

MSSTATS

Meena Choi, Olga Vitek

College of Science

College of Computer and Information Science



Northeastern University

OUTLINE

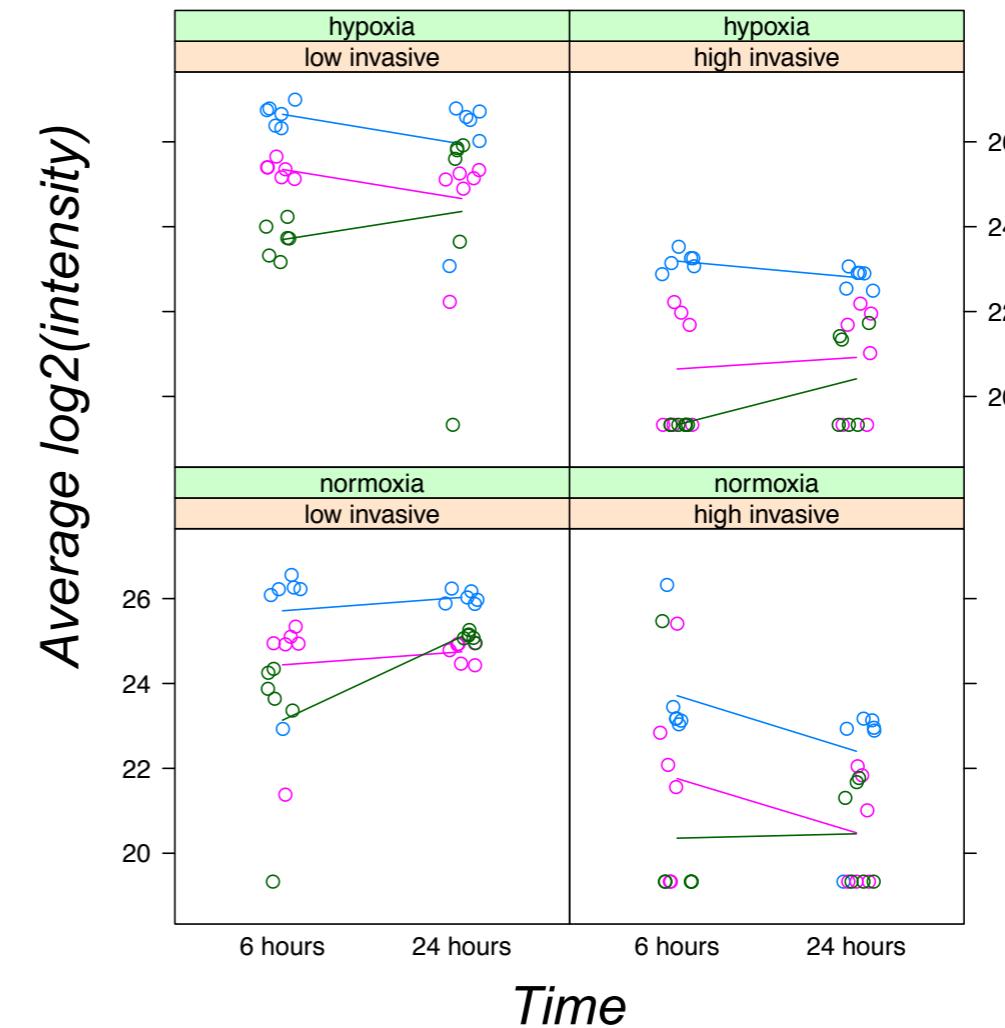
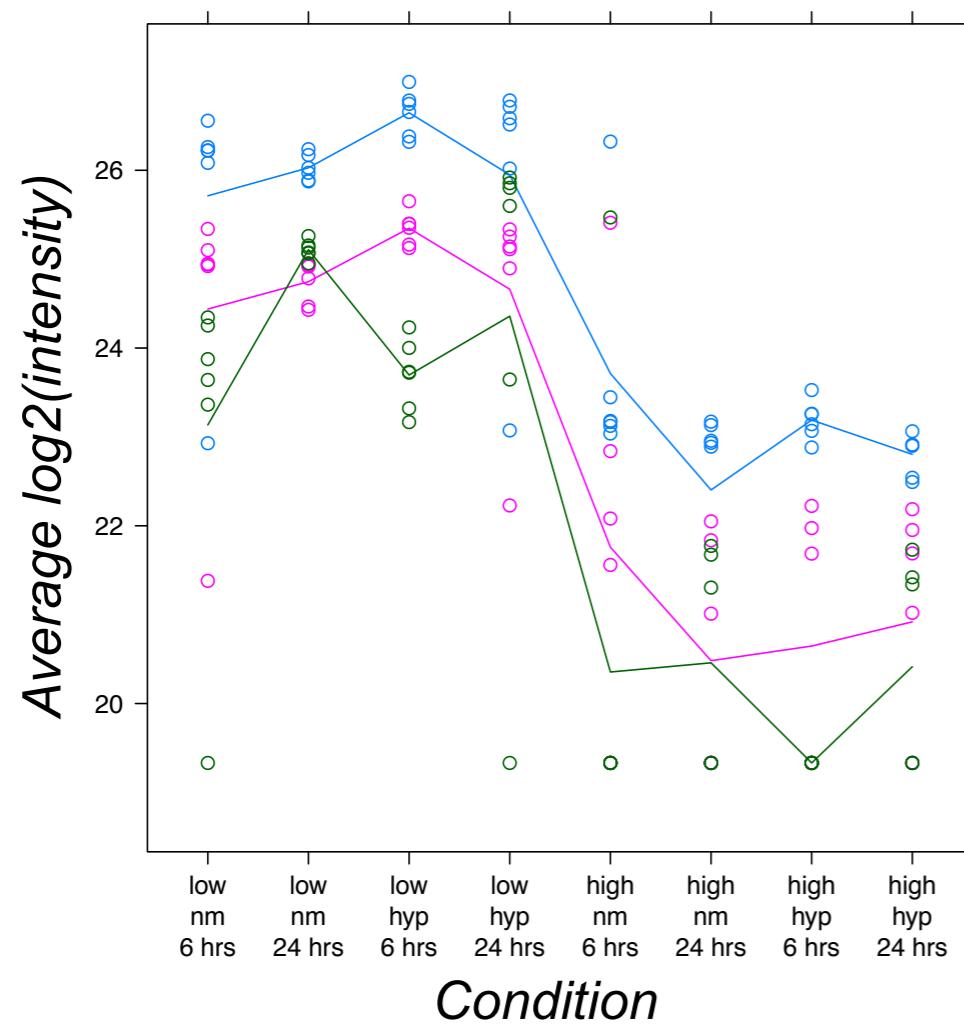
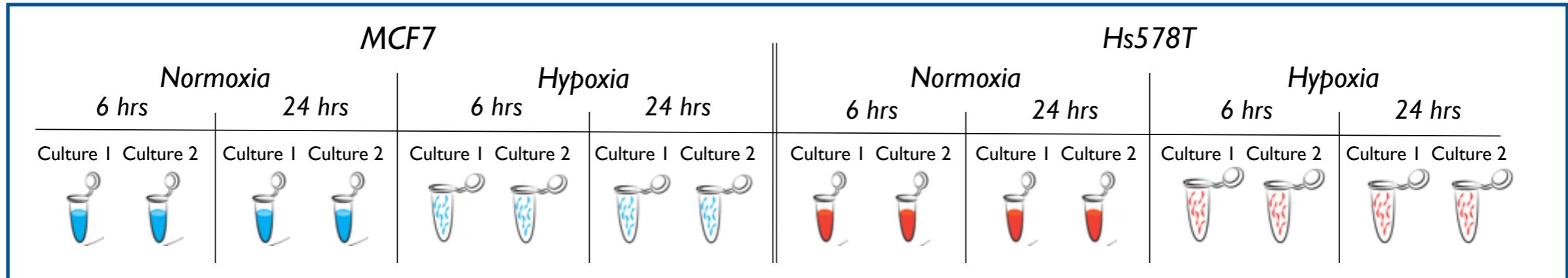
- Methods: background
 - Linear mixed effects models
- Evaluation
 - Spike-in and experimental datasets

MSSTATS

- Statistical relative quantification of proteins and peptides
 - Which protein changes in abundance?
- Complex experimental designs
 - Multiple conditions, factorial experiments, paired designs, time course
- Chromatography-based quantification
 - Shotgun DDA, targeted SRM, data independent DIA/SWATH
- Label-free or label-based
 - Simple summaries and models
- Multiple functionalities
 - Data visualization, statistical modeling and inference, sample size
- Free, open-source and inter-operable with other tools
 - External tool for Skyline, converter from MaxQuant

