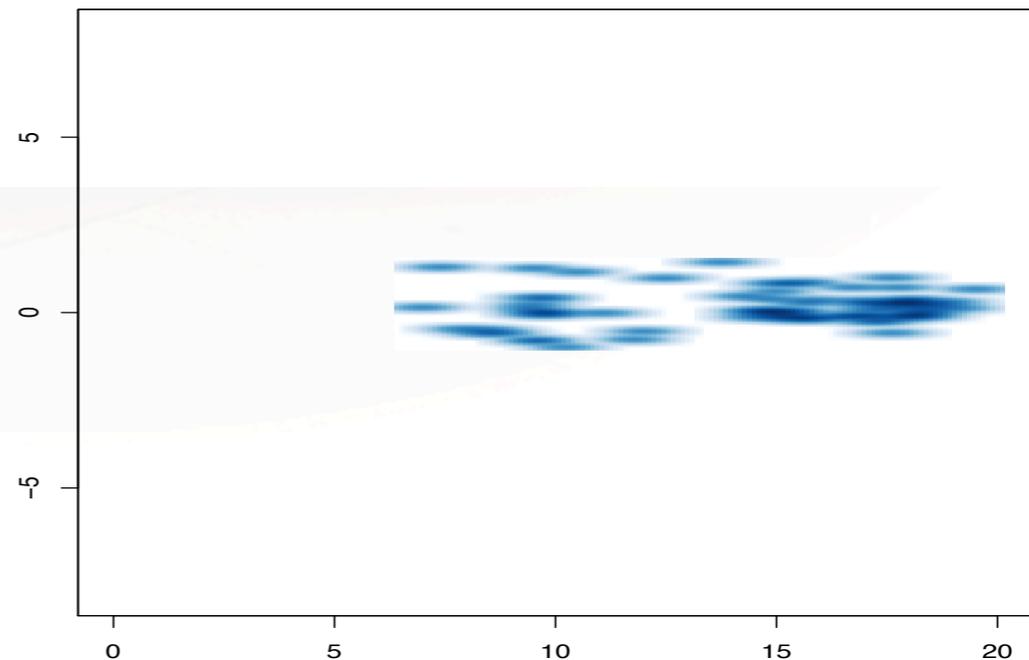
SRM: yeast proteome
(Picotti *et al*, 2009)



Affy: ALL (Acute Leukemia)



SRM: Human ovarian cancer
(Hüttenhain *et al*, 2012)



