



# R Fundamentals and Best Practices for Mass Spectrometry Data Analysis

Sunday, November 15 (12:00-3:15pm Eastern)

**Meena Choi**, Genentech

**Olga Vitek**, Northeastern University

Module #7: Analysis of proteomic data with MSstats



# Introduction to MSSTATS



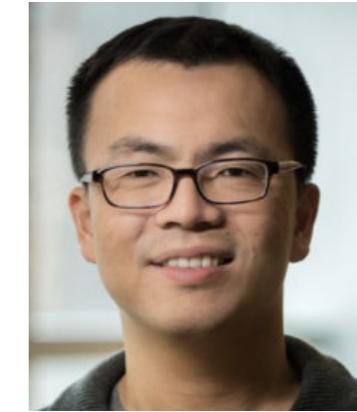
Meena Choi  
*Genentech*



Ting Huang  
*Northeastern*



Mateusz Staniak  
*University of Wrocław*



Tsung-Heng Tsai  
*Kent State*

and

Veavi Chang  
Tim Clough  
Eralp Dogu  
Cyril Galitzine  
Devon Kohler  
Sara Taheri  
Mateusz Staniak  
Akshay Kulkarni  
Sumedh Sankhe



Olga Vitek  
Khoury College of  
Computer Sciences  
**Northeastern University**

## R-based open-source tool for quantitative proteomics

- Which proteins change in abundance?
  - Complex designs: factorial experiments, paired designs, time course
  - Abundance of each subject on a relative scale
- Label-free or label-based
  - Shotgun DDA, data independent DIA/SWATH, PRM, targeted SRM
  - TMT; soon: differential PTMs
- Data processing
  - Free, open-source and inter-operable with other tools
  - Converts to Skyline, MaxQuant, OpenMS, Spectronaut, ProteomeDiscoverer...
- Statistical analysis
  - Flexible models, account for complex designs, missing values and outliers
  - Stable parameter estimation
- Method validation, system suitability, experimental design
  - System suitability/QC, assay characterization
  - Sample size for testing and classification



# Relative protein quantification in MS-based proteomics



Analytical method validation



System suitability testing



Experimental design



Data acquisition  
Quality control



Data processing



Statistical analysis



Data for future designs

**Assay characterization**

LOB/LOD



**MSstatsQC**

*Software package*

System suitability monitoring

**Experimental design**

Sample size for testing and classification

*Meena Choi*

*Ting Huang*

*Mateusz Staniak*

*Tsung-Heng Tsai*

*R software package*

**MSnbase**  
**RforProteomics**

*Open source software*

**Skyline**, **MaxQuant**  
**OpenMS**, **OpenSWATH**  
**DIA-Umpire**

*Commercial software*

**Proteome Discoverer**  
**Spectronaut**, **SpectroMine**  
**Progenesis**

**MSstats**

*Software package*

Significance analysis for DDA, SRM, DIA

**MSstatsTMT**

*Software package*

Significance analysis for TMT

**MSstatsBioData**

*Experiment package*

Published studies with DDA or SRM

**MassIVE**

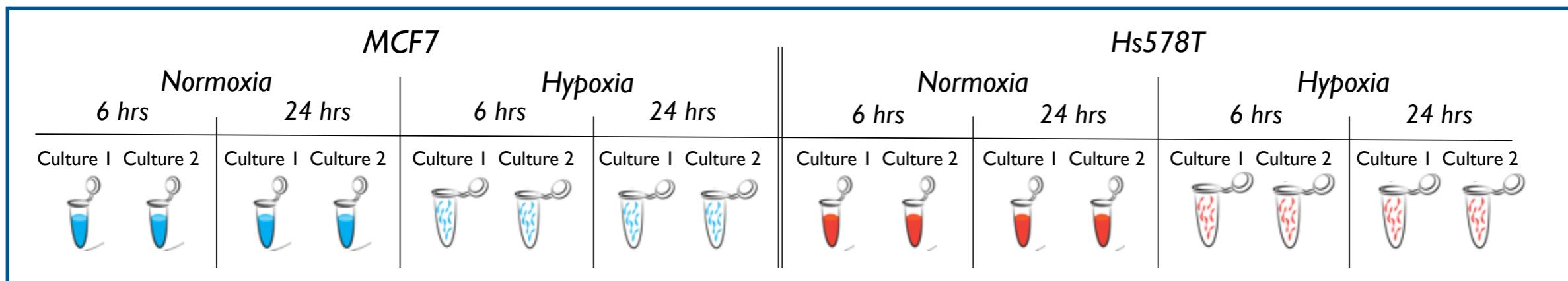
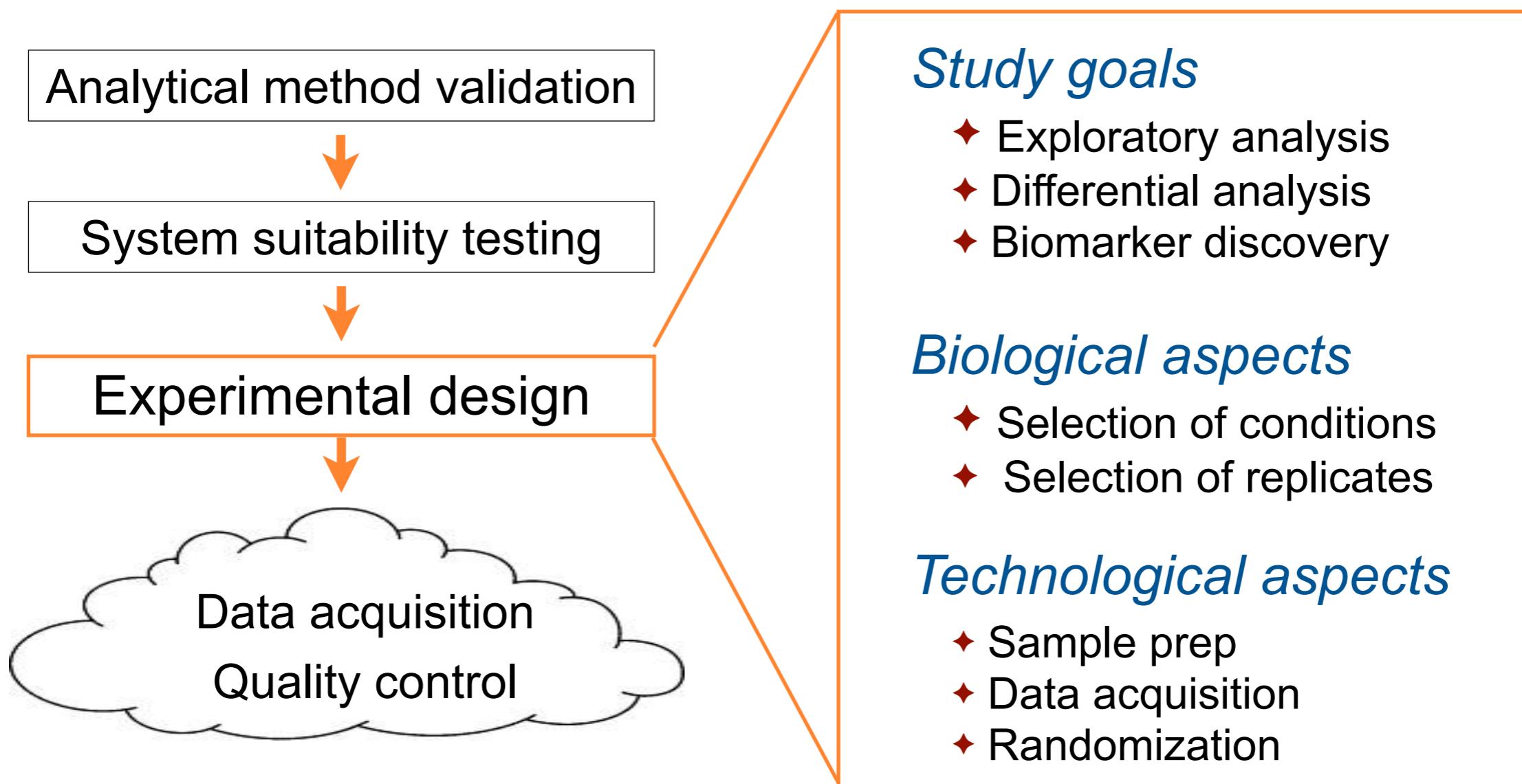
*Mass Spectrometry  
Interactive Virtual Environment*