EXAMPLE: A LABEL-FREE EXPERIMENT

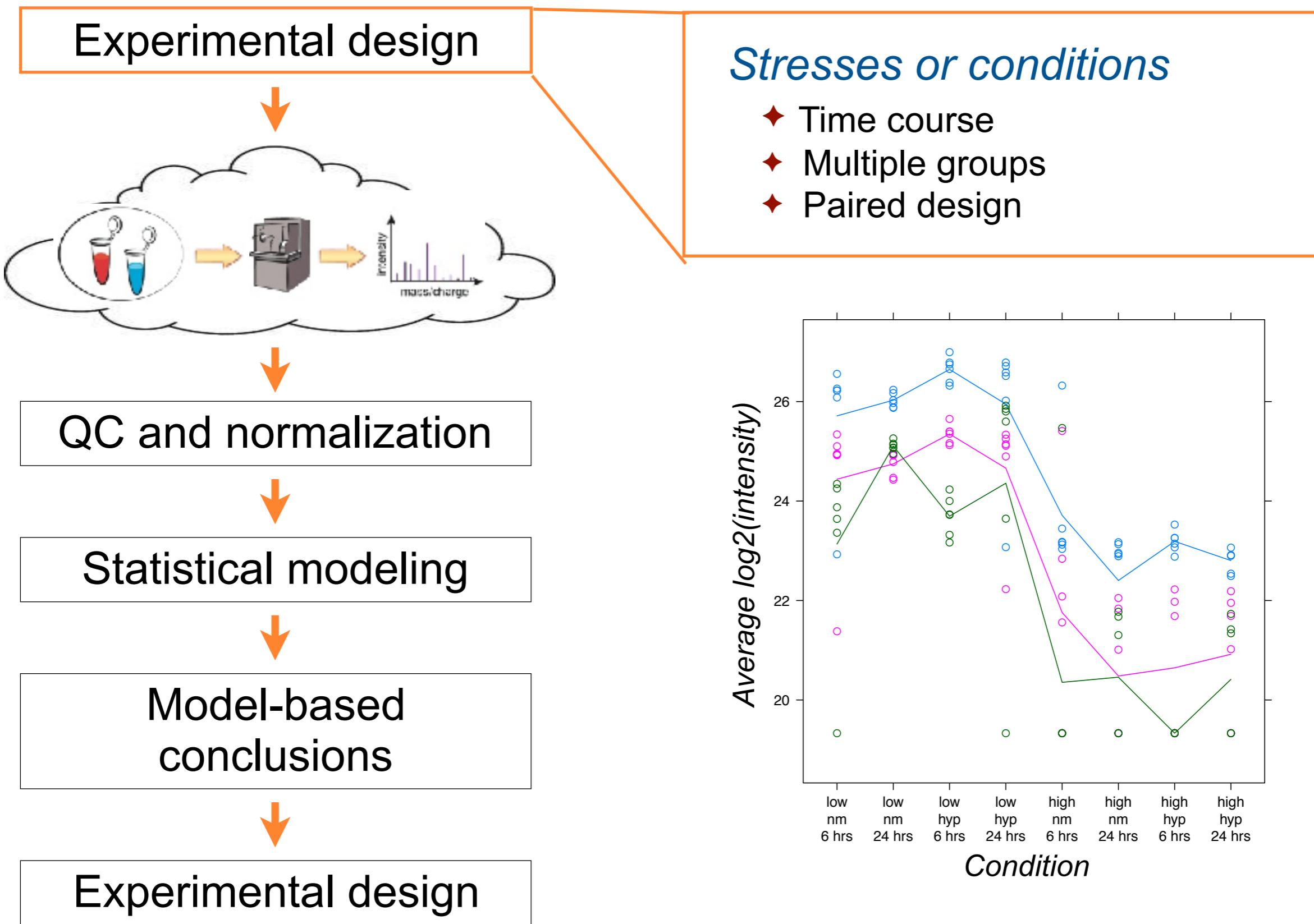
Question: which proteins change in abundance?