Difference of intensities

Outline

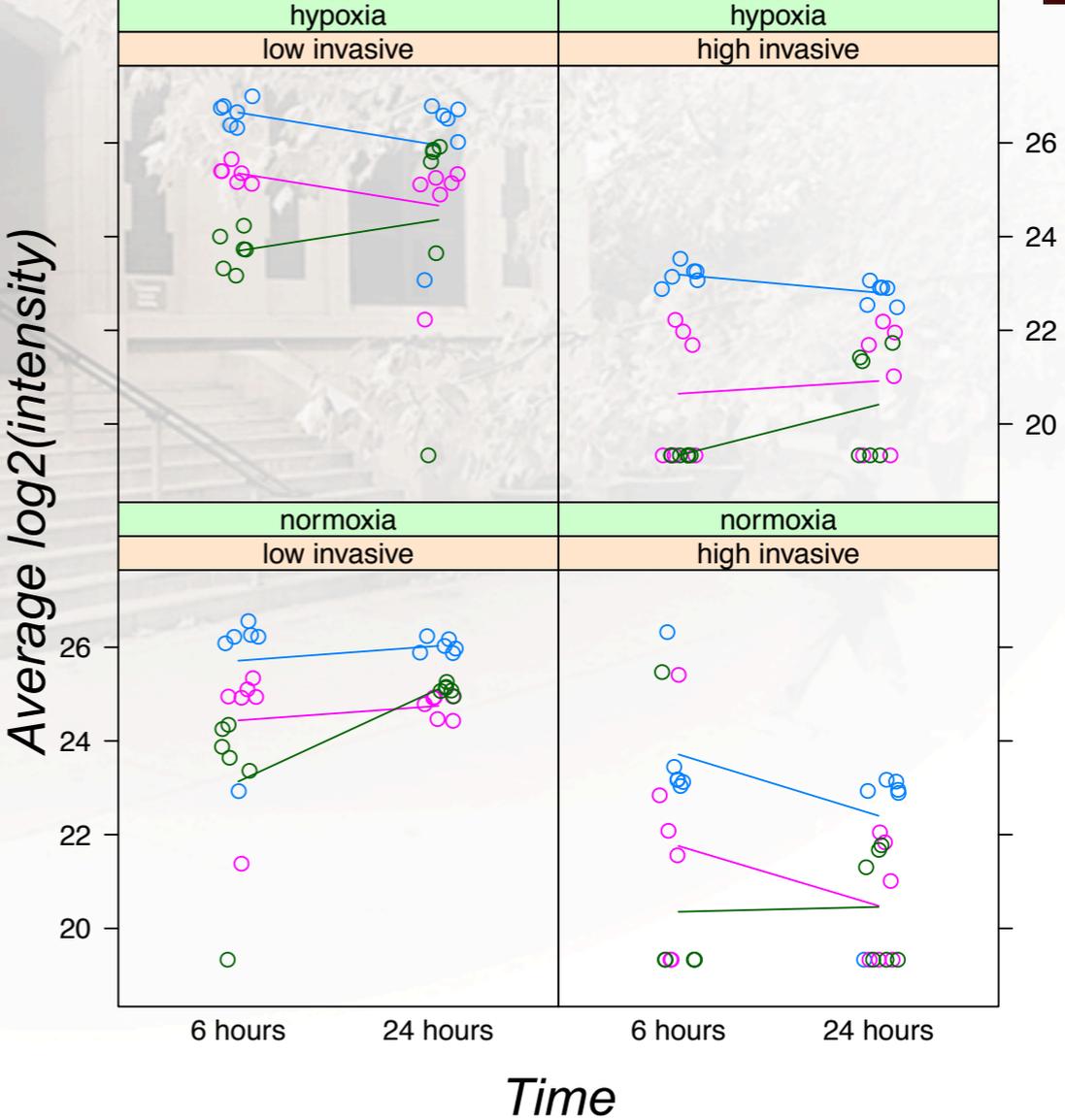
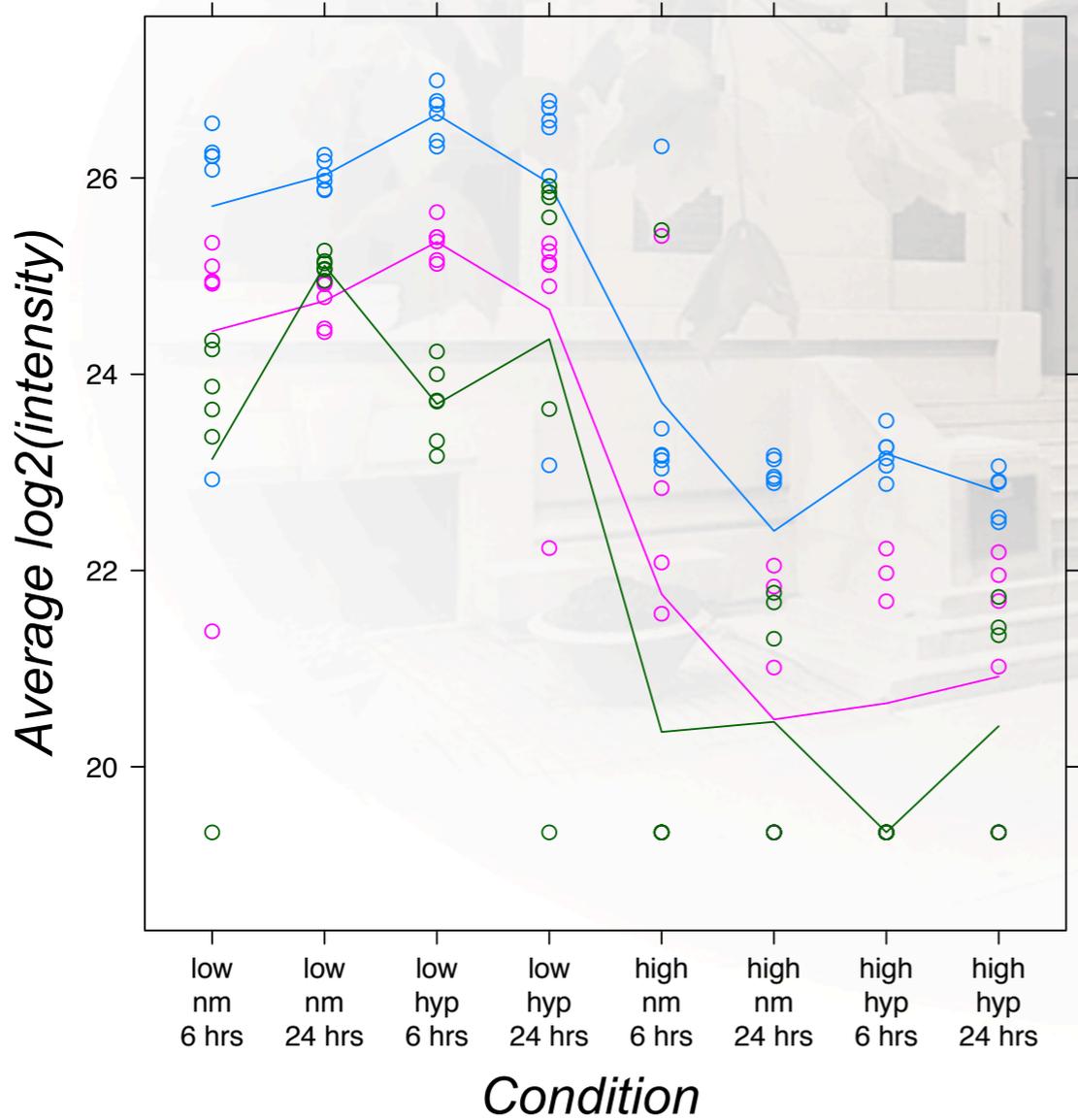
- What is proteomics?
 - ◆ Biological questions and technologies
- Protein quantification in label-free workflows
 - ◆ Joint analysis of multiple features and conditions
- Protein quantification in label-based workflows
 - ◆ Appropriately account for the labeling structure
- Mass spectrometry-based imaging
 - ◆ Account for the spacial heterogeneity of spectral data