42 datasets, > 178 reanalyses

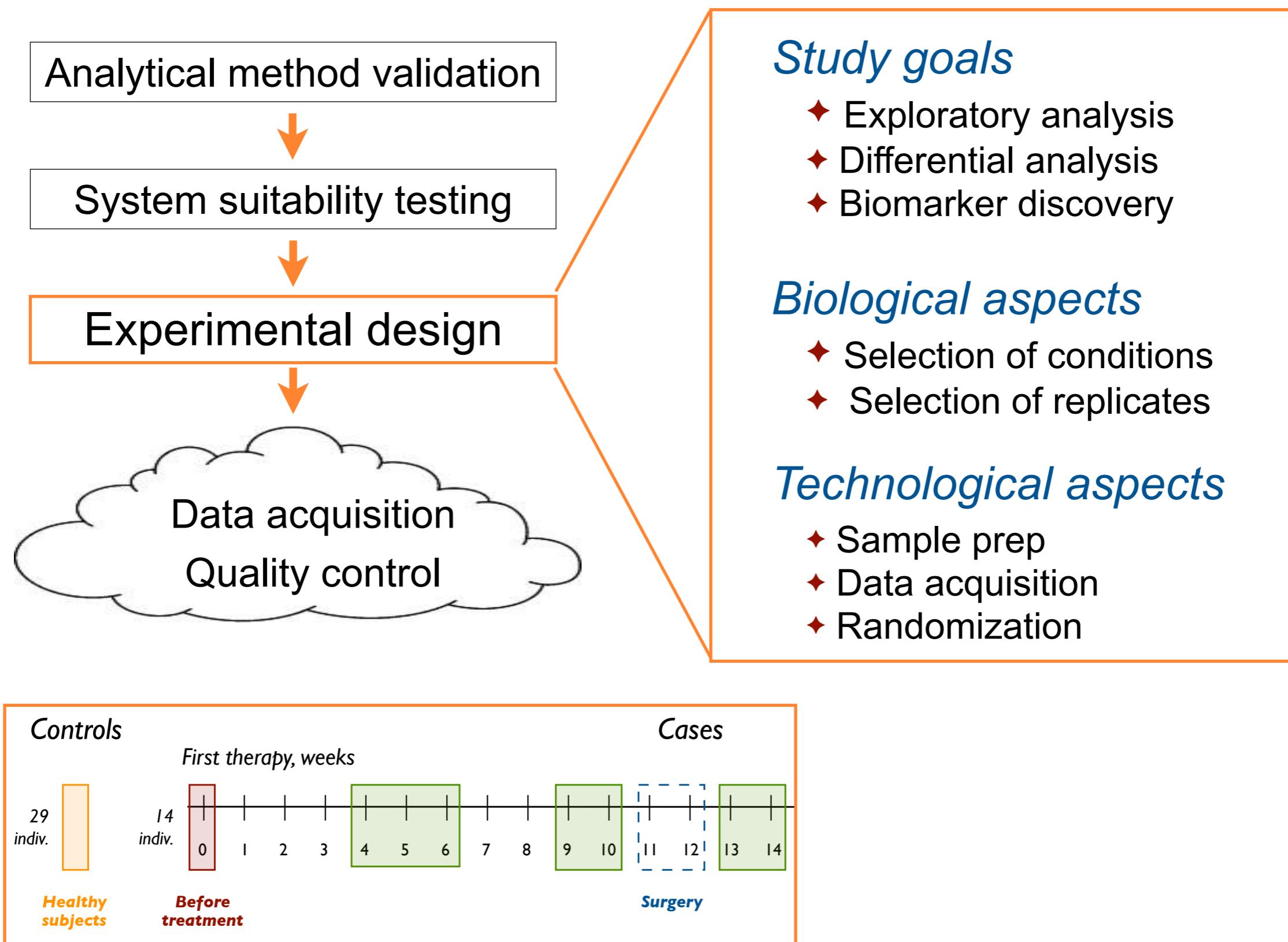
# INTRODUCTION TO MSSTATS

- Overview
- Statistical modeling
- Community resources

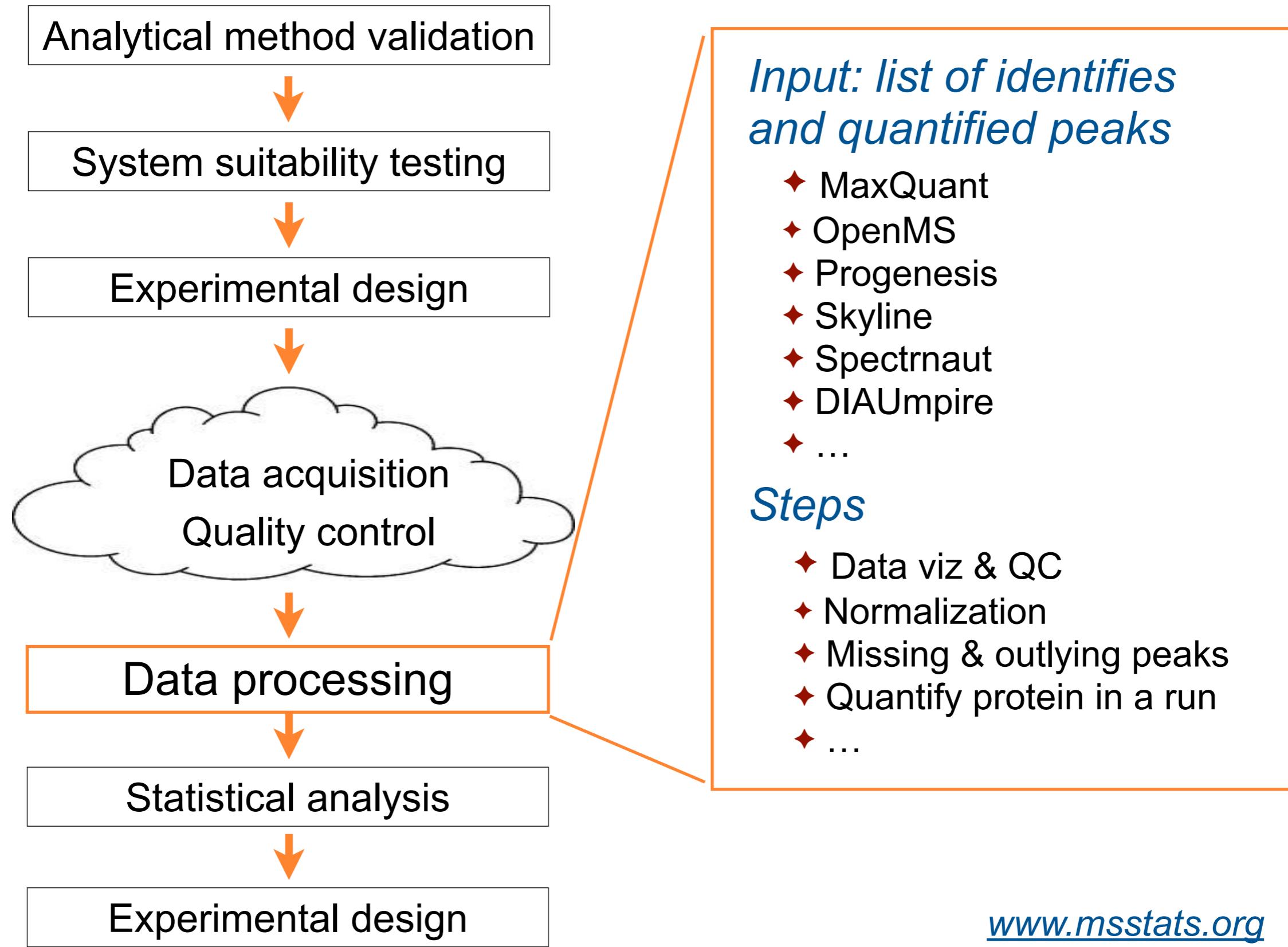
# MS EXPERIMENT: DESIGN



# MS EXPERIMENT: DESIGN



# MS EXPERIMENT: DATA PROCESSING

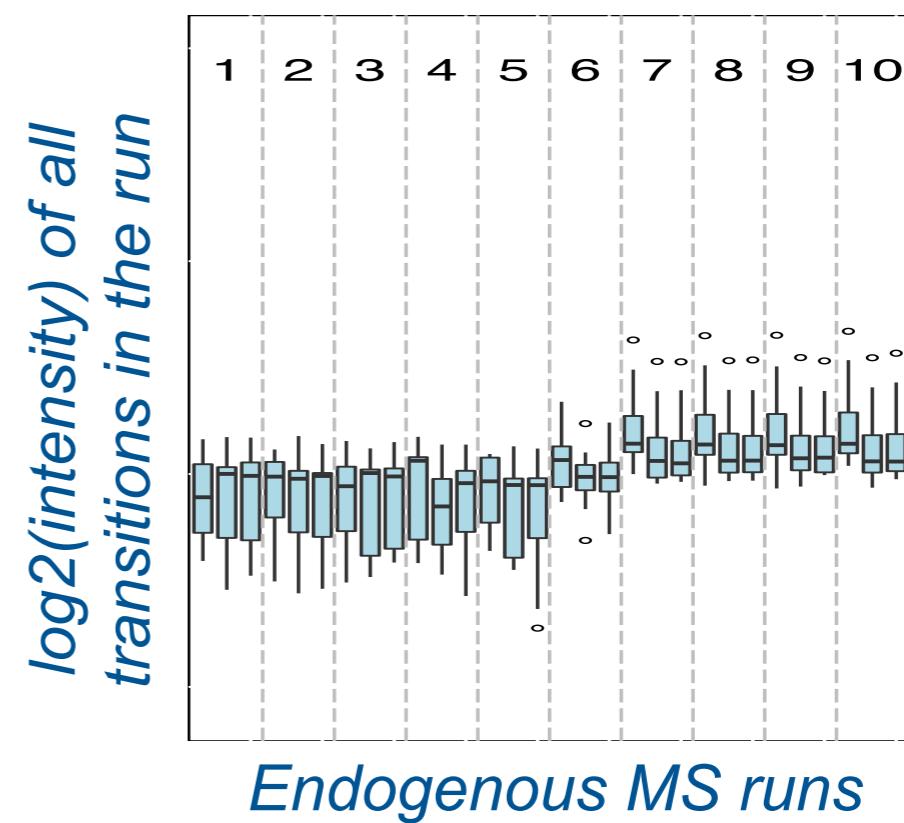
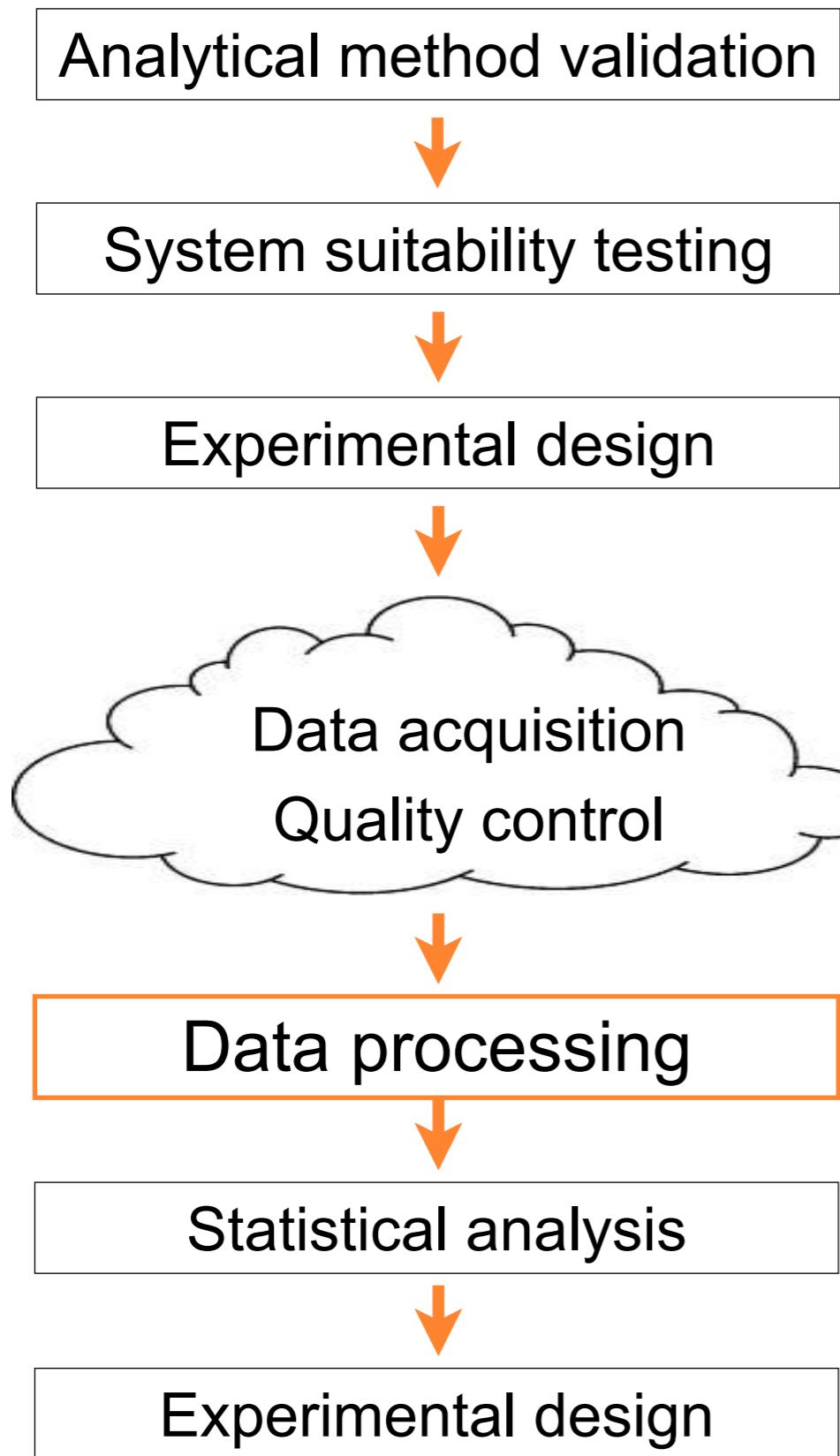


# INPUT DATA FORMAT

	A	B	C	D	E	F	G	H	I	J
1	ProteinName	PeptideSequence	PrecursorCharge	FragmentIon	ProductCharge	IsotopeLabelType	Condition	BioReplicate	Run	Intensity
2	ACEA	EILGHEIFFDWELP	3	y3	0	H	1	ReplA	1	66472.3847
3	ACEA	EILGHEIFFDWELP	3	y3	0	L	1	ReplA	1	5764.16228
4	ACEA	EILGHEIFFDWELP	3	y4	0	H	1	ReplA	1	101005.166
5	ACEA	EILGHEIFFDWELP	3	y4	0	L	1	ReplA	1	61.65238
6	ACEA	EILGHEIFFDWELP	3	y5	0	H	1	ReplA	1	90055.4993
7	ACEA	EILGHEIFFDWELP	3	y5	0	L	1	ReplA	1	472.691803
8	ACEA	TDSEAATLISSTID	2	y10	0	H	1	ReplA	1	43506.5425
9	ACEA	TDSEAATLISSTID	2	y10	0	L	1	ReplA	1	217.203553
10	ACEA	TDSEAATLISSTID	2	y7	0	H	1	ReplA	1	68023.0377
11	ACEA	TDSEAATLISSTID	2	y7	0	L	1	ReplA	1	725.284308
12	ACEA	TDSEAATLISSTID	2	y8	0	H	1	ReplA	1	68276.0489
13	ACEA	TDSEAATLISSTID	2	y8	0	L	1	ReplA	1	243.658527