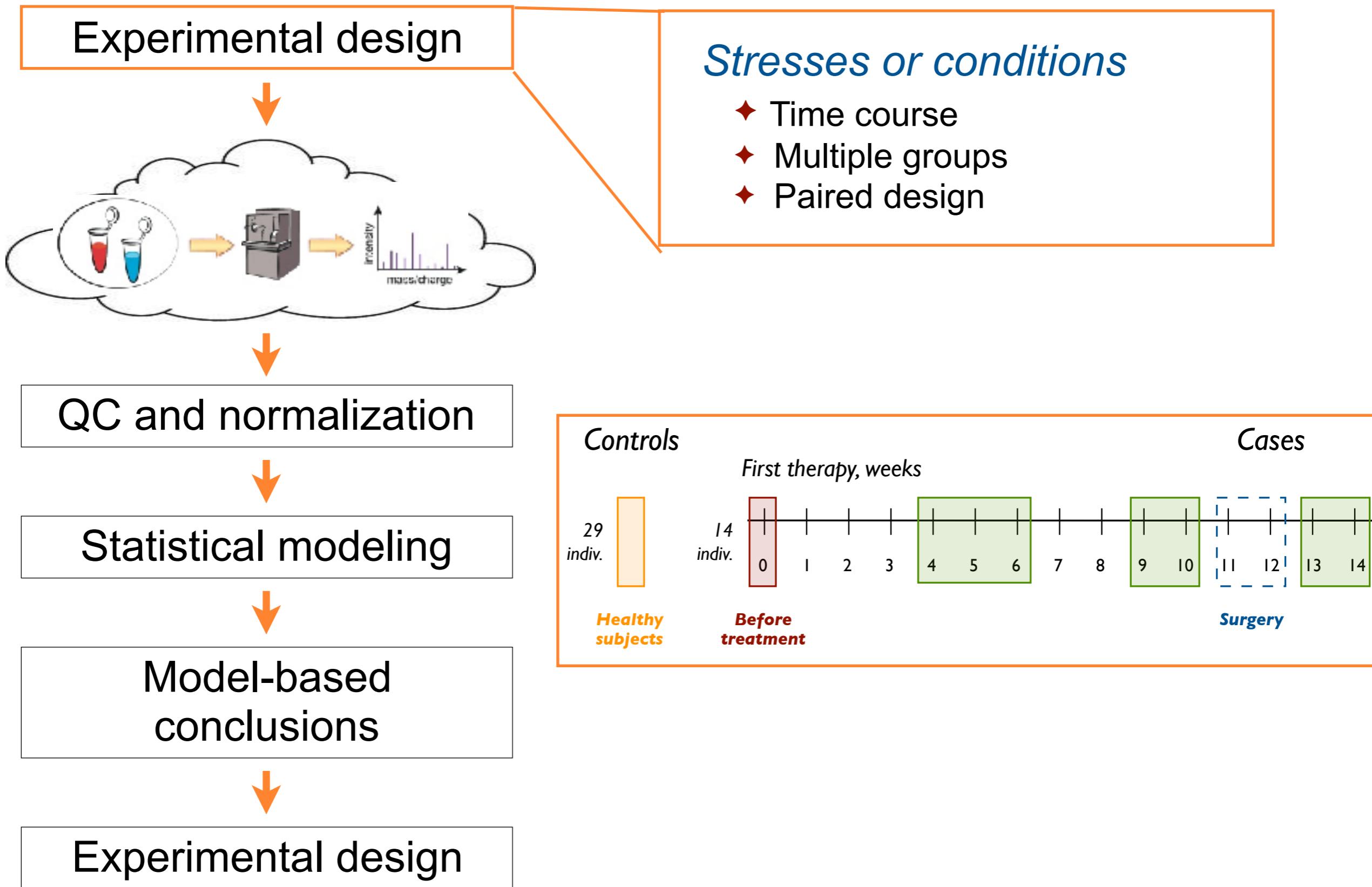
INPUT DATA

	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

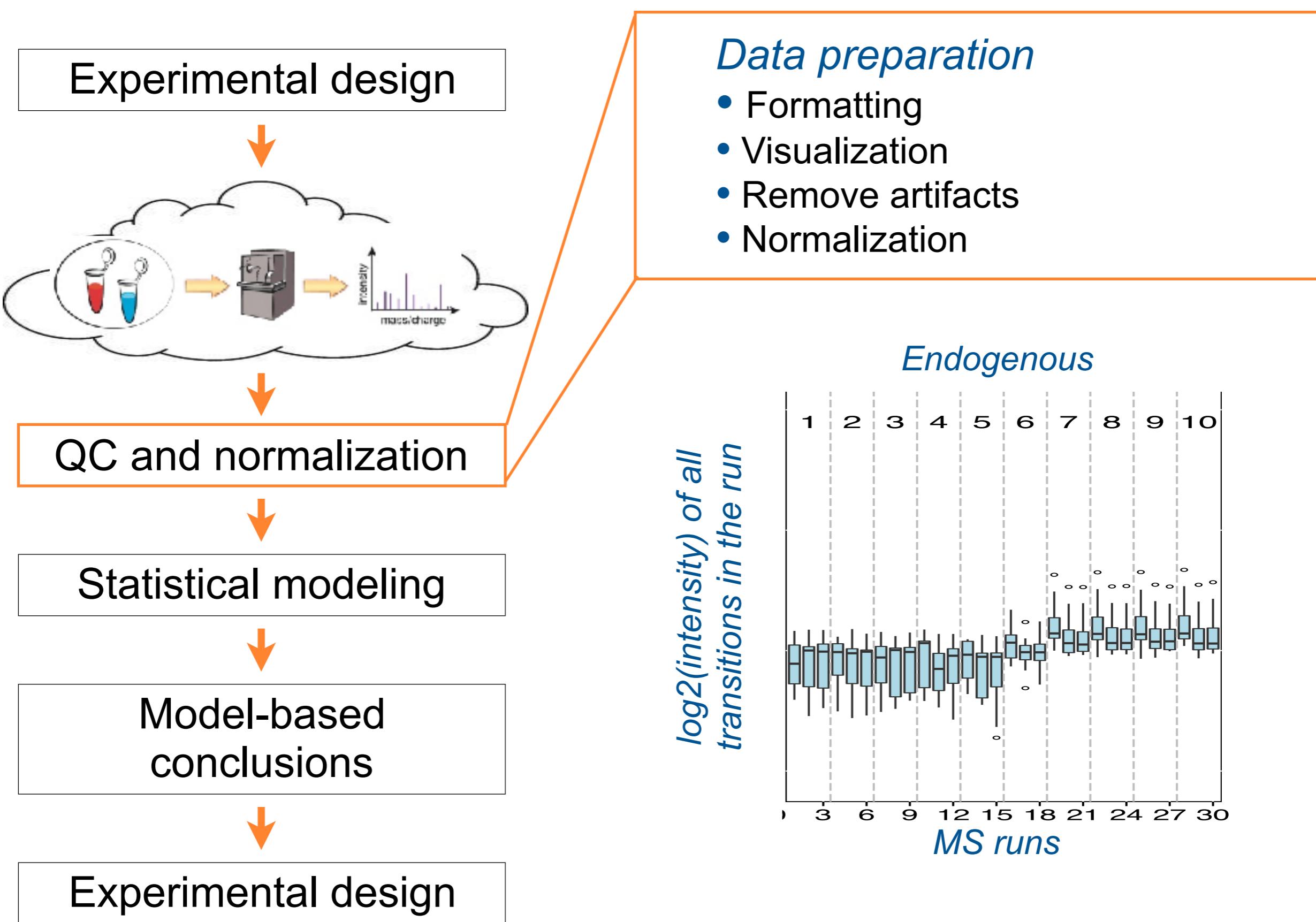
A TYPICAL ANALYSIS WORKFLOW



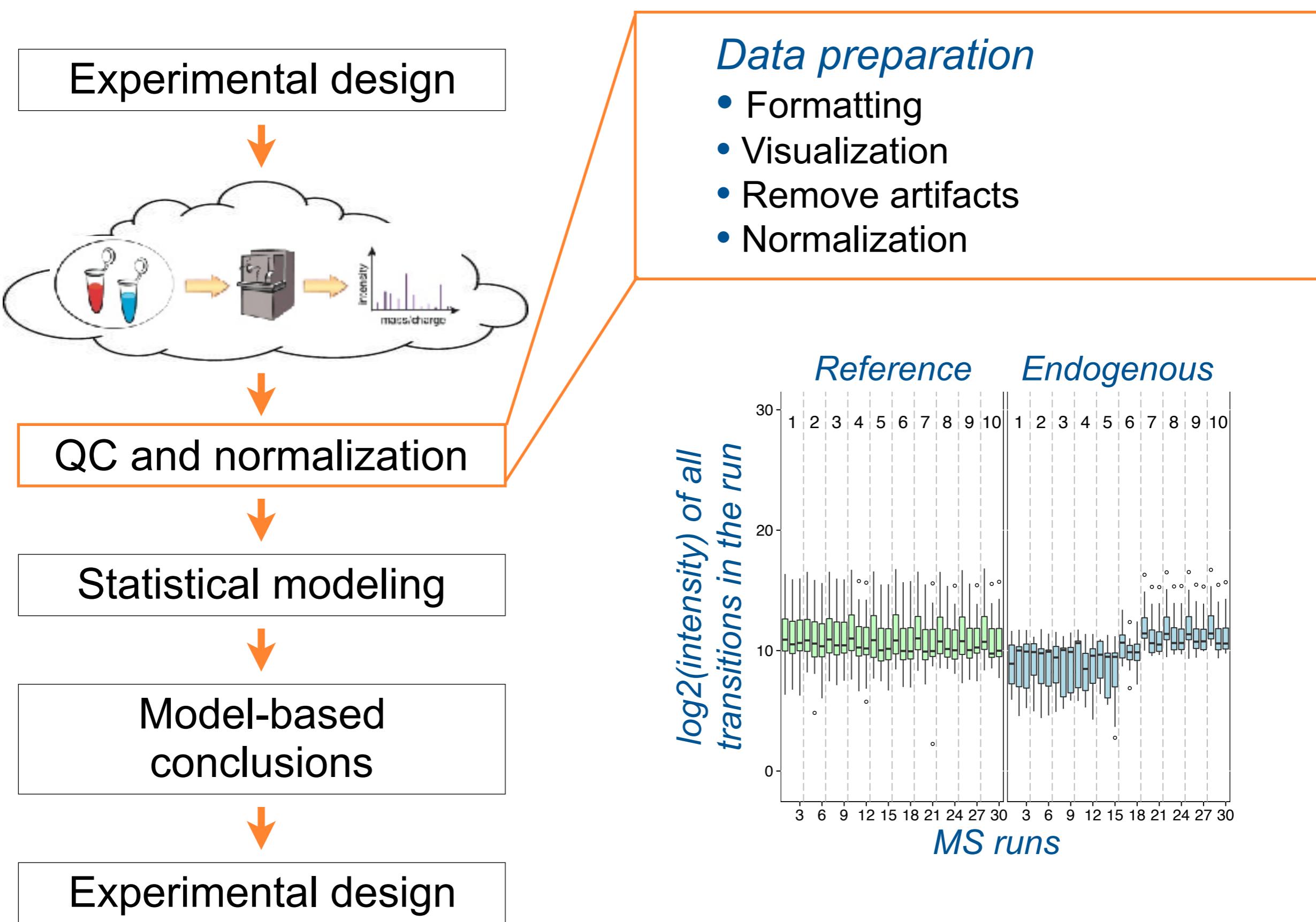
A TYPICAL ANALYSIS WORKFLOW



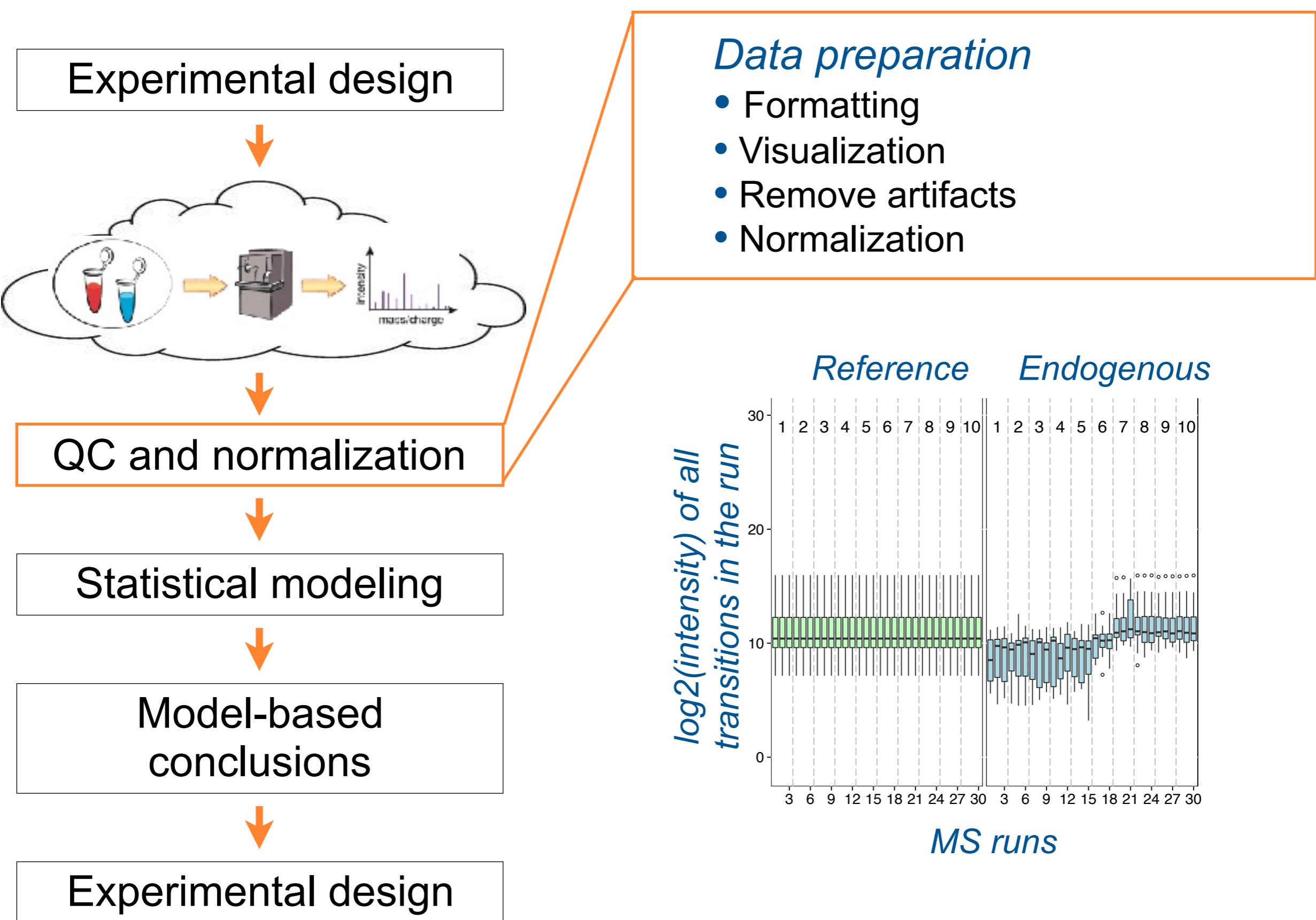
A TYPICAL ANALYSIS WORKFLOW



A TYPICAL ANALYSIS WORKFLOW

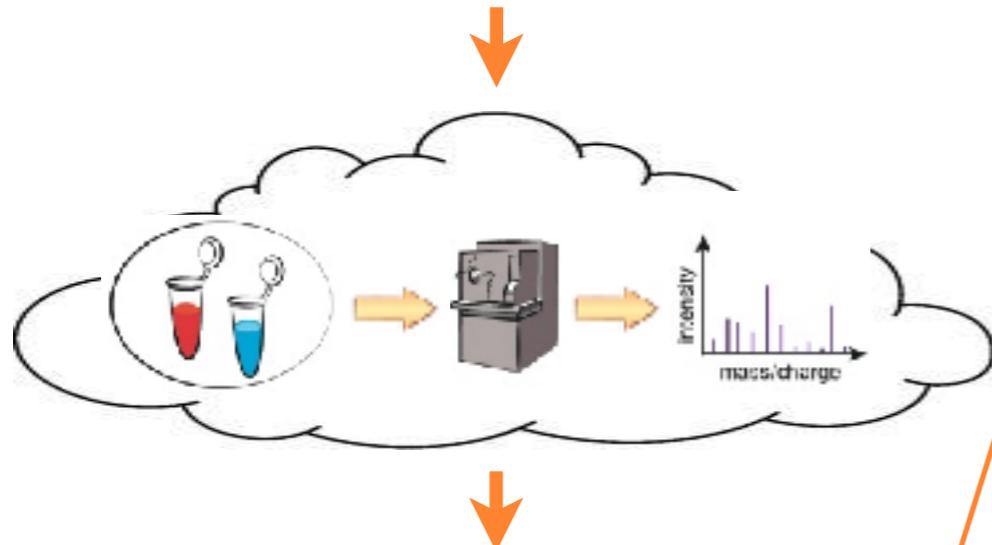


A TYPICAL ANALYSIS WORKFLOW



A TYPICAL ANALYSIS WORKFLOW

Experimental design

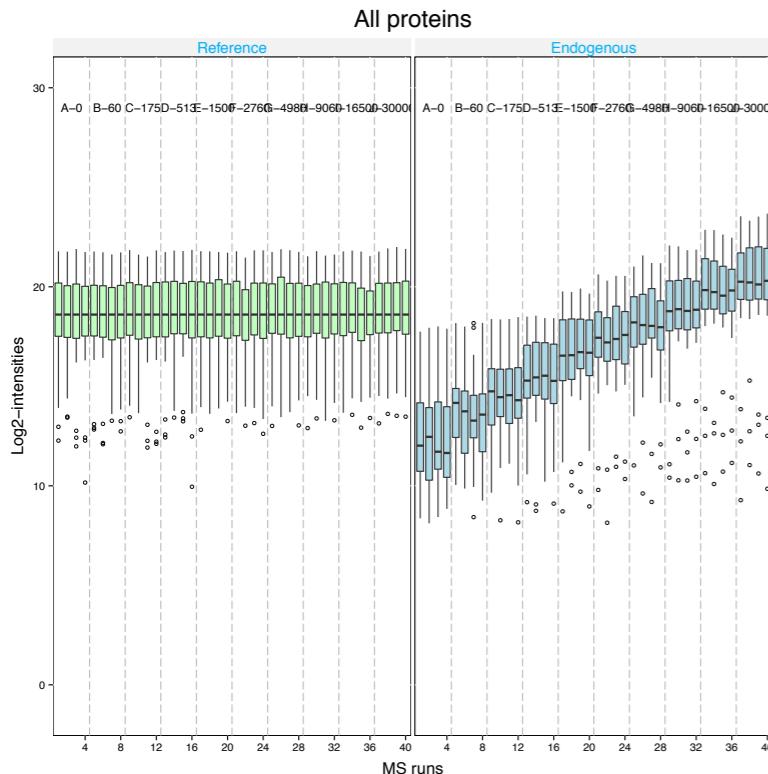


QC and normalization