Example: label-free LC-MS

High/low invasive breast cancer cell lines

Goal: protein-level conclusions

Safia
Thaminy
ETHZ



Important differences with microarrays

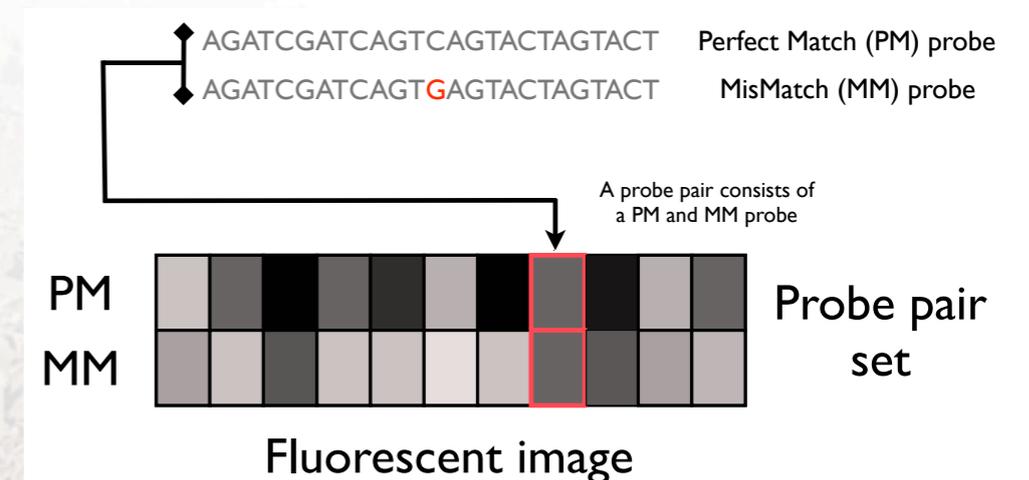
- Oligonucleotide microarrays

- ◆ RMA: Tukey Median Polish
 - Robust averaging of all probes
 - array-specific summary

- Proteomics

- ◆ Number and quality of features per protein vary widely
- ◆ Missing features introduce imbalance
- ◆ Label-based workflows combine multiple samples
 - blocking structure
- ◆ Targeted workflows create a nested structure protein/peptide/transition

- Explicit probabilistic models best represent the data



Sebastiani et al, Statistical Science, 2003.