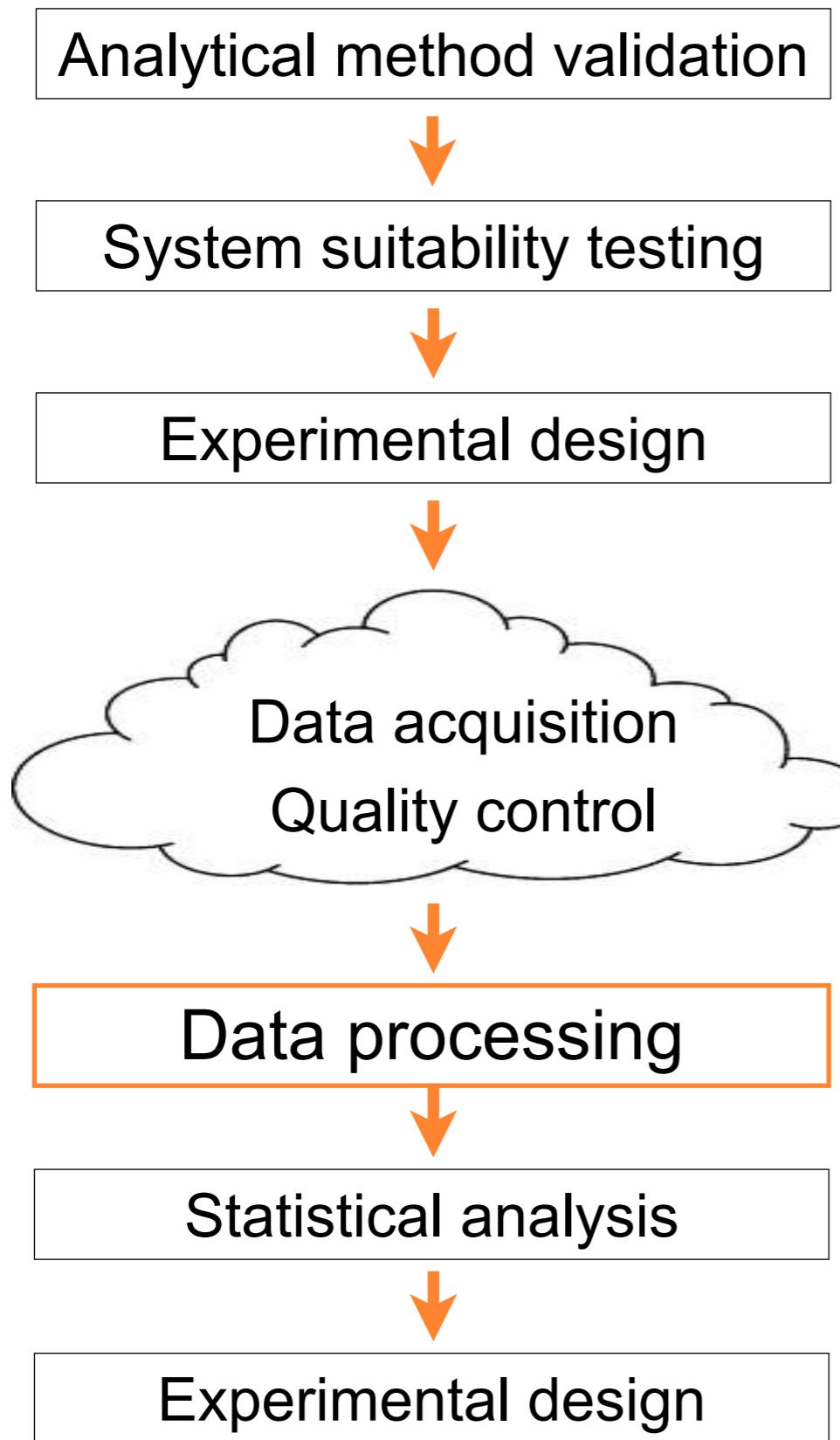
*Reports produced by (or converted from) data processing tools*