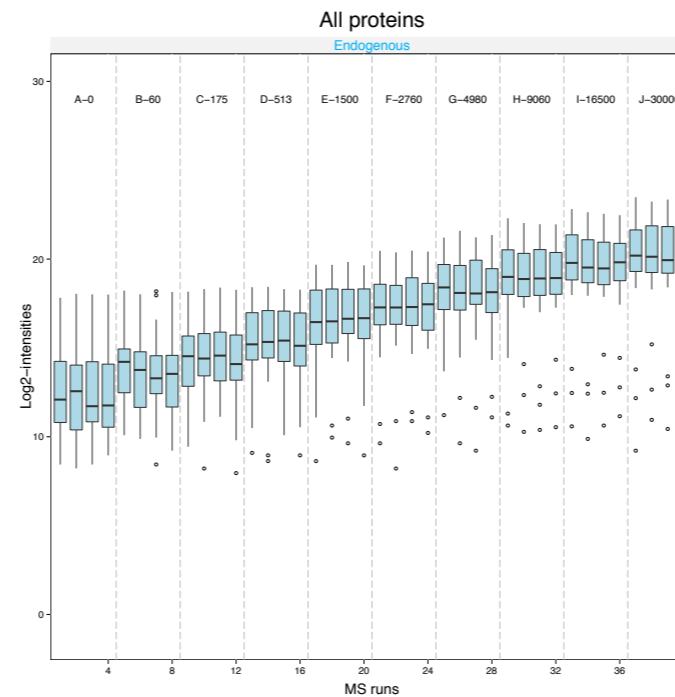
Data preparation

- Formatting
- Visualization
- Remove artifacts
- Normalization

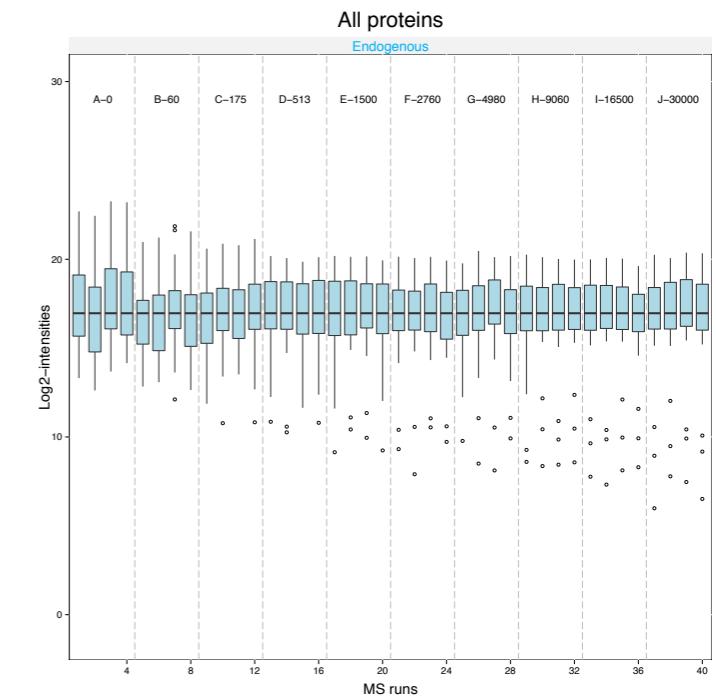
Equalize medians normalization



No normalization

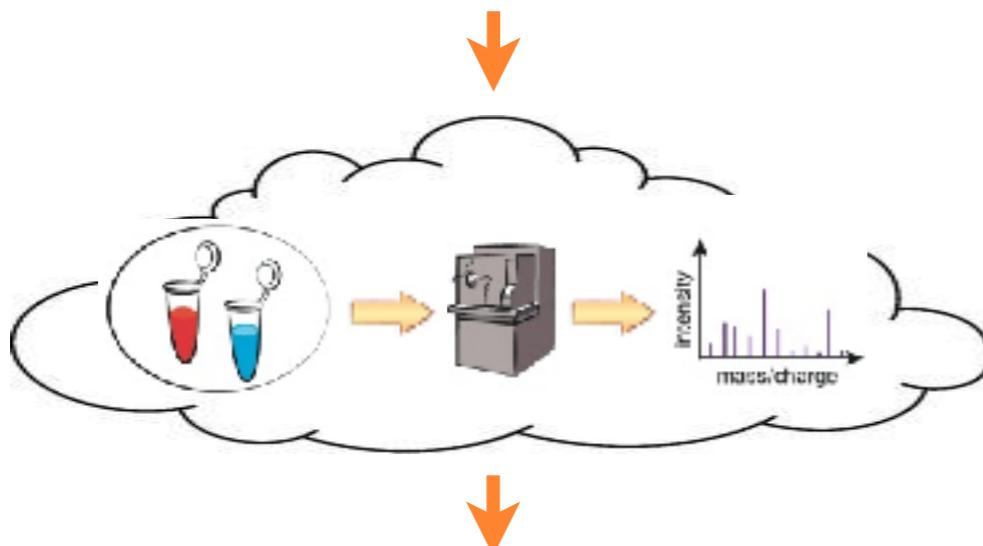


Equalize medians normalization



A TYPICAL ANALYSIS WORKFLOW

Experimental design



QC and normalization

Statistical modeling

Model-based conclusions

Experimental design

Summarize all protein features in a statistical model

- Systematic variation
- Random variation

Verify the assumptions!