Linear mixed models for feature intensities

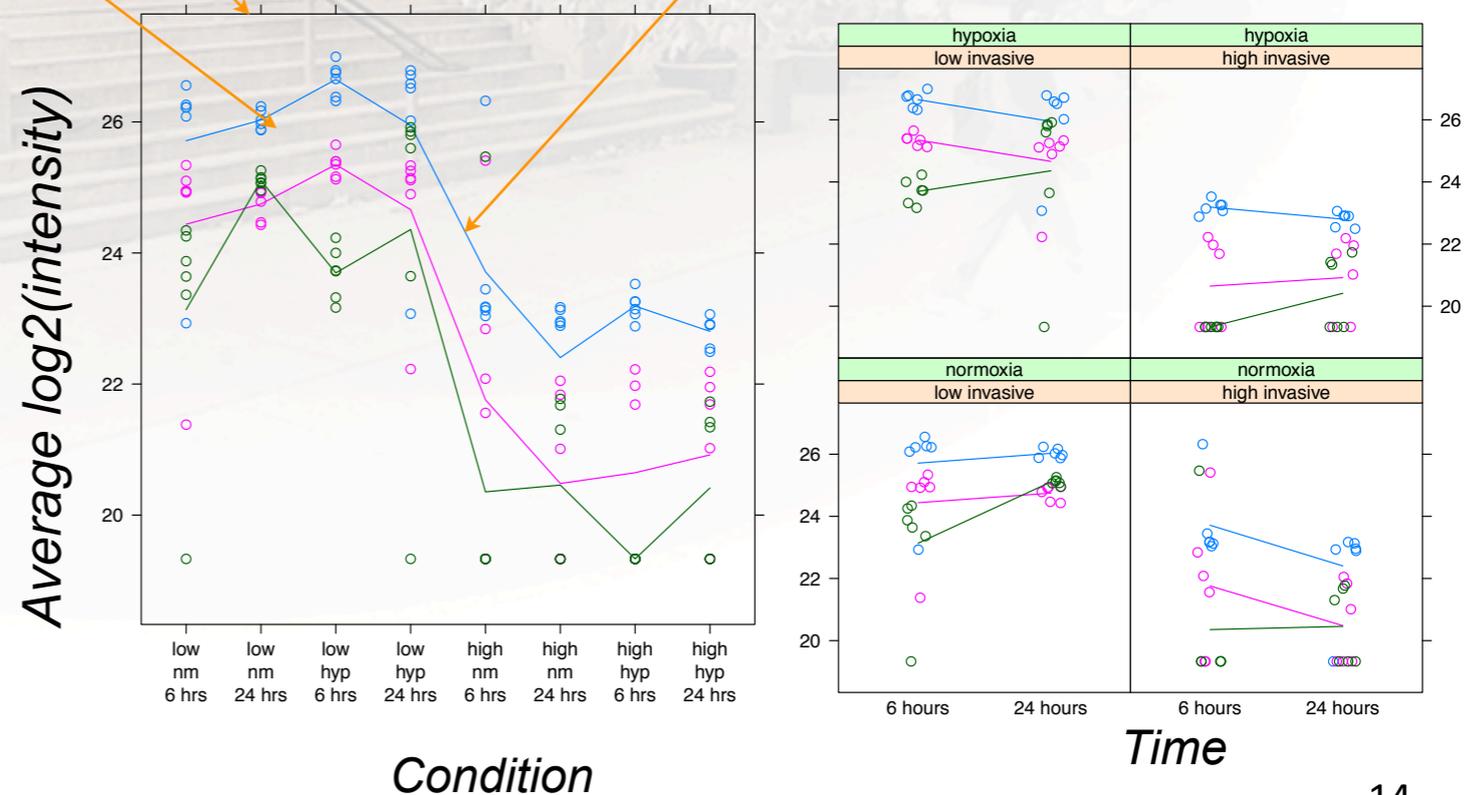
		Deviation from the reference due to									
log(peak intensity)	= Expected reference abundance	+	LC-MS feature	+	condition	+	feature × condition interaction	+	biol. replicate	+	Random meas. error
y_{ijkl}	= μ_{111}	+	F_i	+	C_j	+	$(F \times C)_{ij}$	+	$S(C)_{k(j)}$	+	ϵ_{ijkl}
where					$F_1 = C_1 = (F \times C)_{i1} = (F \times C)_{1j} = 0$						
and	(a) reduced scope of biological replication:				$S(C)_{1(1)} = 0$						
	(b) expanded scope of biological replication:				$S(C)_{k(j)} \stackrel{iid}{\sim} N(0, \sigma_S^2)$					$\epsilon_{ijkl} \stackrel{iid}{\sim} N(0, \sigma_{Error,ijk}^2)$	