# MS EXPERIMENT: NORMALIZATION

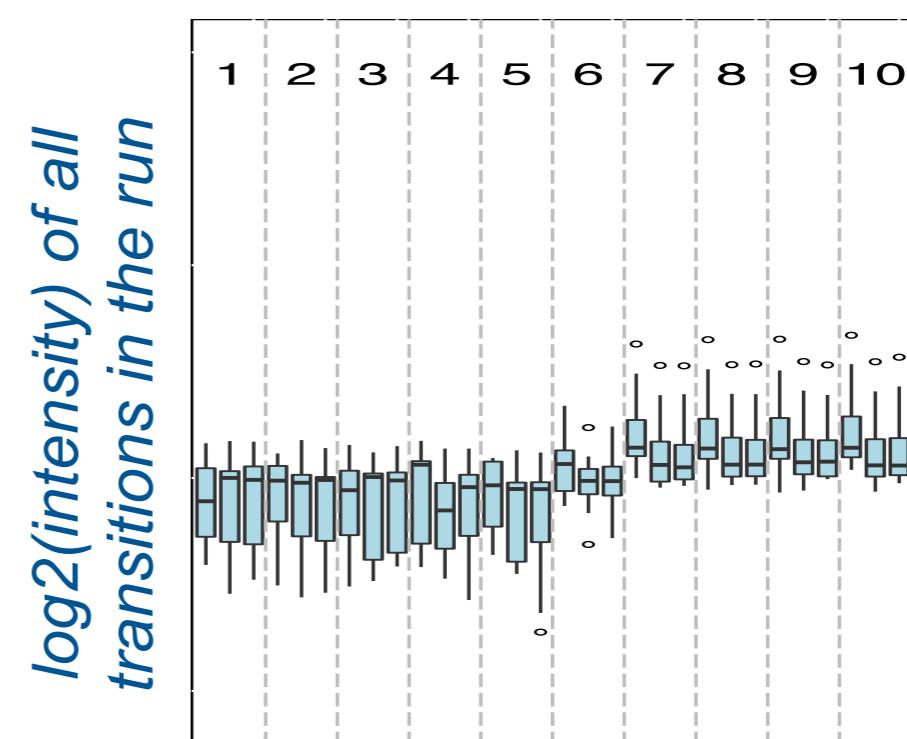


*Endogenous MS runs*

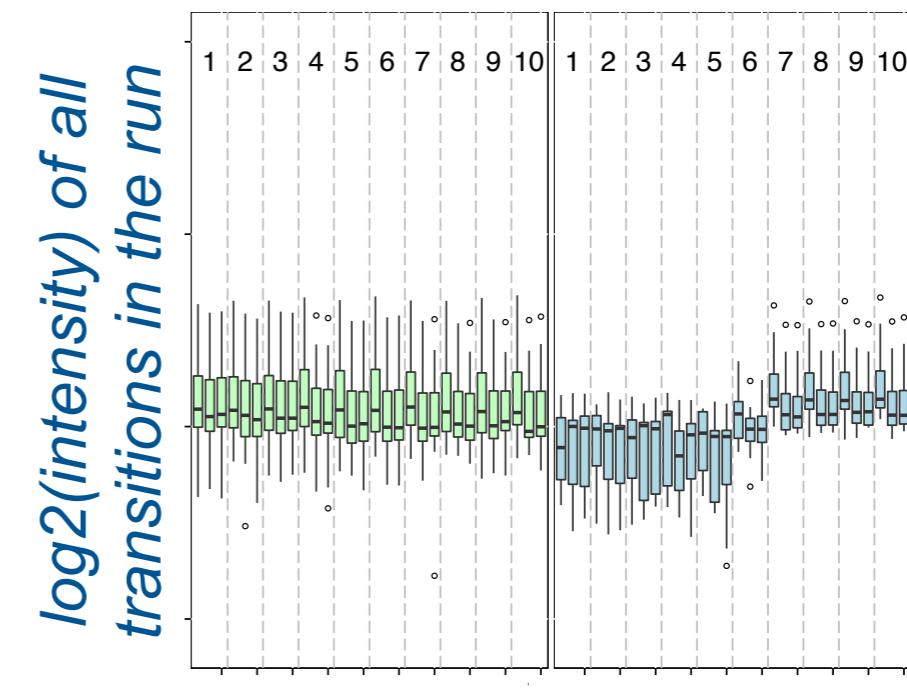
# MS EXPERIMENT: NORMALIZATION



## NORMALIZATION

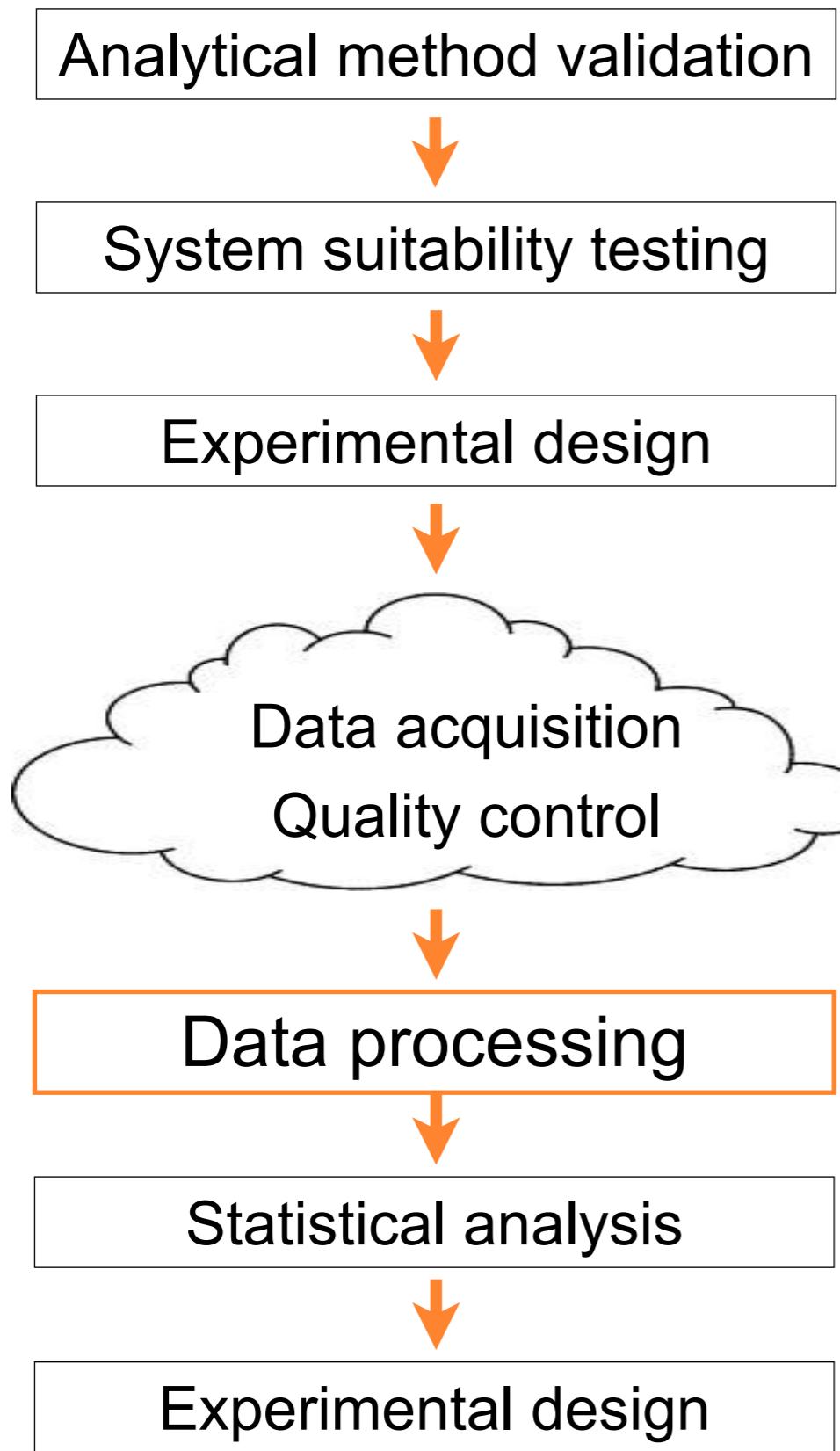


*Endogenous MS runs*



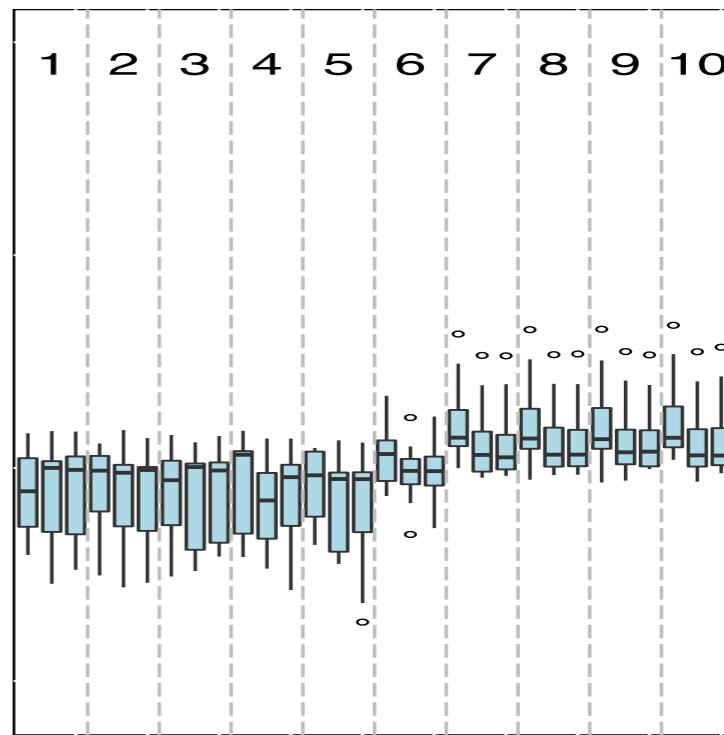
*Reference Endogenous MS runs*

# MS EXPERIMENT: NORMALIZATION



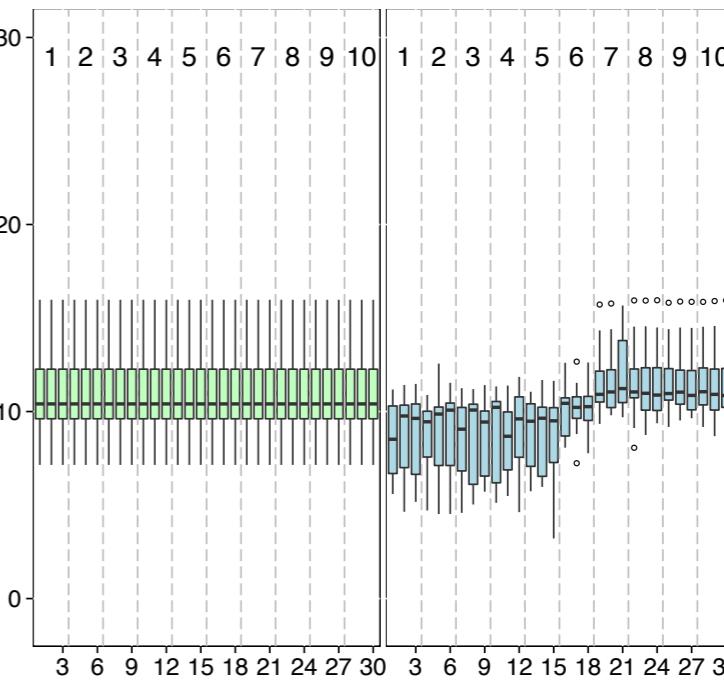
## NORMALIZATION

*log<sub>2</sub>(intensity) of all transitions in the run*



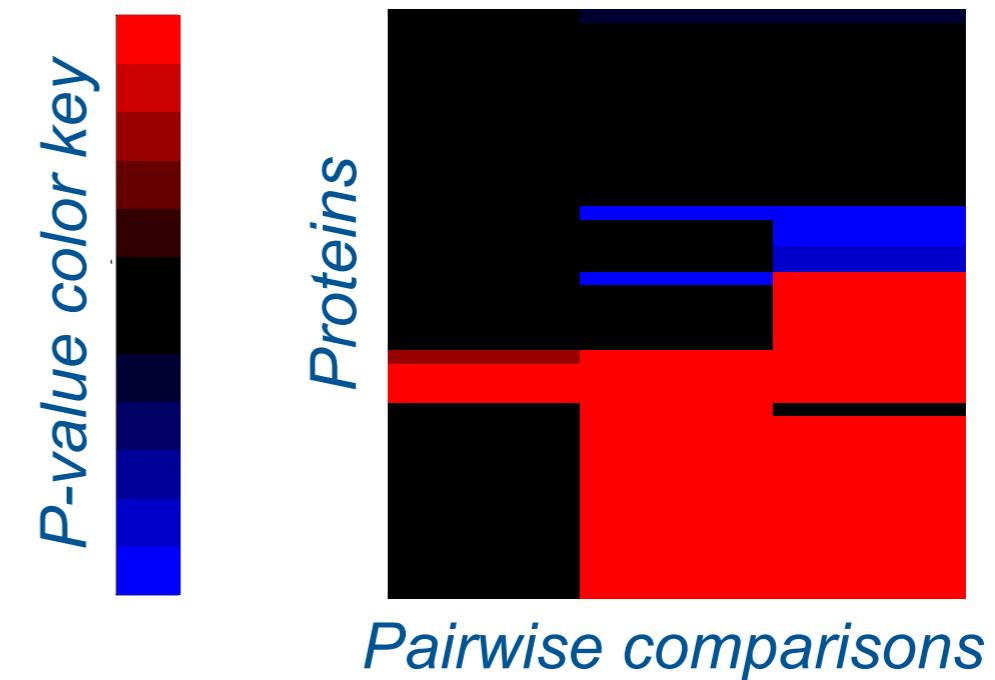
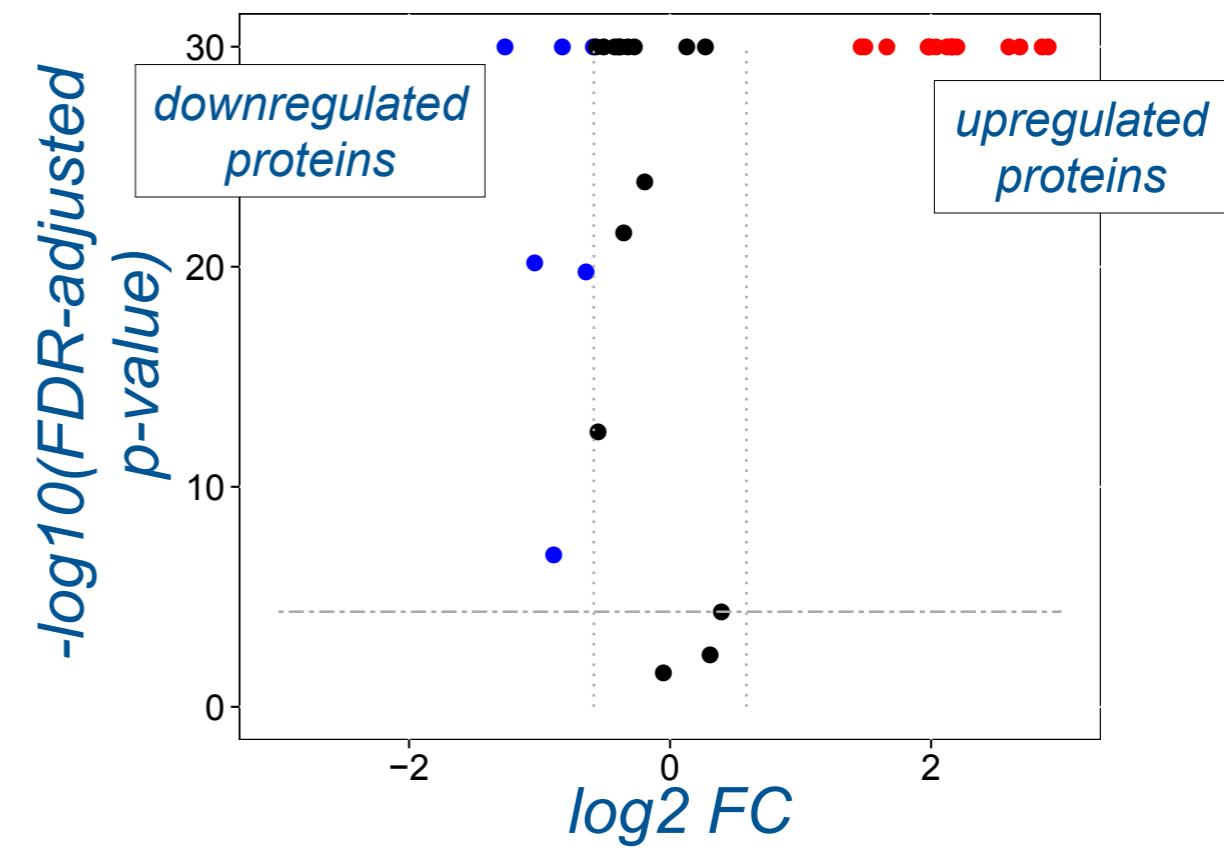
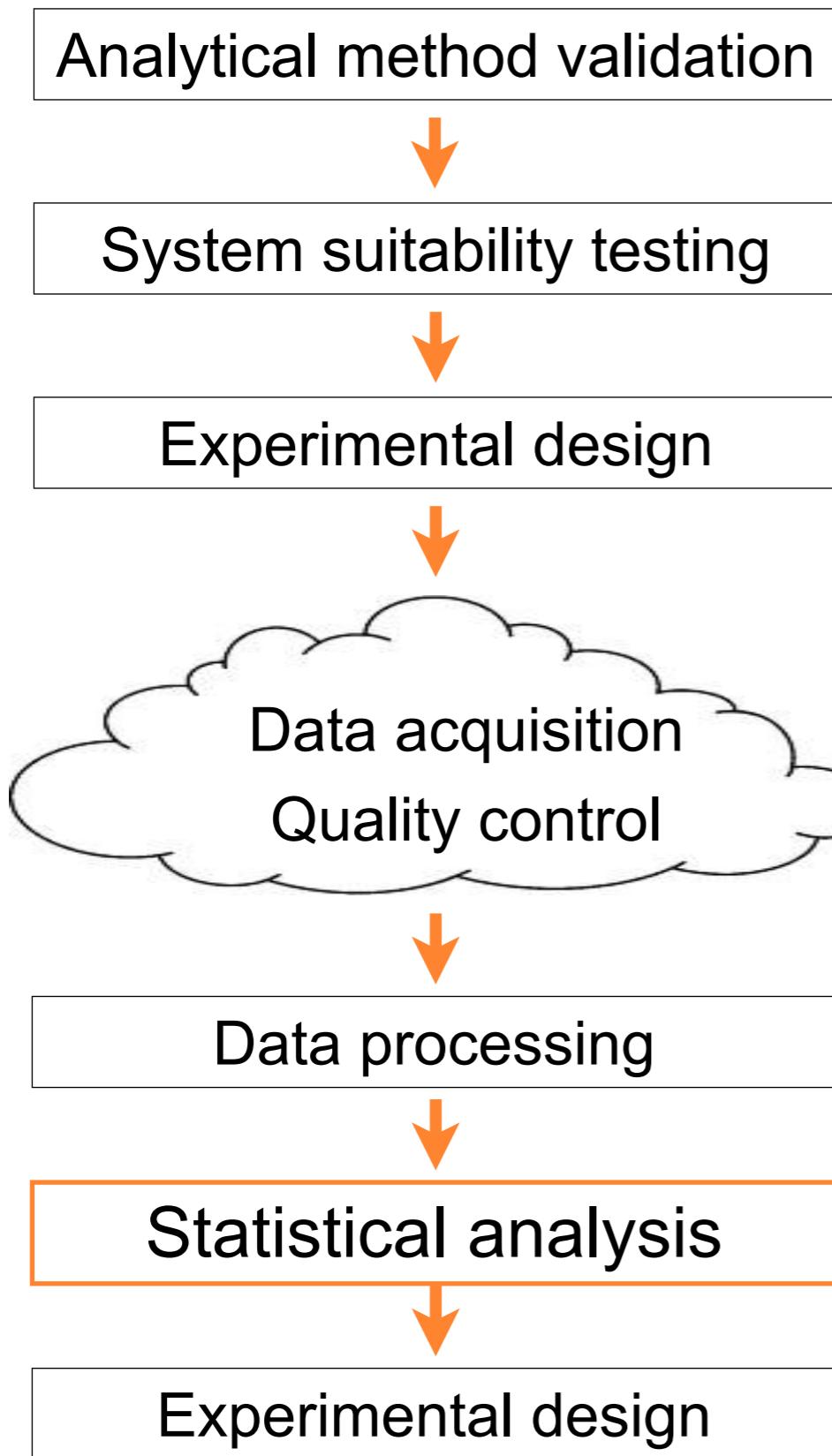
*Endogenous MS runs*