- Describe statistical properties of
 - experimental design
 - biological variation
 - measurement technology

LINEAR MIXED MODELS

A split plot approach

Whole plot

Subplot	Condition ₁						...	Condition _I													
	Subject ₁			Subject ₂			...	Subject _J			...	Subject _{(I-1)J+1}			Subject _{(I-1)J+2}			...	Subject _{IJ}		
	Run ₁	Run ₂	Run ₃	Run ₄	Run ₅	Run ₆	...	Run _{JK-2}	Run _{JK-1}	Run _{JK}	...	Run _{(I-1)JK+1}	Run _{(I-1)JK+2}	Run _{(I-1)JK+3}	Run _{(I-1)JK+4}	Run _{(I-1)JK+5}	Run _{(I-1)JK+6}	...	Run _{IJK-2}	Run _{IJK-1}	Run _{IJK}
Feature ₁	y	y	y	y	y	y	...	y	y	y	...	y	y	y	y	y	y	...	y	y	y
Feature ₂	y	y	y	y	y	NA	...	y	y	y	...	y	y	y	y	y	y	...	y	y	y
...
Feature _L	y	NA	y	NA	NA	y	...	NA	y	y	...	NA	y	y	y	y	y	...	y	NA	y

Whole plot

Subplot

$$y_{ijkl} = \mu + \text{Condition}_i + \text{Subject(Condition)}_{j(i)} + \text{Run}_{ijk} + \text{Feature}_l + \text{Run} \times \text{Feature}_{ijkl}$$

Whole-plot
biological variation Whole-plot
technical variation Subplot
error

where $\sum_{i=1}^I \text{Condition}_i = 0$, $\sum_{j=1}^L \text{Feature}_l = 0$

$$\text{Subject(Condition)}_{j(i)} \stackrel{\text{iid}}{\sim} \mathcal{N}(\mathbf{0}, \sigma_{\text{Subject}}^2)$$

$$\text{Run}_{ijk} = \psi_{ijk} \stackrel{\text{iid}}{\sim} \mathcal{N}(\mathbf{0}, \sigma_{\psi}^2)$$

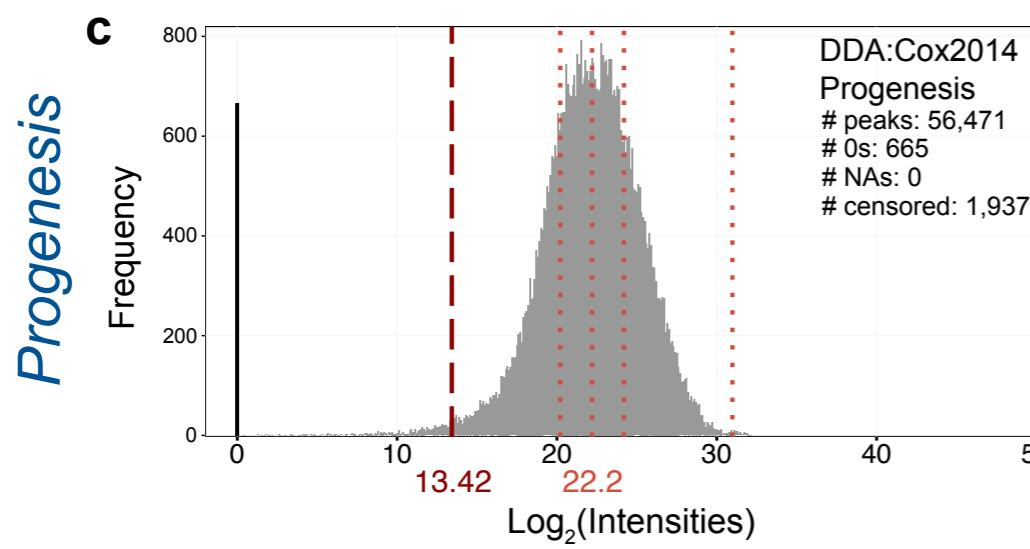
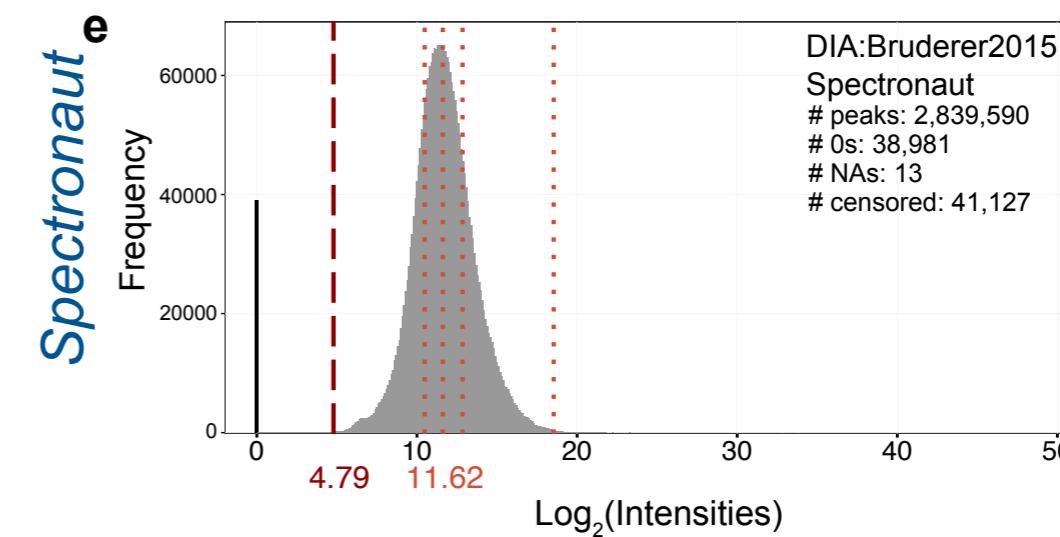
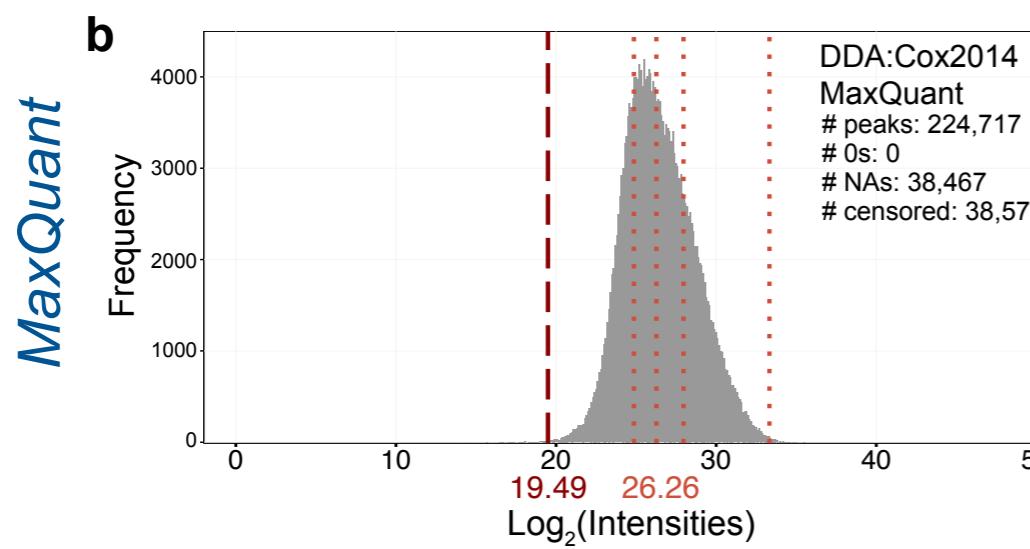
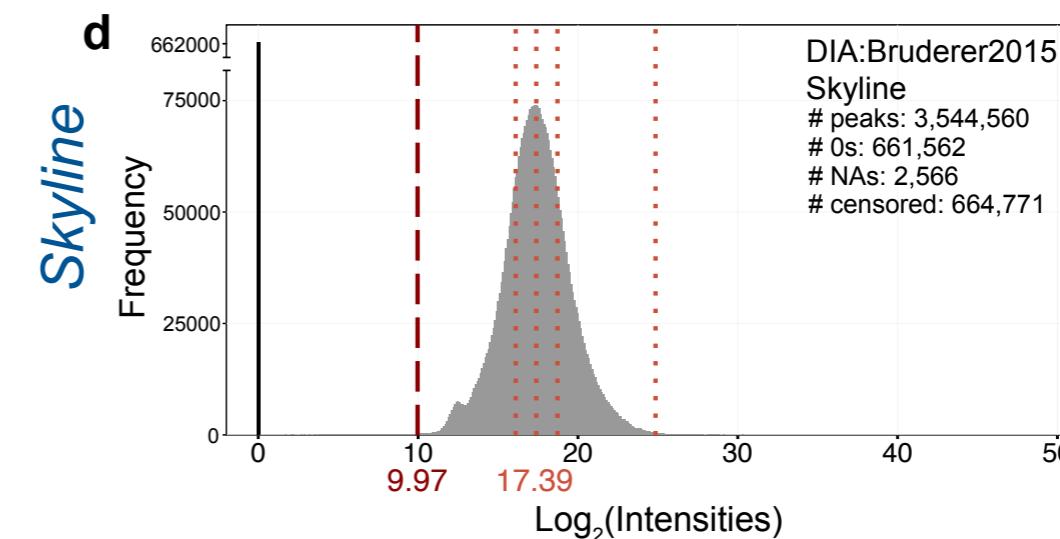
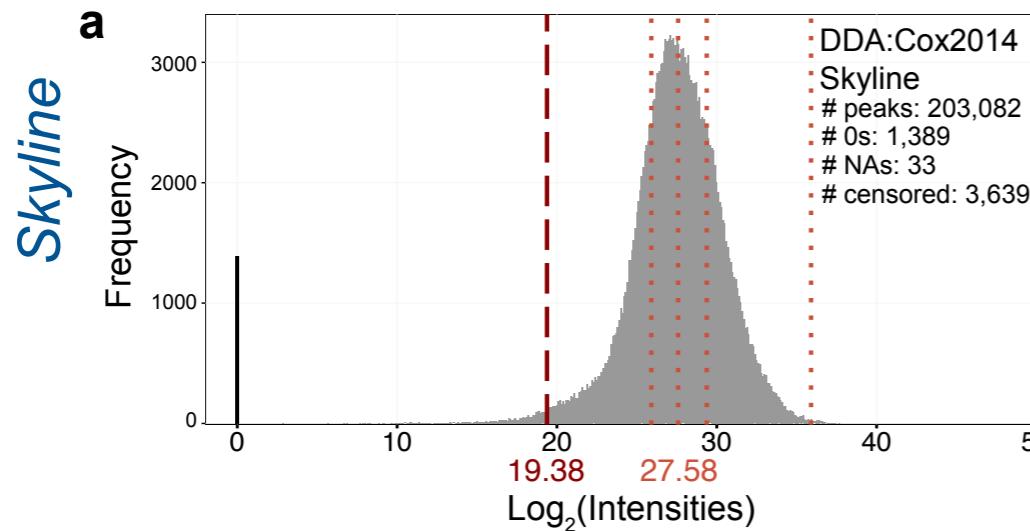
$$\text{Run} \times \text{Feature}_{ijkl} = \epsilon_{ijkl} \stackrel{\text{iid}}{\sim} \mathcal{N}(\mathbf{0}, \sigma_{\epsilon}^2)$$