T. Clough et al.
JPR, 2009

T. Clough et al.
Meth. Mol. Biol., 2011

T. Clough et al.
BMC Bioinformatics., 2012

C.-Y. Chang et al.
Mol. Cel. Proteomics, 2012



All inference is based on problem-specific linear combinations of model terms

Quantity of interest:

$$H_0 : L = \bar{\mu}_{[\text{high, nm, 6}]} - \bar{\mu}_{[\text{low, nm, 6}]} = 0$$

Model-based estimate and test statistic:

$$\hat{L} = \hat{C}_{[\text{high, nm, 6}]} + \frac{1}{I} \sum_{i=1}^I (\widehat{F \times C})_{i, [\text{high, nm, 6}]} + \frac{1}{K} \sum_{k=1}^K \widehat{S(C)}_{k([\text{high, nm, 6}]})$$

$$- \left(\hat{C}_{[\text{low, nm, 6}]} + \frac{1}{I} \sum_{i=1}^I (\widehat{F \times C})_{i, [\text{low, nm, 6}]} + \frac{1}{K} \sum_{k=1}^K \widehat{S(C)}_{k([\text{low, nm, 6}]}) \right)$$

$$t = \frac{\hat{L}}{SE\{\hat{L}\}} \sim \text{Student distribution}$$

In balanced datasets:

$$\hat{L} = \bar{Y}_{\cdot[\text{high, nm, 6}]\cdot} - \bar{Y}_{\cdot[\text{low, nm, 6}]\cdot}$$

$$t = \frac{\hat{L}}{\sqrt{\frac{2}{IKL} \hat{\sigma}_{Error}^2}} \sim \text{Student}_{IJK(L-1)+(I-1)J(K-1)} \text{ distribution}$$

Also methods for model diagnostics, data visualization etc

Tools for quantitative proteomics



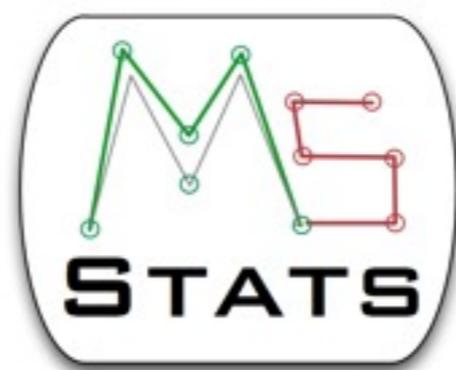
Veavi Chang



Meena Choi



Tim Clough



- A variety of quantitative workflow
 - *Global, targeted, data-independent*
 - *Label-free and label-based*
- Accounts for experimental designs
 - *Group comparison, time course*
- Data visualization & QC
- Model-based analysis
 - *Model fitting and inference*
- Planning future experiments
 - *Sample size, resource allocations*

Since 2010:

- *extensive documentation*
- *published case studies*
- *protocols for typical analyses*
- *13 tutorials and workshops*

Since 2011:

- *370 unique visitors*
- *over 50 unique downloads*
- *over 50 mailing list members*



Collaboration:



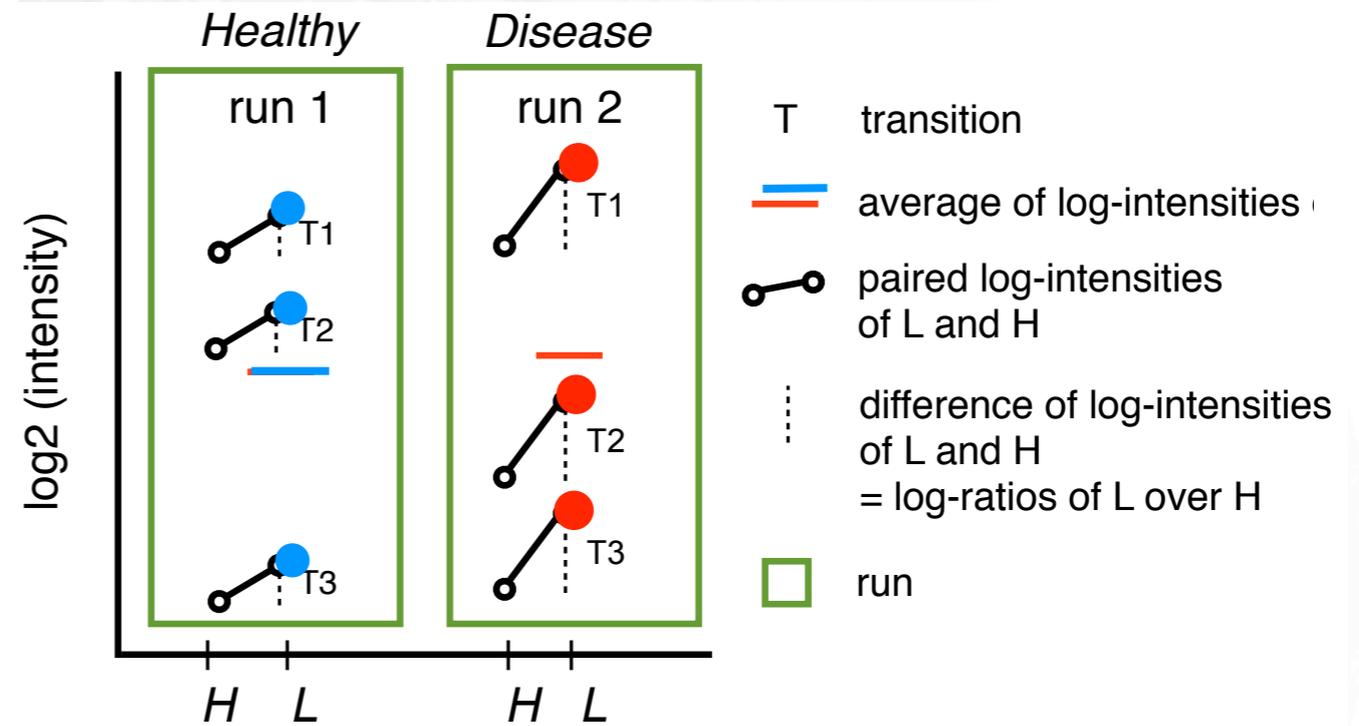
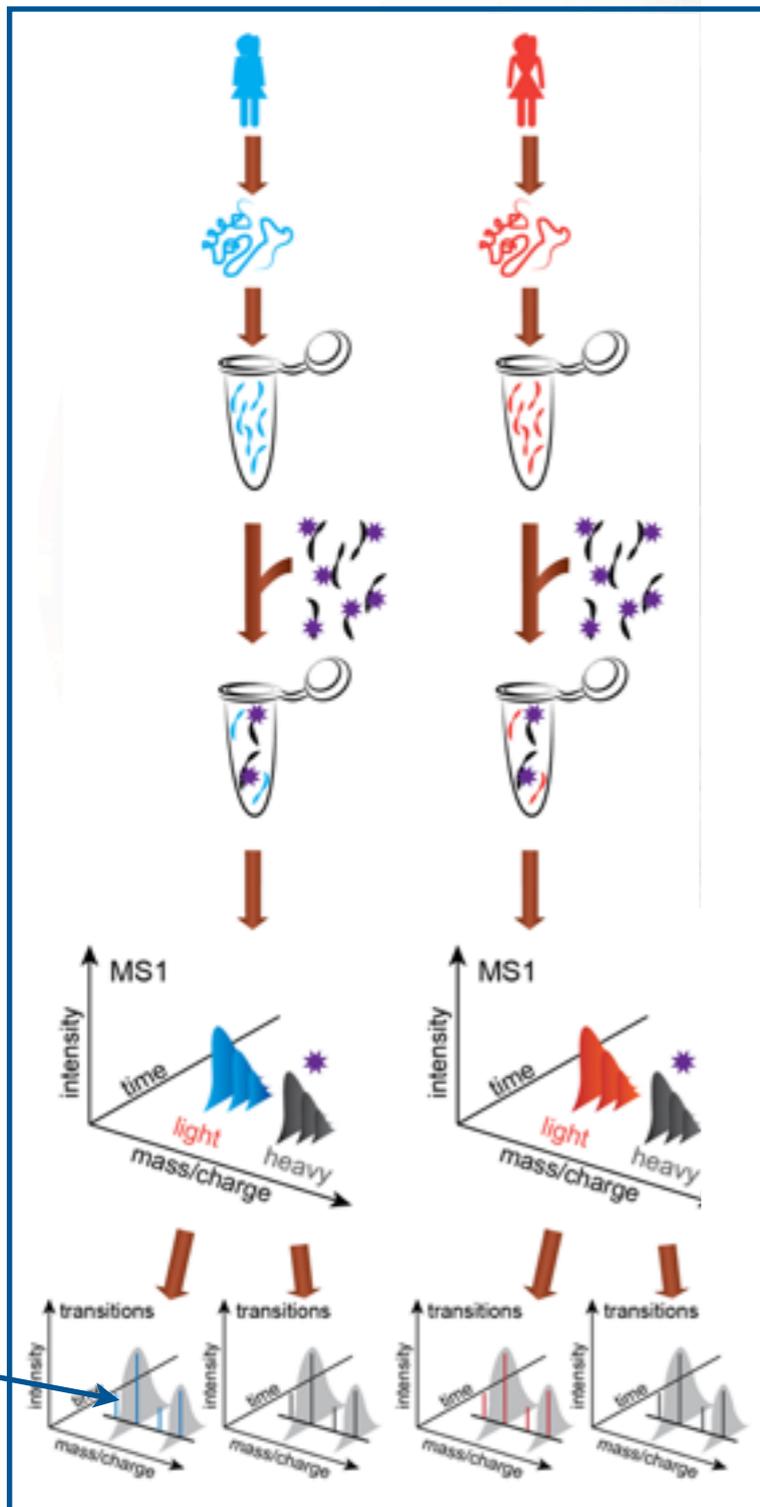
Michael
MacCoss