*log<sub>2</sub>(intensity) of all transitions in the run*

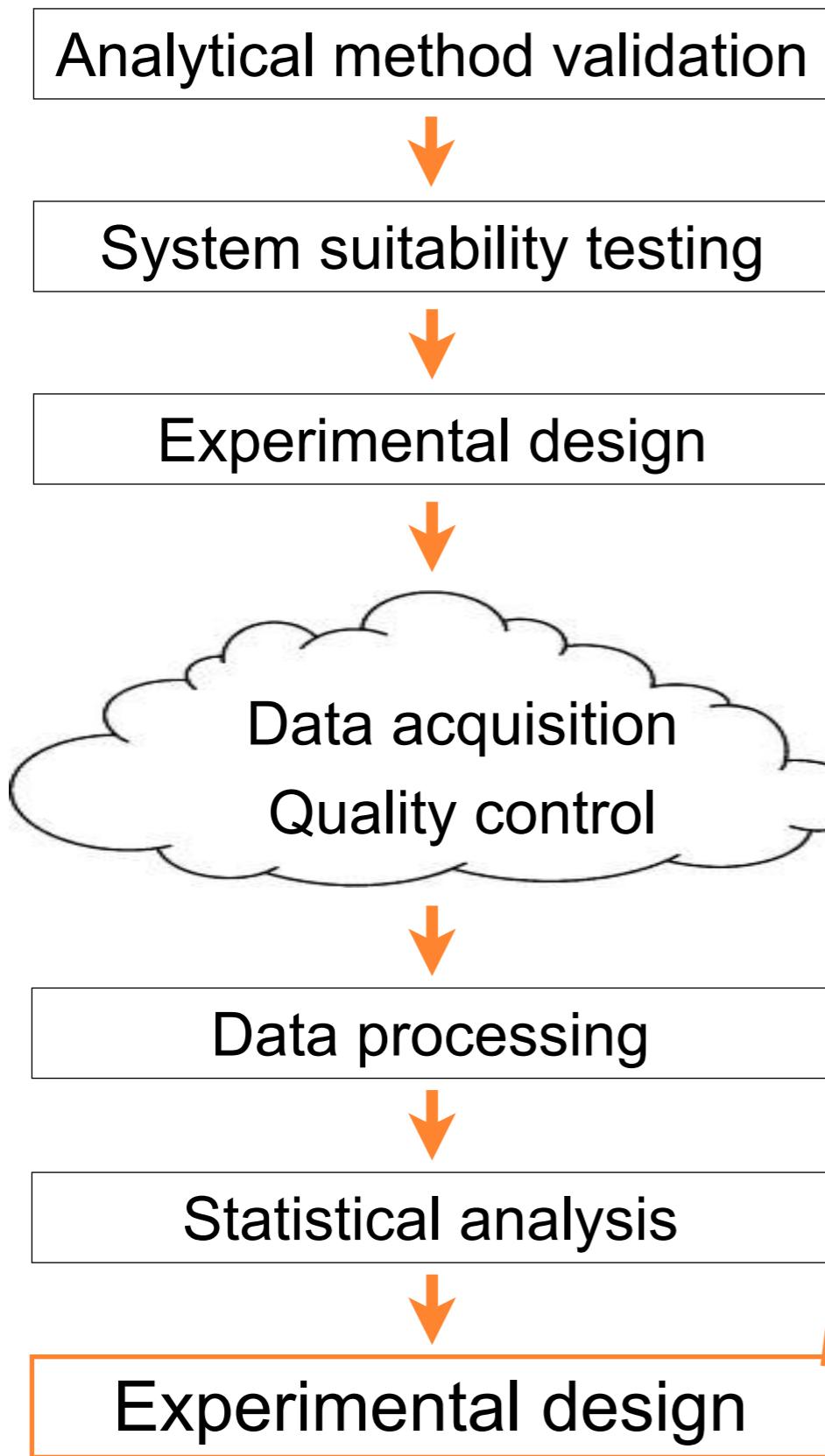


*Reference Endogenous MS runs*

# MS EXPERIMENT: INFERENCE

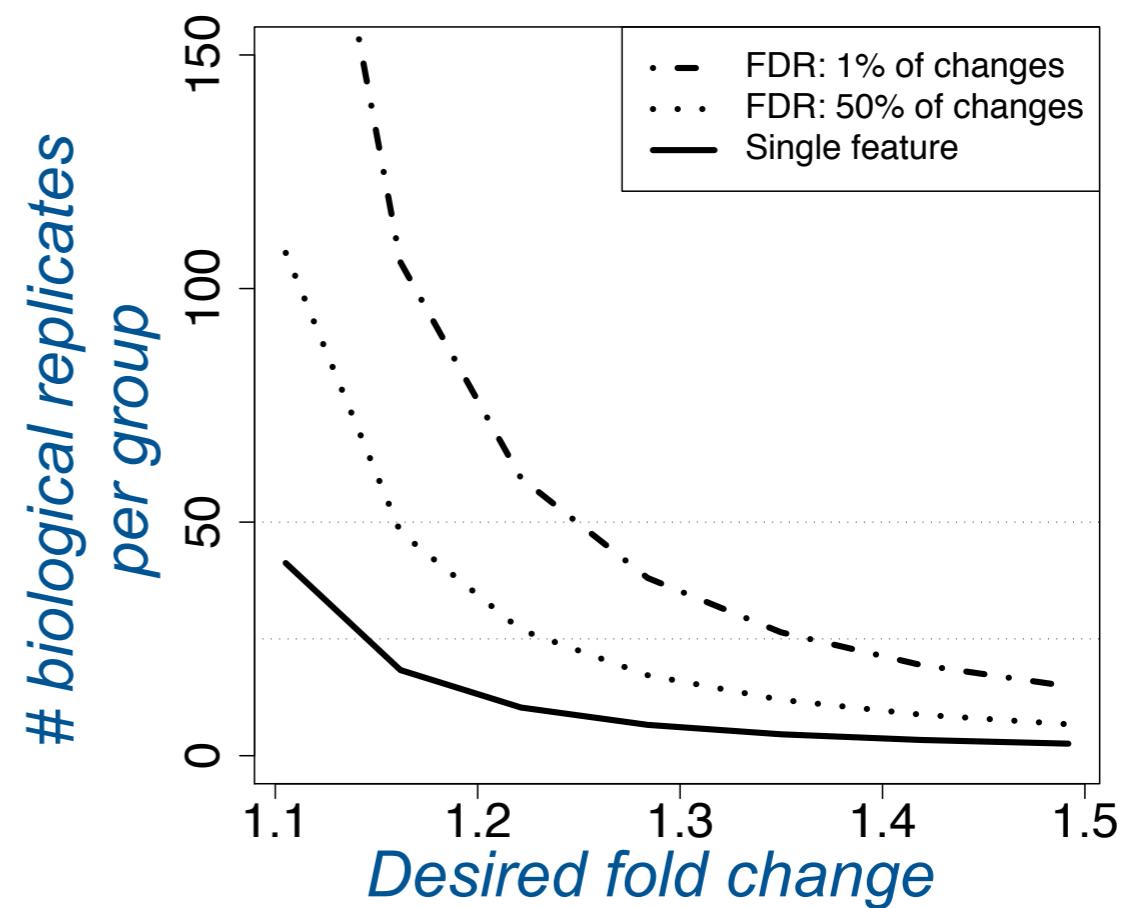


# MS EXPERIMENT: DESIGN



*Use the dataset to improve:*

- Subject selection: matching
- Resource allocation: blocking
- Calculation of sample size

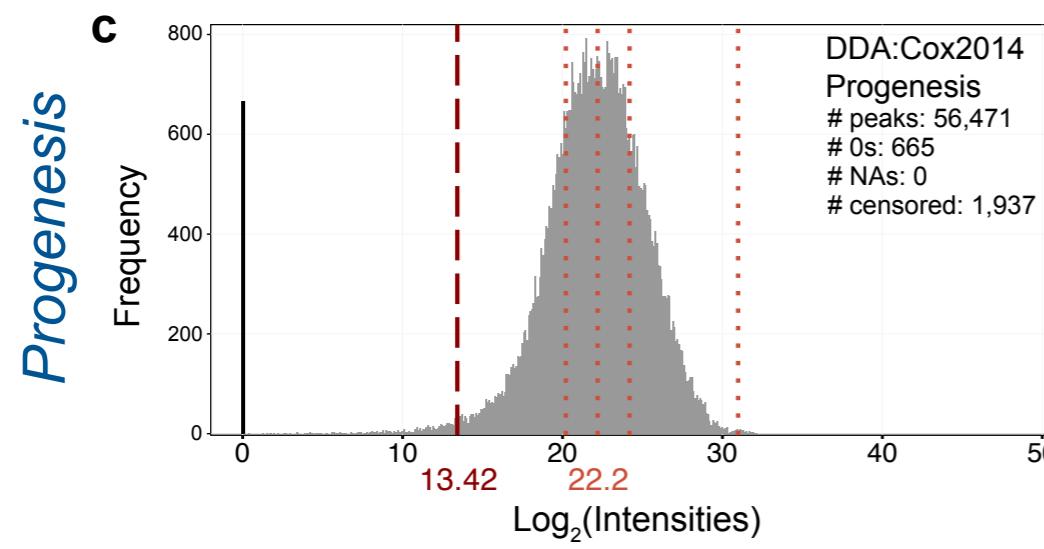
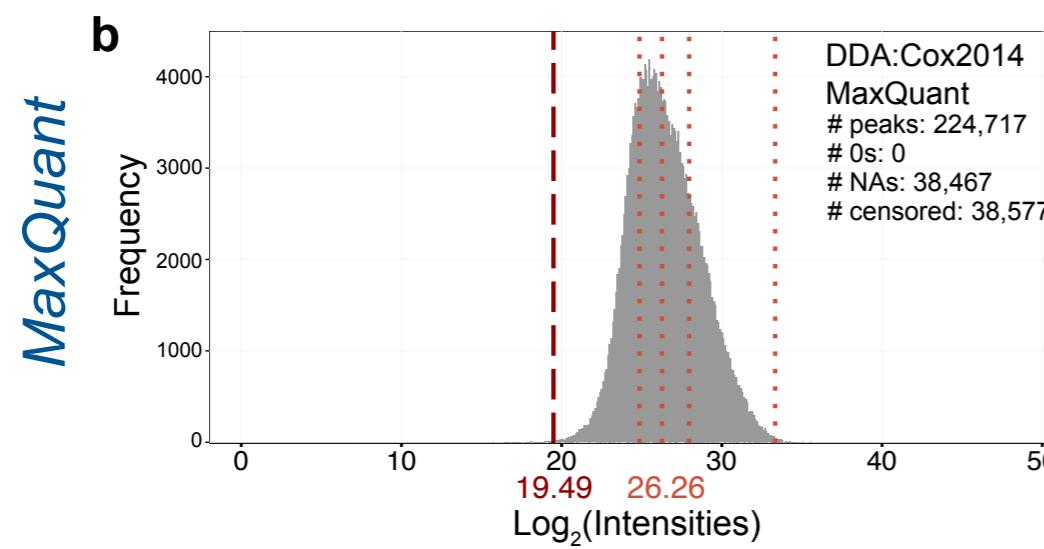
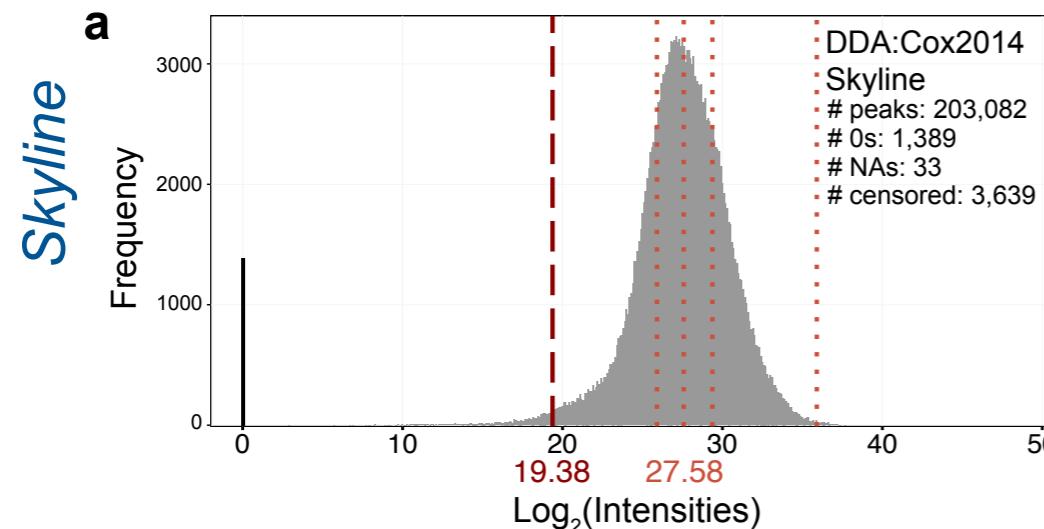


# INTRODUCTION TO MSSTATS

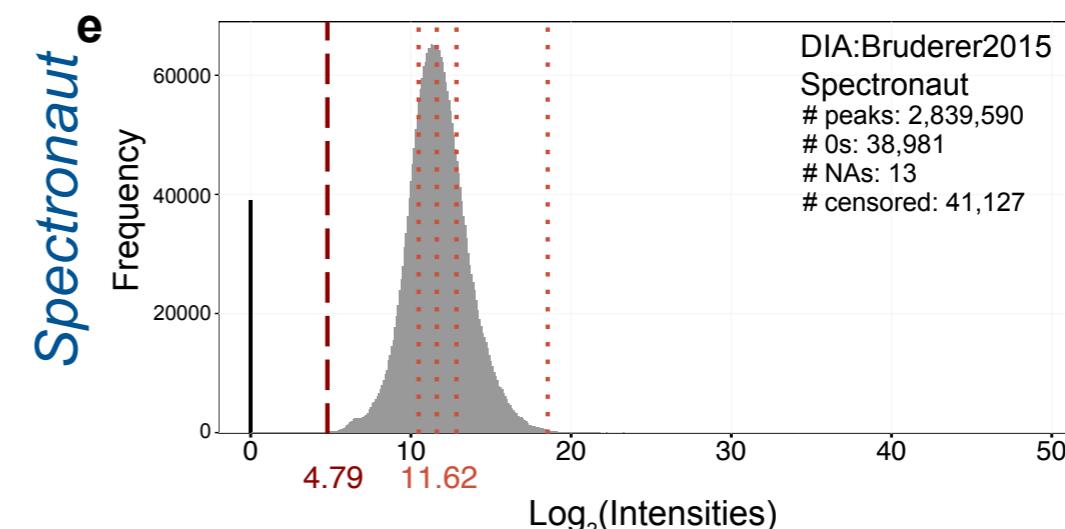
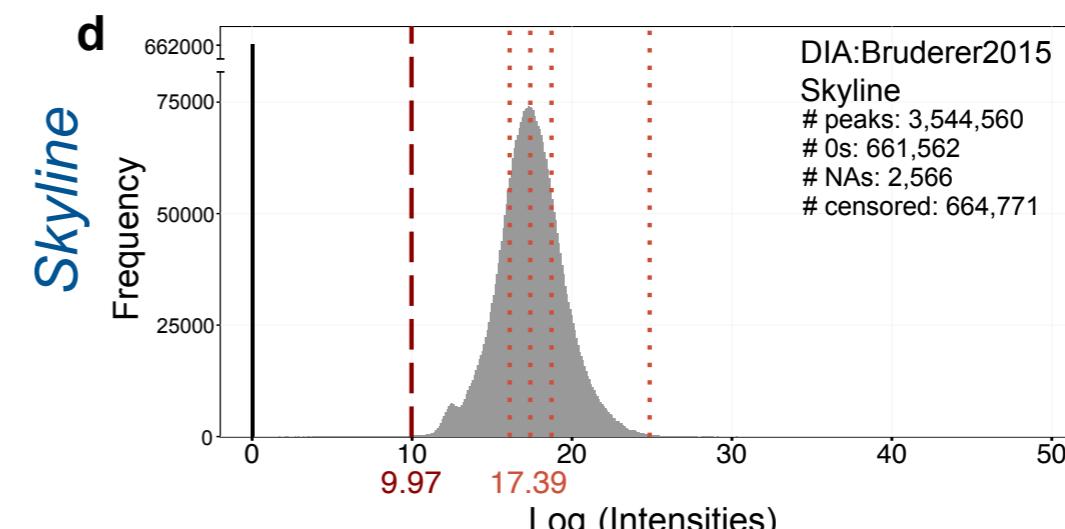
- Overview
- Statistical modeling
- Community resources