PROPERTIES OF PEAK INTENSITIES VARY BETWEEN DATA PROCESSING TOOLS

14

DDA: Cox 2014

DIA: Bruderer 2015



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

INTERPRETING CENSORED VALUES

A

	Run ₁	Run ₂	Run ₃	Run ₄	Run ₅	Run ₆	...	Run _{JK-2}	Run _{JK-1}	Run _{JK}	...	Run _{(I-1)JK+1}	Run _{(I-1)JK+2}	Run _{(I-1)JK+3}	Run _{(I-1)JK+4}	Run _{(I-1)JK+5}	Run _{(I-1)JK+6}	...	Run _{IJK-2}	Run _{IJK-1}	Run _{IJK}
Feature ₁	y	y	y	y	y	y	...	y	y	y	...	y	y	y	y	y	y	...	y	y	y
Feature ₂	y	y	y	y	y	NA_{rand}	...	y	y	y	...	y	y	y	y	y	y	...	y	y	y
...
Feature _L	y	NA_{cen}	y	NA_{cen}	NA_{cen}	y	...	NA_{cen}	y	y	...	NA_{cen}	y	y	y	y	y	...	y	NA_{cen}	y

IMPUTING CENSORED VALUES

A

	Run ₁	Run ₂	Run ₃	Run ₄	Run ₅	Run ₆	...	Run _{JK-2}	Run _{JK-1}	Run _{JK}	...	Run _{(I-1)JK+1}	Run _{(I-1)JK+2}	Run _{(I-1)JK+3}	Run _{(I-1)JK+4}	Run _{(I-1)JK+5}	Run _{(I-1)JK+6}	...	Run _{IJK-2}	Run _{IJK-1}	Run _{IJK}	
Feature ₁	y	y	y	y	y	y	...	y	y	y	...	y	y	y	y	y	y	y	...	y	y	y
Feature ₂	y	y	y	y	y	NA_{rand}	...	y	y	y	...	y	y	y	y	y	y	...	y	y	y	
...
Feature _L	y	NA_{cen}	y	NA_{cen}	NA_{cen}	y	...	NA_{cen}	y	y	...	NA_{cen}	y	y	y	y	y	...	y	NA_{cen}	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 ₁	Run ₂	Run ₃	Run ₄	Run ₅	Run ₆	...	Run _{JK-2}	Run _{JK-1}	Run _{JK}	...	Run _{(I-1)JK+1}	Run _{(I-1)JK+2}	Run _{(I-1)JK+3}	Run _{(I-1)JK+4}	Run _{(I-1)JK+5}	Run _{(I-1)JK+6}	...	Run _{IJK-2}	Run _{IJK-1}	Run _{IJK}
Feature ₁	y	y	y	y	y	y	...	y	y	y	...	y	y	y	y	y	y	...	y	y	y
Feature ₂	y	y	y	y	y		...	y	y	y	...	y	y	y	y	y	y	...	y	y	y
...
Feature _L	y	y_{imp}	y	y_{imp}	y_{imp}	y	...	y_{imp}	y	y	...	y_{imp}	y	y	y	y	y	...	y	y_{imp}	y

ROBUST RUN SUMMARIZATION

A

	Run ₁	Run ₂	Run ₃	Run ₄	Run ₅	Run ₆	...	Run _{JK-2}	Run _{JK-1}	Run _{JK}	...	Run _{(I-1)JK+1}	Run _{(I-1)JK+2}	Run _{(I-1)JK+3}	Run _{(I-1)JK+4}	Run _{(I-1)JK+5}	Run _{(I-1)JK+6}	...	Run _{IJK-2}	Run _{IJK-1}	Run _{IJK}	
Feature ₁	y	y	y	y	y	y	...	y	y	y	...	y	y	y	y	y	y	y	...	y	y	y
Feature ₂	y	y	y	y	y	NA_{rand}	...	y	y	y	...	y	y	y	y	y	y	y	...	y	y	y
...
Feature _L	y	NA_{cen}	y	NA_{cen}	NA_{cen}	y	...	NA_{cen}	y	y	...	NA_{cen}	y	y	y	y	y	y	...	y	NA_{cen}	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 ₁	Run ₂	Run ₃	Run ₄	Run ₅	Run ₆	...	Run _{JK-2}	Run _{JK-1}	Run _{JK}	...	Run _{(I-1)JK+1}	Run _{(I-1)JK+2}	Run _{(I-1)JK+3}	Run _{(I-1)JK+4}	Run _{(I-1)JK+5}	Run _{(I-1)JK+6}	...	Run _{IJK-2}	Run _{IJK-1}	Run _{IJK}	
Feature ₁	y	y	y	y	y	y	...	y	y	y	...	y	y	y	y	y	y	y	...	y	y	y
Feature ₂	y	y	y	y	y		...	y	y	y	...	y	y	y	y	y	y	y	...	y	y	y
...
Feature _L	y	y_{imp}	y	y_{imp}	y_{imp}	y	...	y_{imp}	y	y	...	y_{imp}	y	y	y	y	y	y	...	y	y_{imp}	y



TMP : Parameter estimation by robust method

$$y_{ijkl} = \mu + Run_{ijk} + Feature_l + \epsilon_{ijkl}, \text{ where}$$

$$\text{median}_{ijk}(Run_{ijk}) = 0, \text{ median}_l(Feature_l) = 0, \text{ and } \text{median}_{ijk}(\epsilon_{ijkl}) = \text{median}_l(\epsilon_{ijkl}) = 0$$

	Run ₁	Run ₂	Run ₃	Run ₄	Run ₅	Run ₆	...	Run _{JK-2}	Run _{JK-1}	Run _{JK}	...	Run _{(I-1)JK+1}	Run _{(I-1)JK+2}	Run _{(I-1)JK+3}	Run _{(I-1)JK+4}	Run _{(I-1)JK+5}	Run _{(I-1)JK+6}	...	Run _{IJK-2}	Run _{IJK-1}	Run _{IJK}
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}

SUB-PLOT

Summarization over all features in a run

		<i>Run</i>															
		1	2	...	12												
<i>log(Feature int)</i>	1	x_{11}	x_{12}	...	$x_{1\ 12}$												
	2	x_{21}	x_{22}	...	$x_{2\ 12}$												
	...																
n	x_{n1}	x_{n2}	...	$x_{n\ 12}$													

Tukey median polish
Represent features and runs in a sub-plot as 2-way Analysis of Variance

$$x_{ij} = \text{feature}_i + \text{run}_j + \text{error}_{ij}$$

Addition: censored data
Impute missing values by assuming that they have intensities below detection threshold