Outline

- What is proteomics?
 - ◆ Biological questions and technologies
- Protein quantification in label-free workflows
 - ◆ Joint analysis of multiple features and conditions
- Protein quantification in label-based workflows
 - ◆ Appropriately account for the labeling structure
- Mass spectrometry-based imaging
 - ◆ Account for the spacial heterogeneity of spectral data

Label-based workflows help separate the biological and the technological variation



			Group 1		...	Group I		
		Run 1	Subject 1	Subject J	Subject 1	Subject J	Run M	Subject J
Endogenous: light labeled peptide	Peptide 1	Transition 1	10.21	10.57	15.64	15.03		
	...	Transition L	10.52	10.92	15.29	15.68		
	Peptide K	Transition 1	11.76	11.92	16.22	16.71		
	...	Transition L	11.65	11.09	16.27	16.51		
Reference: heavy labeled peptide	Peptide 1	Transition 1	19.46	19.77	19.82	19.03		
	...	Transition L	19.13	19.25	19.67	19.80		
	Peptide K	Transition 1	19.26	19.33	19.58	19.61		
	...	Transition L	19.73	19.09	19.84	19.55		

Legend : Label Feature: Transition/Peptide Group Run Subject

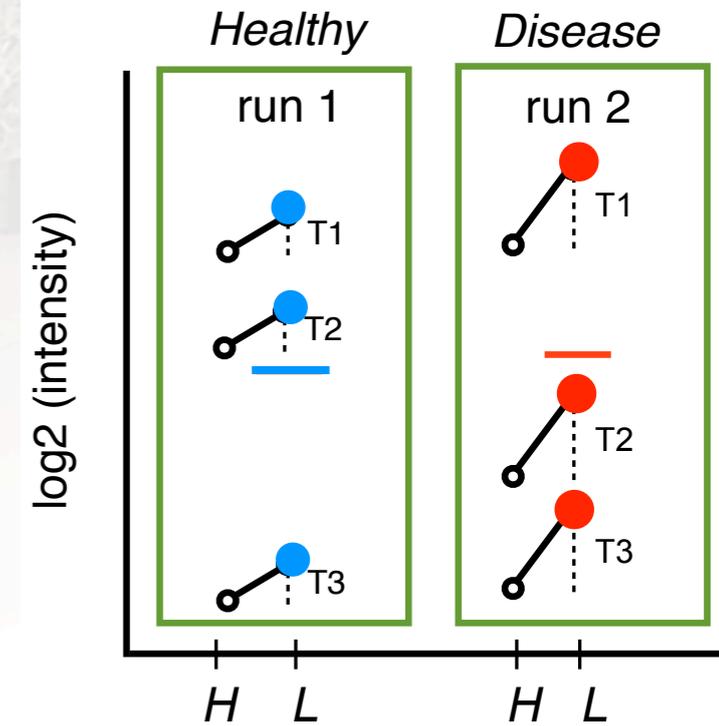
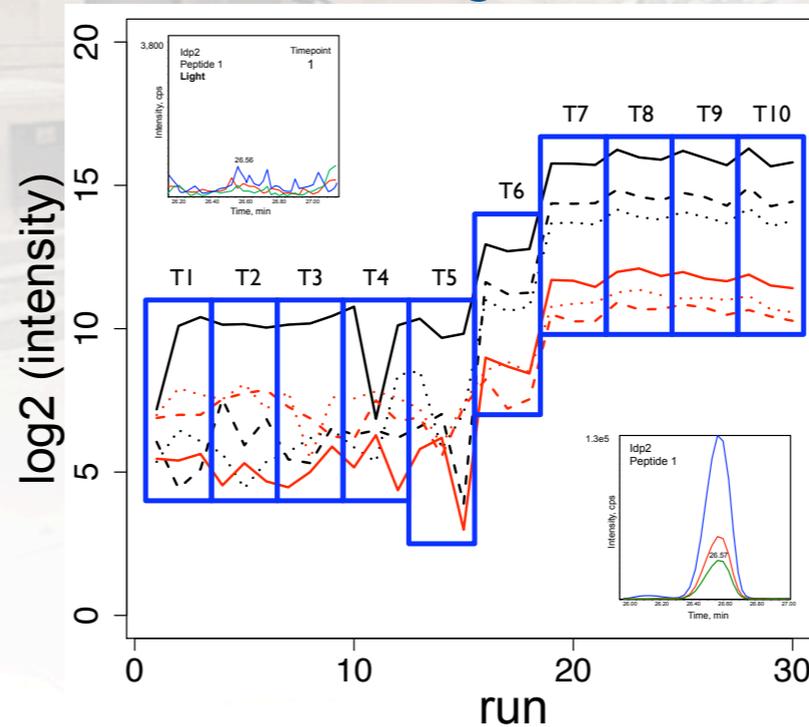
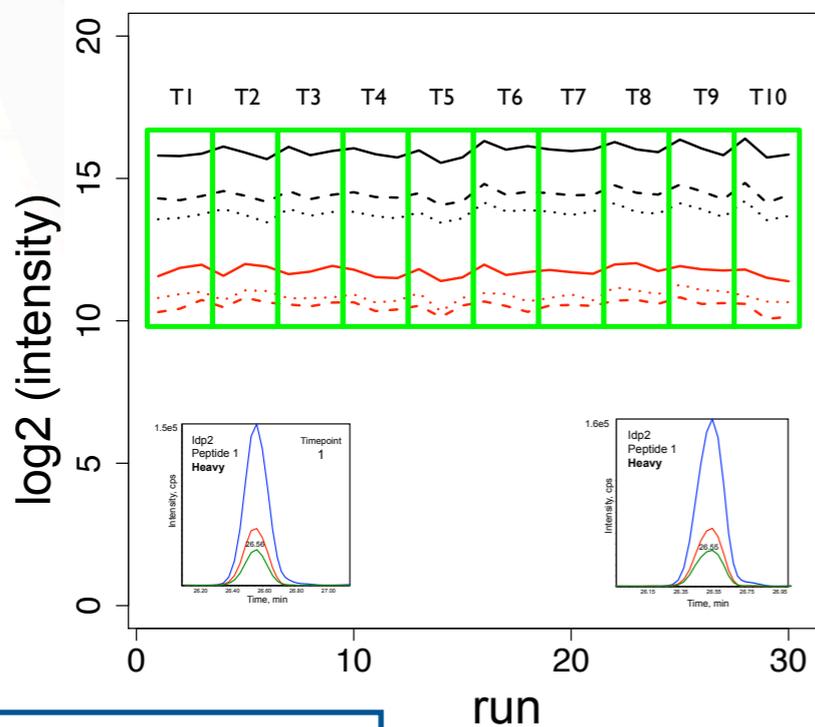
Transitions

More complex but similar linear mixed effects models

Systematic deviations from the mean						Statistical interactions													
Observed $\log_2(\text{int of peak})$	=	Overall mean	+	Group or time	+	Subject	+	Feature	+	Run	+	Time by subject	+	Time by feature	+	Run by feature	+	Random error	
y_{ijklm}	=	μ	+	T_i	+	$S(G)_{j(i)}$	+	F_{kl}	+	R_m	+	$(T \times S)_{ij}$	+	$(T \times F)_{ikl}$	+	$(R \times F)_{klm}$	+	ε_{ijklm}	
Fixed/Random		F		F		F/R		F		F/R		F/R		F		F/R		F/R	R: $N(0, \sigma^2)$

Reference

Endogenous



C.-Y. Chang et al.
MCP, 2012

Purdue University



Kyle Bemis Meena Choi Tim Clough Veavi Chang Robert Ness Danni Yu Cheng Zheng

ETH Zürich



Ruedi Aebersold and lab



Ruth Hüttenhain

Purdue



Graham Cooks and lab

Stanford



Mark Stolowitz



Parag Mallick

CGR, Spain



Eduard Sabidó



Safia Thaminy



Silvia Surinova



- ◆ NSF-SI2-SSE
- ◆ NSF-BIO-DBI