# PROPERTIES OF PEAK INTENSITIES VARY BETWEEN DATA PROCESSING TOOLS

*DDA: Cox 2014*



*DIA: Bruderer 2015*



— — —	Estimated censoring threshold
... . .	Quantiles of log <sub>2</sub> (intensity)
— — —	Frequency of peaks with intensity reported as between 0 and 1

# SCHEMATIC DATA REPRESENTATION

17

Repeat for every protein

**Whole plot**

Subplot	Condition <sub>I</sub>						...			Condition <sub>J</sub>																	
	Subject <sub>1</sub>			Subject <sub>2</sub>			...			Subject <sub>J</sub>			...			Subject <sub>(J-1)J+1</sub>			Subject <sub>(J-1)J+2</sub>			...			Subject <sub>IJ</sub>		
	Run <sub>1</sub>	Run <sub>2</sub>	Run <sub>3</sub>	Run <sub>4</sub>	Run <sub>5</sub>	Run <sub>6</sub>	...	Run <sub>JK-2</sub>	Run <sub>JK-1</sub>	Run <sub>JK</sub>	...	Run <sub>(I-1)JK+1</sub>	Run <sub>(I-1)JK+2</sub>	Run <sub>(I-1)JK+3</sub>	Run <sub>(I-1)JK+4</sub>	Run <sub>(I-1)JK+5</sub>	Run <sub>(I-1)JK+6</sub>	...	Run <sub>IJK-2</sub>	Run <sub>IJK-1</sub>	Run <sub>IJK</sub>						
Feature <sub>1</sub>	y	y	y	y	y	y	...	y	y	y	...	y	y	y	y	y	y	...	y	y	y	...	y	y	y		
Feature <sub>2</sub>	y	y	y	y	y	NA	...	y	y	y	...	y	y	y	y	y	y	...	y	y	y	...	y	y	y		
...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...		
Feature <sub>L</sub>	y	Cens	y	Cens	Cens	y	...	Cens	y	y	...	NA	y	y	y	y	y	...	y	Cens	y	...	...	...	...		

*Spectral features:  
technological aspects of the  
experiment*

*Missing values*

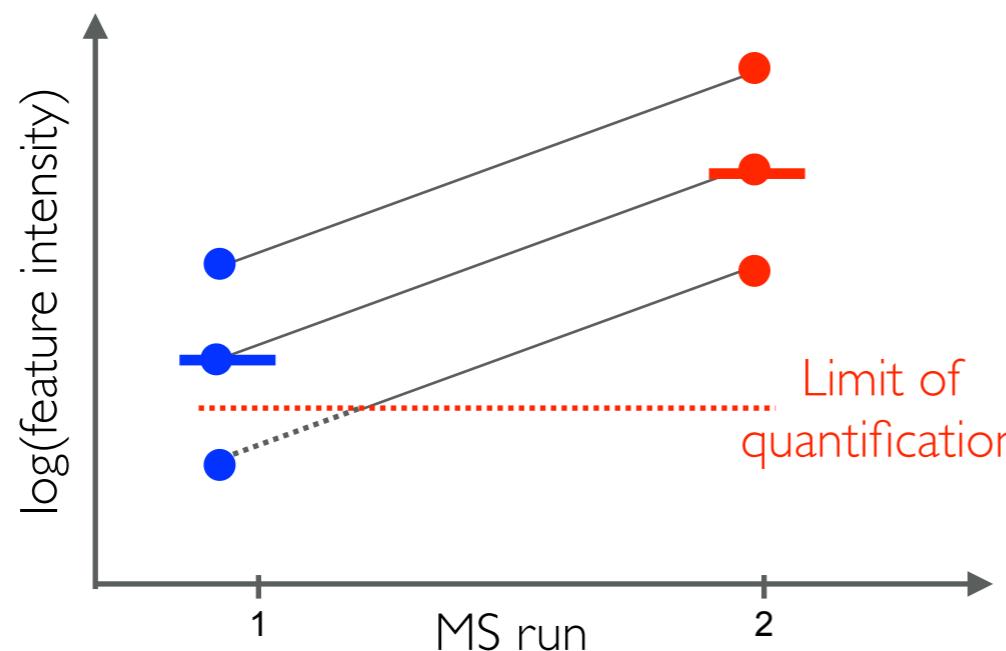
*View feature intensities  
below the threshold as  
censored*

*Log(feature  
intensities)*

## Challenges

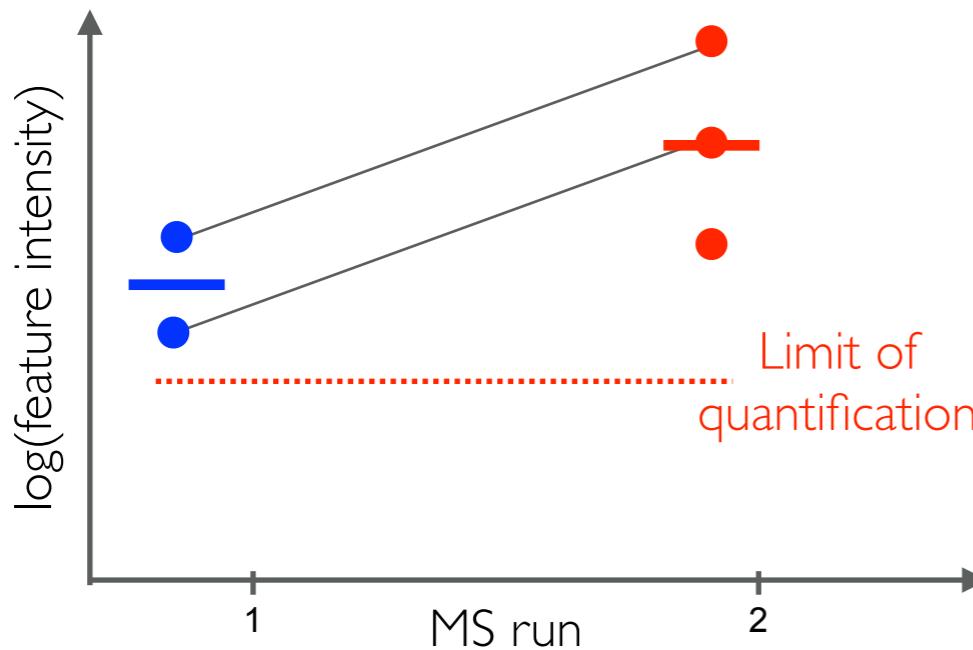
- Represent experimental design
- Represent peak intensities
- Estimate variance components & inference

# IMPUTATION OF MISSING VALUES

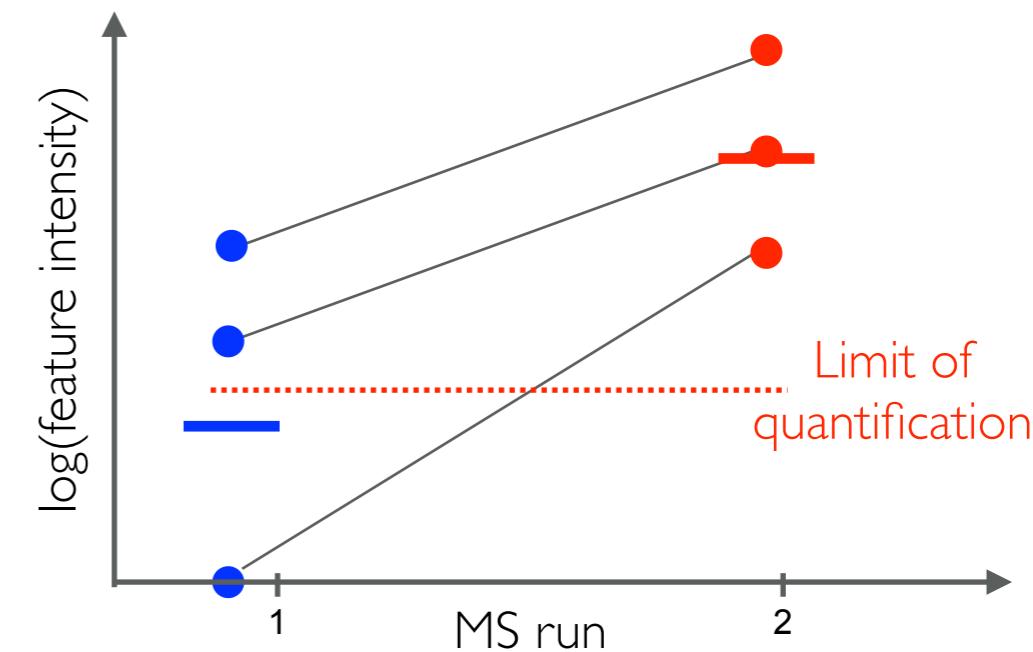


MSstats  
*Relies on information from other features and runs*

No imputation  
*Underestimates the regulation*



Imputation with 0  
*Overestimates the regulation*



# MSSTATS: MODEL AND INFERENCE

A

	Run <sub>1</sub>	Run <sub>2</sub>	Run <sub>3</sub>	Run <sub>4</sub>	Run <sub>5</sub>	Run <sub>6</sub>	...	Run <sub>JK-2</sub>	Run <sub>JK-1</sub>	Run <sub>JK</sub>	...	Run <sub>(I-1)JK+1</sub>	Run <sub>(I-1)JK+2</sub>	Run <sub>(I-1)JK+3</sub>	Run <sub>(I-1)JK+4</sub>	Run <sub>(I-1)JK+5</sub>	Run <sub>(I-1)JK+6</sub>	...	Run <sub>IJK-2</sub>	Run <sub>IJK-1</sub>	Run <sub>IJK</sub>	
Feature <sub>1</sub>	y	y	y	y	y	y	...	y	y	y	...	y	y	y	y	y	y	y	...	y	y	y
Feature <sub>2</sub>	y	y	y	y	y	<b>NA<sub>rand</sub></b>	...	y	y	y	...	y	y	y	y	y	y	...	y	y	y	
...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	
Feature <sub>L</sub>	y	<b>NA<sub>cen</sub></b>	y	<b>NA<sub>cen</sub></b>	<b>NA<sub>cen</sub></b>	y	...	<b>NA<sub>cen</sub></b>	y	y	...	<b>NA<sub>cen</sub></b>	y	y	y	y	y	...	y	<b>NA<sub>cen</sub></b>	y	



Step 1 : Run-level subplot summarization

AFT model : Impute censored missing values by accelerated failure model

$$y_{ijkl} = \mu + Run_{ijk} + Feature_l + \epsilon_{ijkl}, \text{ where } \sum_{ijk} Run_{ijk} = 0, \sum_l Feature_l = 0, \epsilon_{ijkl} \stackrel{iid}{\sim} \mathcal{N}(0, \sigma_\epsilon^2)$$

	Run <sub>1</sub>	Run <sub>2</sub>	Run <sub>3</sub>	Run <sub>4</sub>	Run <sub>5</sub>	Run <sub>6</sub>	...	Run <sub>JK-2</sub>	Run <sub>JK-1</sub>	Run <sub>JK</sub>	...	Run <sub>(I-1)JK+1</sub>	Run <sub>(I-1)JK+2</sub>	Run <sub>(I-1)JK+3</sub>	Run <sub>(I-1)JK+4</sub>	Run <sub>(I-1)JK+5</sub>	Run <sub>(I-1)JK+6</sub>	...	Run <sub>IJK-2</sub>	Run <sub>IJK-1</sub>	Run <sub>IJK</sub>	
Feature <sub>1</sub>	y	y	y	y	y	y	...	y	y	y	...	y	y	y	y	y	y	y	...	y	y	y
Feature <sub>2</sub>	y	y	y	y	y		...	y	y	y	...	y	y	y	y	y	y	...	y	y	y	
...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	
Feature <sub>L</sub>	y	<b>y<sub>imp</sub></b>	y	<b>y<sub>imp</sub></b>	<b>y<sub>imp</sub></b>	y	...	<b>y<sub>imp</sub></b>	y	y	...	<b>y<sub>imp</sub></b>	y	y	y	y	y	...	y	<b>y<sub>imp</sub></b>	y	