- Robust parameter estimation

- ◆ subtract column median from each value
- ◆ subtract row median from each value
- ◆ continue until no change
- ◆ obtain fitted values
 - subtract the resulting residuals from the original values
- ◆ obtain array-based summary
 - average fitted values over the column

LINEAR MODELS FOR THE DESIGN

	Run ₁	Run ₂	Run ₃	Run ₄	Run ₅	Run ₆	...	Run _{JK-2}	Run _{JK-1}	Run _{JK}	...	Run _{(I-1)JK+1}	Run _{(I-1)JK+2}	Run _{(I-1)JK+3}	Run _{(I-1)JK+4}	Run _{(I-1)JK+5}	Run _{(I-1)JK+6}	...	Run _{IJK-2}	Run _{IJK-1}	Run _{IJK}
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}



C Step 2 : Model-based inference by whole plot

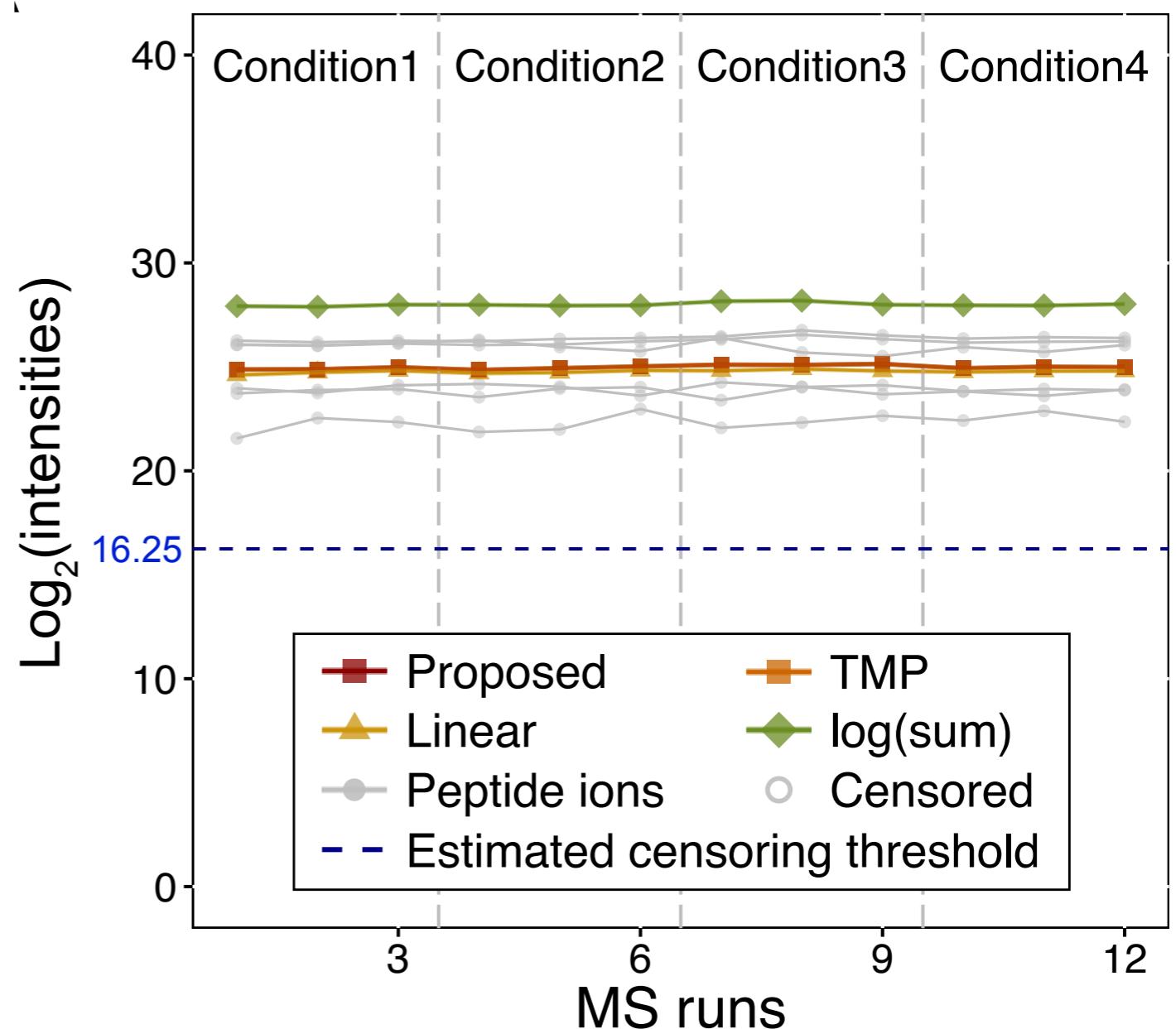
$$\hat{y}_{ijk} = \mu + Condition_i + Subject(Condition)_{j(i)} + \psi_{ijk}, \text{ where}$$

$$\sum_i Condition_i = 0, \quad Subject(Condition)_{j(i)} \stackrel{iid}{\sim} \mathcal{N}(0, \sigma_{Subject}^2), \quad \psi_{ijk} \stackrel{iid}{\sim} \mathcal{N}(0, \sigma_\psi^2)$$

Condition ₁										...	Condition _I										
Subject ₁			Subject ₂			...	Subject _J			...	Subject _{(I-1)J+1}			Subject _{(I-1)+2}			...	Subject _{IJ}			
Run ₁	Run ₂	Run ₃	Run ₄	Run ₅	Run ₆	...	Run _{JK-2}	Run _{JK-1}	Run _{JK}	...	Run _{(I-1)JK+1}	Run _{(I-1)JK+2}	Run _{(I-1)JK+3}	Run _{(I-1)JK+4}	Run _{(I-1)JK+5}	Run _{(I-1)JK+6}	...	Run _{IJK-2}	Run _{IJK-1}	Run _{IJK}	
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

Methods perform similarly with high quality data



Condition2-Condition1 : True fold change=1

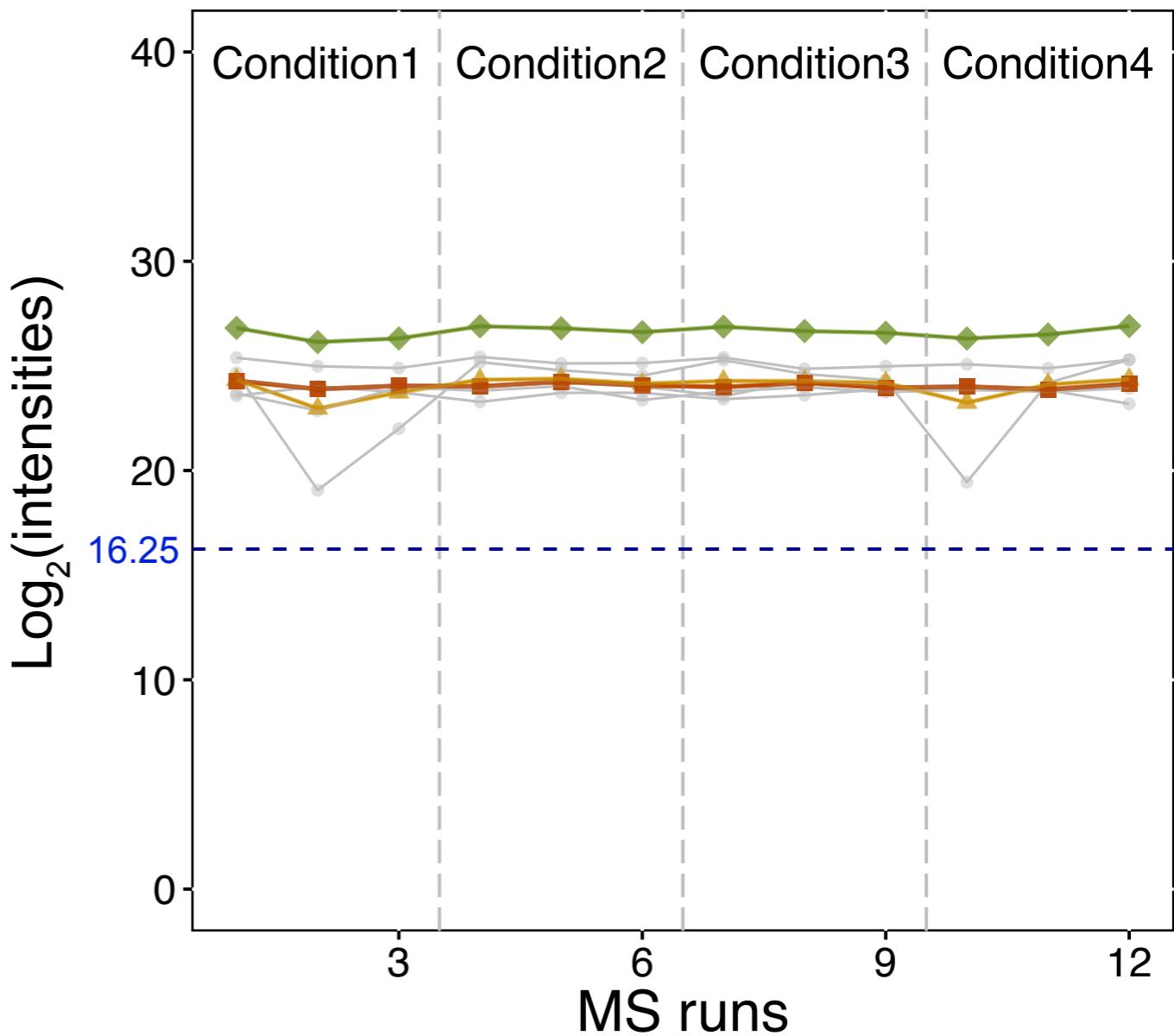
EstimatedFC Adj.pvalue

Proposed	1.016	0.999
TMP	1.016	0.999
Linear model	1.020	0.999
log(sum)	1.019	0.999