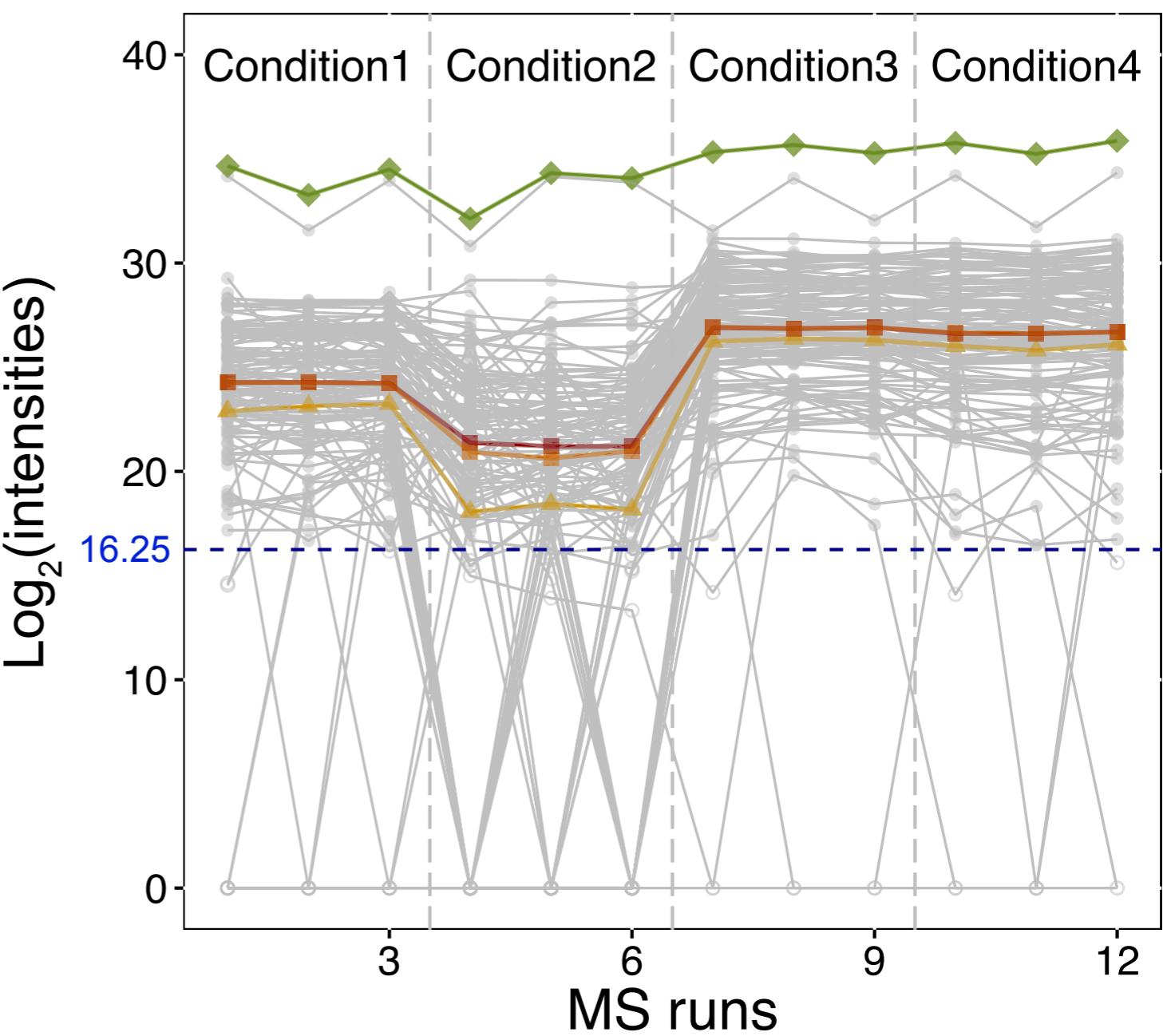
TMP : Parameter estimation by robust method

**Whole plot model**

	Condition <sub>I</sub>						...	Condition <sub>J</sub>						...	Condition <sub>K</sub>						...	Subject <sub>K</sub>	
	Subject <sub>I</sub>		Subject <sub>J</sub>		...	Subject <sub>J</sub>		...	Subject <sub>(I-1)J+1</sub>		Subject <sub>(I-1)J+2</sub>		...	Subject <sub>J</sub>		...	Subject <sub>(I-1)K+1</sub>		Subject <sub>(I-1)K+2</sub>		...	Subject <sub>K</sub>	
	Run <sub>1</sub>	Run <sub>2</sub>	Run <sub>3</sub>	Run <sub>4</sub>	Run <sub>5</sub>	Run <sub>6</sub>	...	Run <sub>JK-2</sub>	Run <sub>JK-1</sub>	Run <sub>JK</sub>	...	Run <sub>(I-1)JK+1</sub>	Run <sub>(I-1)JK+2</sub>	Run <sub>(I-1)JK+3</sub>	Run <sub>(I-1)JK+4</sub>	Run <sub>(I-1)JK+5</sub>	Run <sub>(I-1)JK+6</sub>	...	Run <sub>IJK-2</sub>	Run <sub>IJK-1</sub>	Run <sub>IJK</sub>		
Summarized	$\hat{y}$	$\hat{y}$	$\hat{y}$	$\hat{y}$	$\hat{y}$	$\hat{y}$	...	$\hat{y}$	$\hat{y}$	$\hat{y}$	...	$\hat{y}$	$\hat{y}$	$\hat{y}$	$\hat{y}$	$\hat{y}$	$\hat{y}$	...	$\hat{y}$	$\hat{y}$	$\hat{y}$		

# ROBUSTNESS TO OUTLIERS

*Outliers in both high and low intensities: TMP improves upon linear model and log(sum)*



Condition1-Condition2 : True fold change=7.5

**EstimatedFC Adj.pvalue**

Peptide ions
MSstats
TMP

Linear model

log(sum)

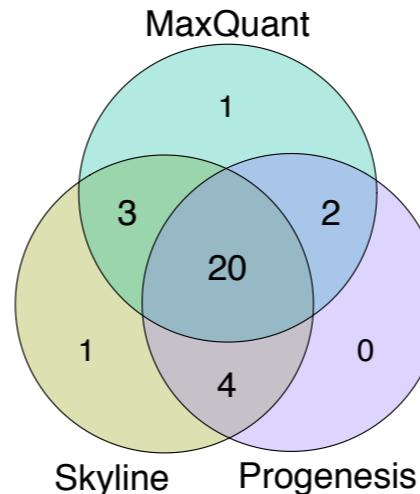
MSstats	8.015	< 0.001
TMP	10.605	< 0.001
Linear model	29.106	< 0.001
log(sum)	1.552	0.999

# STATISTICAL METHODS ENHANCE REPRODUCIBLE RESEARCH

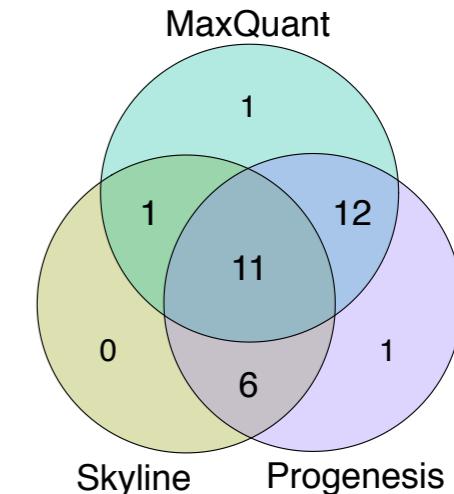
Better agreement  
in correct  
detection of  
differentially  
abundant proteins  
between tools

DDA: *iPRG2015*

*MSstats*

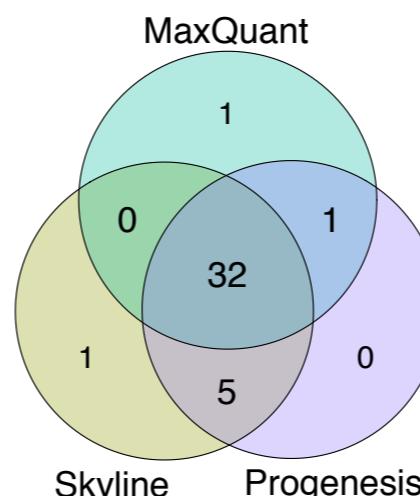


*Log(sum)*

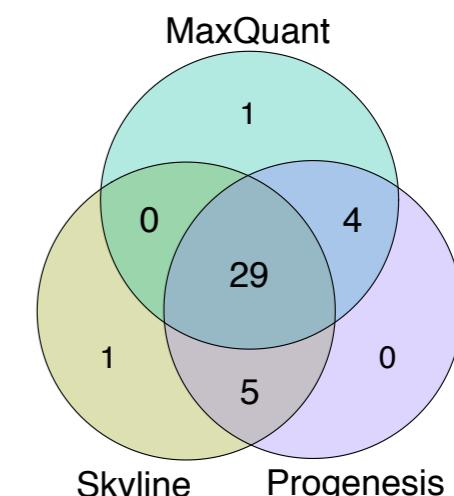


DDA: *Cox 2014*

*MSstats*

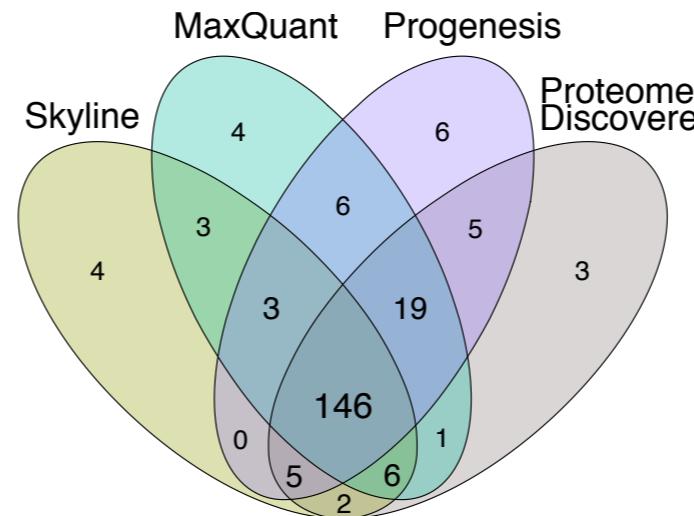


*Log(sum)*

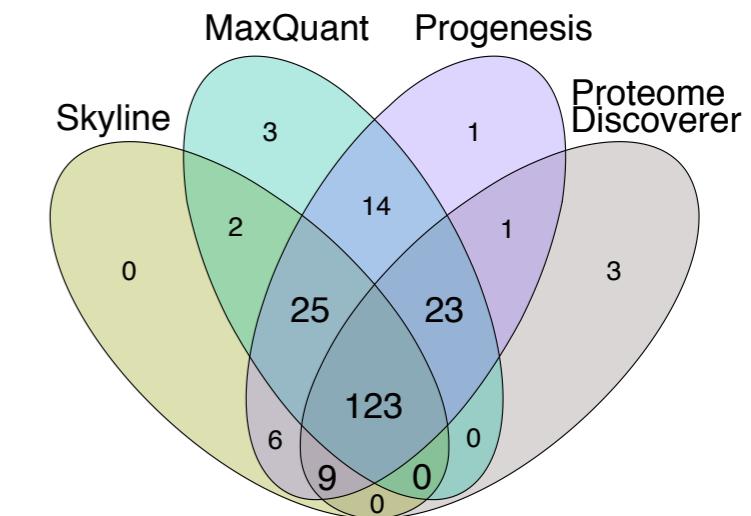


DDA: *Spike-in*

*MSstats*



*Log(sum)*



# INTRODUCTION TO MSSTATS

- Overview
- Statistical modeling
- Community resources

# DATA PROCESSING IMPACTS STATISTICAL ANALYSIS

Publicly available data are key for method development and evaluation

<https://massive.ucsd.edu/ProteoSAFe/static/massive.jsp>

User:  Pass:  Sign in  
Don't have an account? [Register!](#)



nature methods

BRIEF COMMUNICATION

<https://doi.org/10.1038/s41592-020-0955-0>

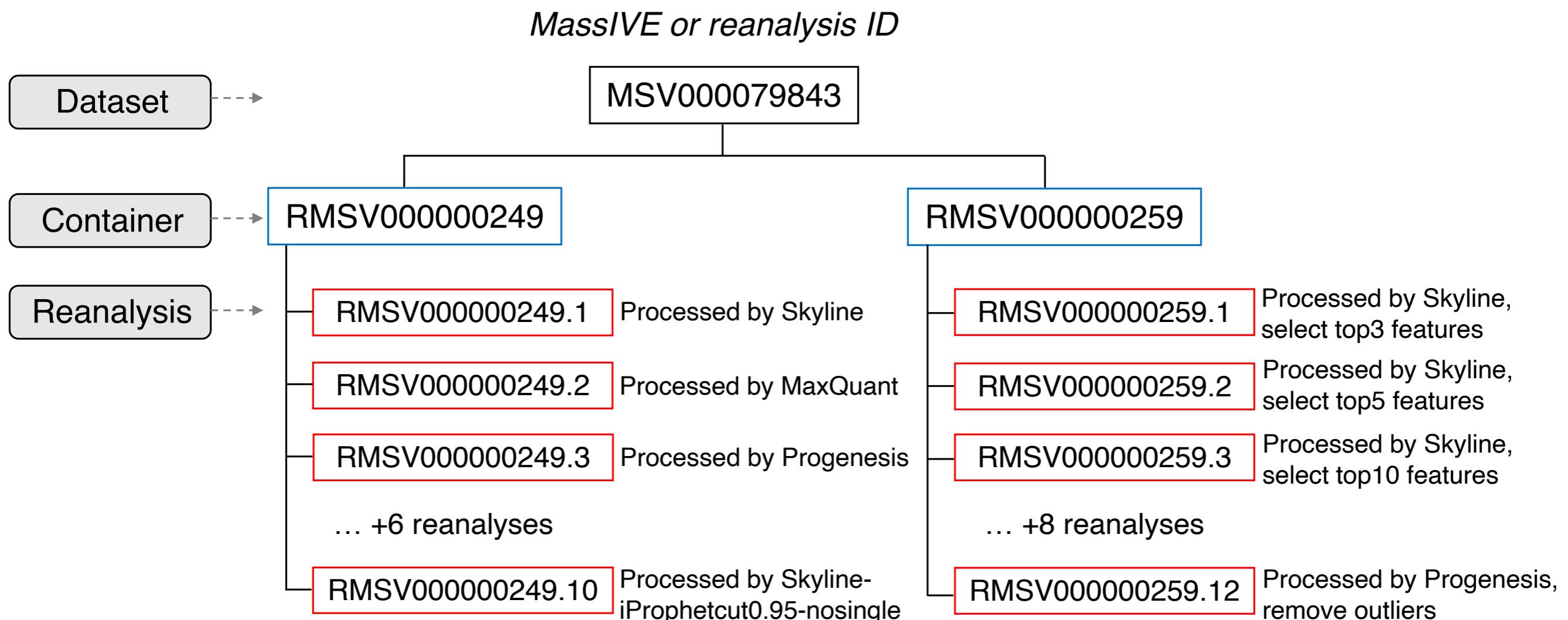
**MassIVE.quant: a community resource of quantitative mass spectrometry-based proteomics datasets**

MassIVE.quant Statistics	
Public Datasets:	<a href="#">46</a>
Quant Reanalyses:	<a href="#">188</a>
Number of Datasets by Species:	
Homo sapiens	<a href="#">23</a>
Mus musculus	<a href="#">12</a>
Saccharomyces cerevisiae	<a href="#">5</a>
Human immunodeficiency virus 1	<a href="#">4</a>
Others	<a href="#">14</a>

Meena Choi<sup>ID 1</sup>, Jeremy Carver<sup>2</sup>, Cristina Chiva<sup>3,4</sup>, Manuel Tzouros<sup>ID 5</sup>, Ting Huang<sup>1</sup>, Tsung-Heng Tsai<sup>ID 1</sup>, Benjamin Pullman<sup>ID 2</sup>, Oliver M. Bernhardt<sup>6</sup>, Ruth Hüttenhain<sup>ID 7</sup>, Guo Ci Teo<sup>ID 8</sup>, Yasset Perez-Riverol<sup>ID 9</sup>, Jan Muntel<sup>ID 6</sup>, Maik Müller<sup>ID 10</sup>, Sandra Goetze<sup>ID 10,11</sup>, Maria Pavlou<sup>10</sup>, Erik Verschueren<sup>7</sup>, Bernd Wollscheid<sup>ID 10,11</sup>, Alexey I. Nesvizhskii<sup>ID 8</sup>, Lukas Reiter<sup>ID 6</sup>, Tom Dunkley<sup>5</sup>, Eduard Sabidó<sup>3,4</sup>, Nuno Bandeira<sup>ID 2</sup>✉ and Olga Vitek<sup>ID 1</sup>✉

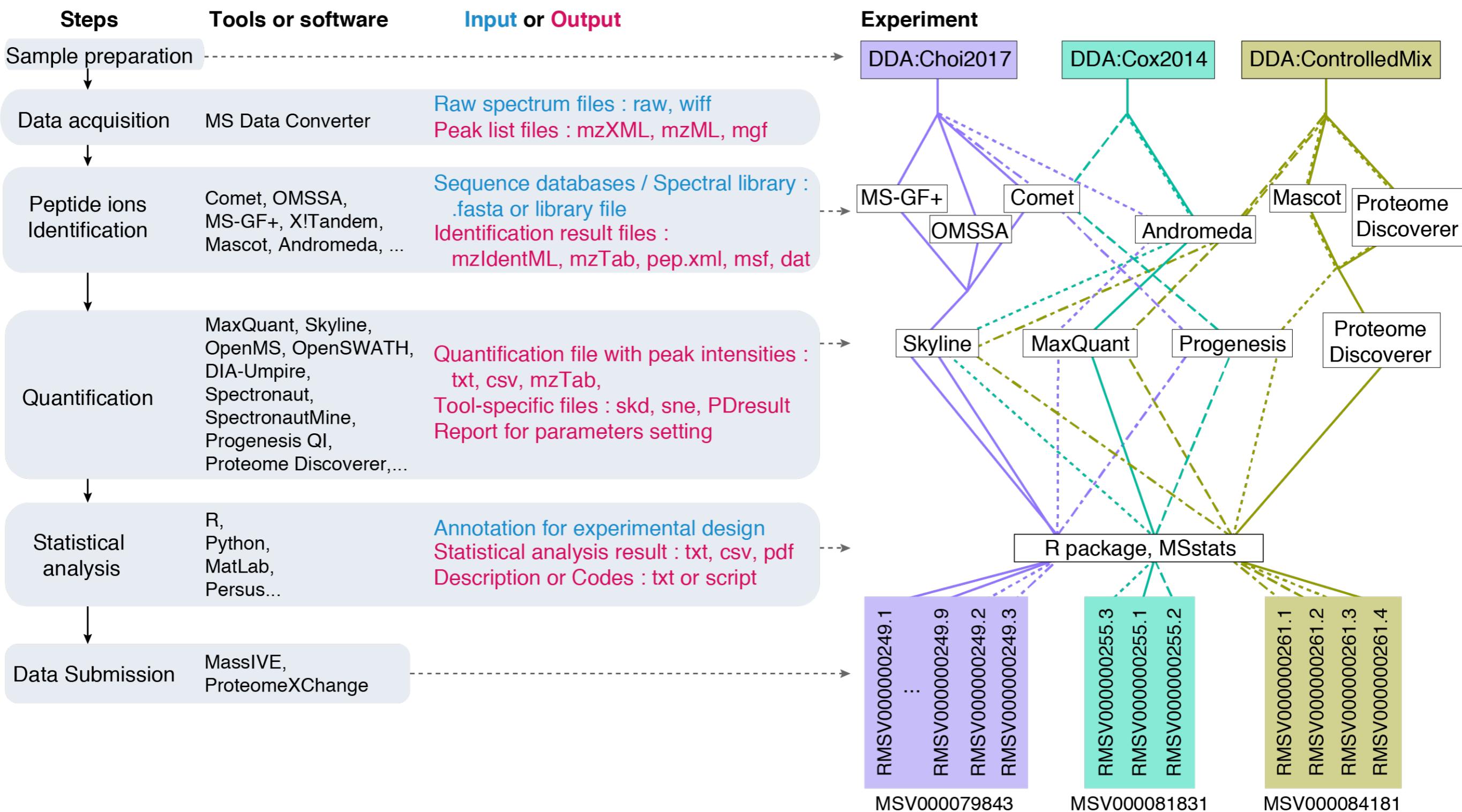
# MASSIVE.QUANT

## Repository for quantitative data and reanalysis



# MASSIVE.QUANT

## Repository for quantitative data and reanalysis



# May Institute goes ONLINE!

## Computation and statistics for mass spectrometry and proteomics

April 27 – May 8, 2020, Northeastern University, Boston MA

Organizers : Meena Choi and Olga Vitek

### Week 1:

April 27, 11:00am- 1:00pm	<b>Alicia Williams:</b> Scientific writing workshop (live)-registration closed	<a href="#">Docs</a> <a href="#">Video</a>
April 28, 11:00am- 12:30pm	<b>Lindsay Pino:</b> Targeted analysis with <a href="#">Skyline</a> , a PRM perspective (live)-registration closed	<a href="#">Docs</a> <a href="#">Video</a>
April 28, 12:45pm- 2:15pm	<b>Brian Searle:</b> Introduction to DIA (live)-registration closed	<a href="#">Docs</a> <a href="#">Video</a>
April 29, 11:00am- 12:30pm	<b>Bernhard Kuster:</b> Large-scale proteomics with TMT (live)-registration closed	<a href="#">Docs</a> <a href="#">Video</a>
April 29, 12:45pm- 2:15pm	<b>Ruedi Aebersold:</b> Modern technology for modern biology (live)-registration closed	<a href="#">Docs</a> <a href="#">Video</a>
April 30, 11:00am- 1:00pm	<b>Brendan MacLean:</b> <a href="#">Skyline</a> for DIA: Intro with a controlled mixture and exploring Prosit spectrum and retention time prediction (video+live)-registration closed	<a href="#">Docs</a> <a href="#">Video</a>
May 1, 11:00am- 1:00pm	<b>Brendan MacLean:</b> Comparing DIA with a deep fractionation library versus a refined subset in <a href="#">Skyline</a> (video+live)-registration closed	<a href="#">Docs</a> <a href="#">Video</a>



Sue  
AbbatIELLO



Ruedi  
Aebersold



Nuno  
Bandeira



Kylie  
Bemis



Meena  
Choi



Laurent  
Gatto



Andy  
Hoofnagle



Oliver  
Kohlbacher



Bernhard  
Kuster



Brendan  
MacLean



Brian  
Searle



Olga  
Vitek

# ACKNOWLEDGEMENTS

## Northeastern University

Kylie Bemis  
Meena Choi  
Eralp Dogu  
Dan Guo  
April Harry  
Ting Huang  
Cyril Galitzine  
Devon Kohler  
Akshay Kulkarni  
Sai Lakkimsetty  
Robert Ness  
Sumedh Sankhe  
Sara Taheri  
Tsung-Heng Tsai

## University of Wroclaw

Mateusz Staniak  
Malgorzata Bogdan

## ETH Zurich

Ruedi Aebersold  
Tiannan Guo  
Ruth Huttenhain  
Paola Picotti  
Silvia Surinova  
Bernd Wollscheid

## University of Washington

Michael MacCoss  
Brendan MacLean  
Jarrett Egertson

## Mugla University

Eralp Dogu



## Hoffman LaRoche

Tom Dunkley  
Balazs Banfa

## Genentech

Don Kirkpatrick  
Erik Verschueren

## UCSD

Nuno Bandeira  
Jeremy Carver

## Support:

NSF  
NIH  
Sternberg Chair  
Canary Center  
Roche  
Genentech  
Eli Lilly  
Chan Zuckerberg foundation

# Thank you to the instructors and to the teaching assistants!

Ryan Benz  
Meena Choi  
Niyati Chopra  
Miguel Cosenza  
Matthias Fahrner  
Amanda Figueroa-Navedo  
Melanie Foell  
Omkar Reddy Gojala  
Dan Guo  
Shubhanshu Gupta  
Ting Huang  
Maanasa Kaza  
Smit Anish Kiri  
Devon Kohler  
Sai Srikanth Lakkimsetty  
Danielle LaMay

Ajeya Makanahalli Kempegowda  
Yogesh Nizzer  
Harish Ramani  
Ruthvik Ravindra  
Abdul Rehman  
Sai Divya Sangeetha Bhagavatula  
Siddarth Sathyanarayanan  
Gopalika Shama  
Rishabh Rajesh Shanbhag  
Sagar Singh  
Mateusz Staniak  
Sara Taheri  
Anuska Tak  
Derrie Susan Varghese  
Amrutha Vempati

Video of the presentation: <https://www.youtube.com/channel/UCnbUMFIIRLaY7fwfSintWuQ/>