ROBUSTNESS TO OUTLIERS

*Outliers in low intensities:
robust summarization with
TMP improves upon linear
model*



Condition3-Condition1 : True fold change=1

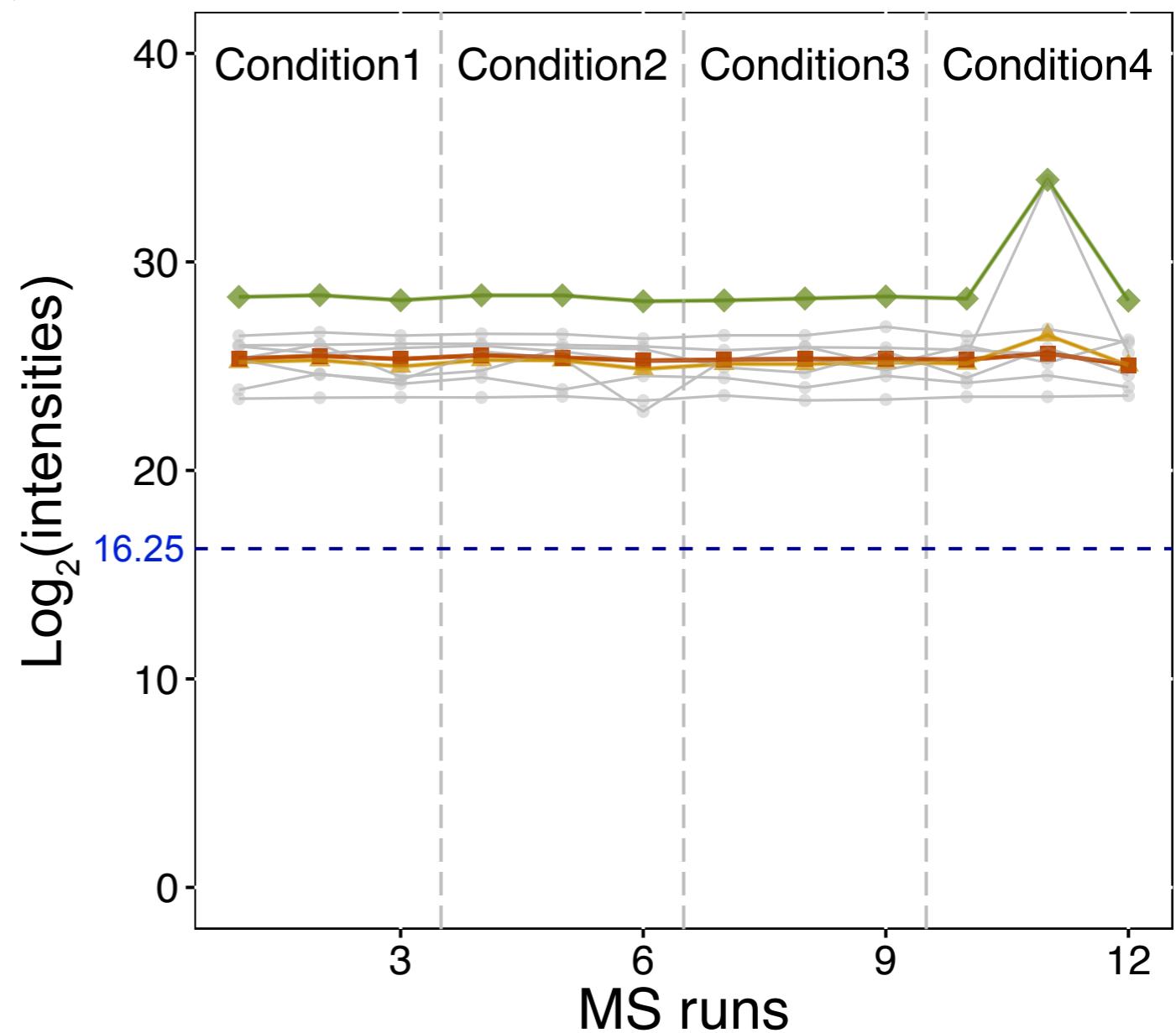
EstimatedFC Adj.pvalue

Peptide ions	
Proposed	
TMP	

	EstimatedFC	Adj.pvalue
Proposed	0.979	0.952
TMP	0.979	0.956
Linear model	1.488	0.815
log(sum)	1.218	0.734

ROBUSTNESS TO OUTLIERS

*Outliers in high intensities:
robust summarization with
TMP improves upon log(sum)*



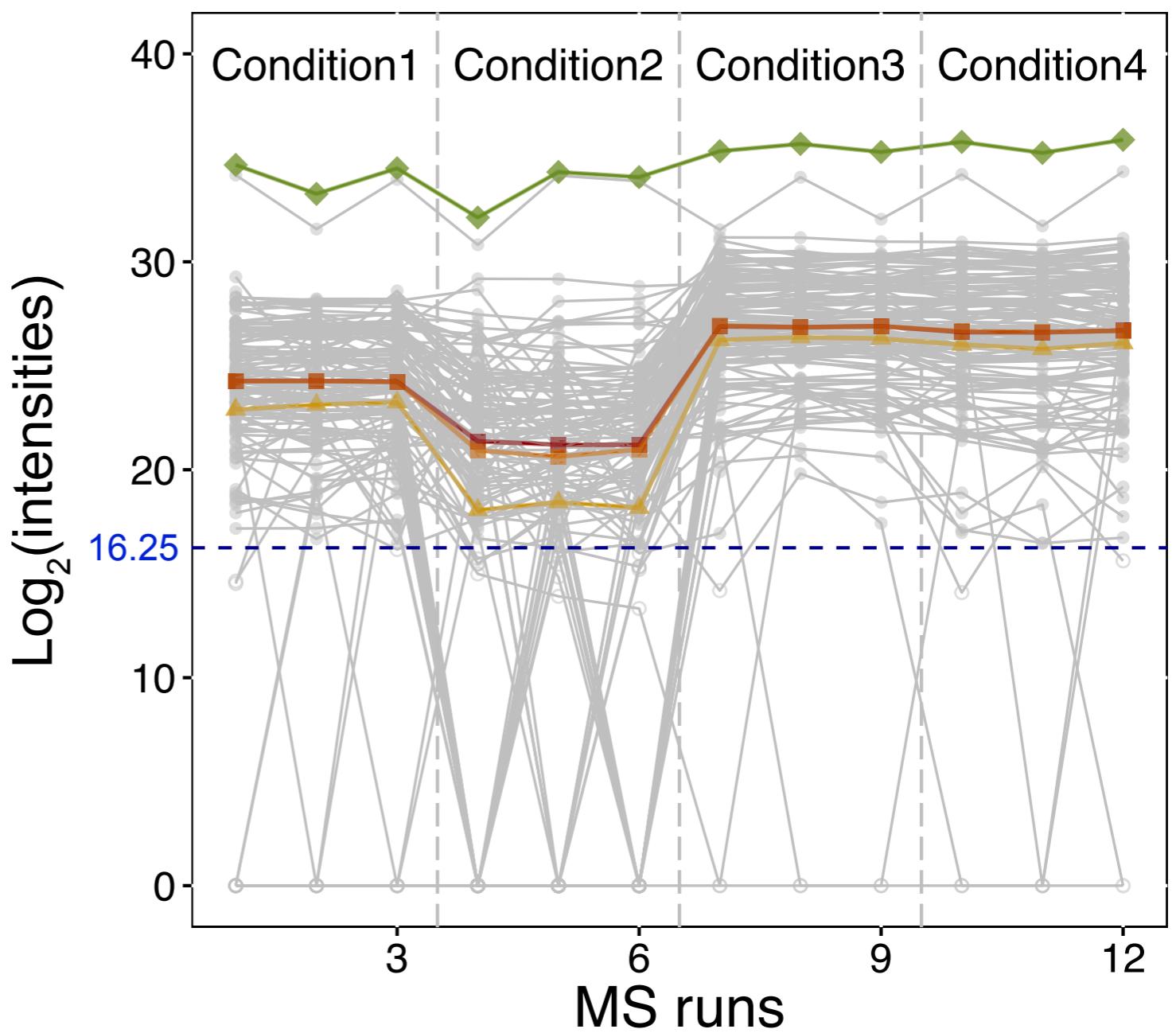
Condition4-Condition1 : True fold change=1
EstimatedFC Adj.pvalue

Peptide ions	
Proposed	
TMP	

	EstimatedFC	Adj.pvalue
Proposed	0.951	0.948
TMP	0.951	0.948
Linear model	1.317	0.881
log(sum)	3.514	0.741

ROBUSTNESS TO OUTLIERS

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

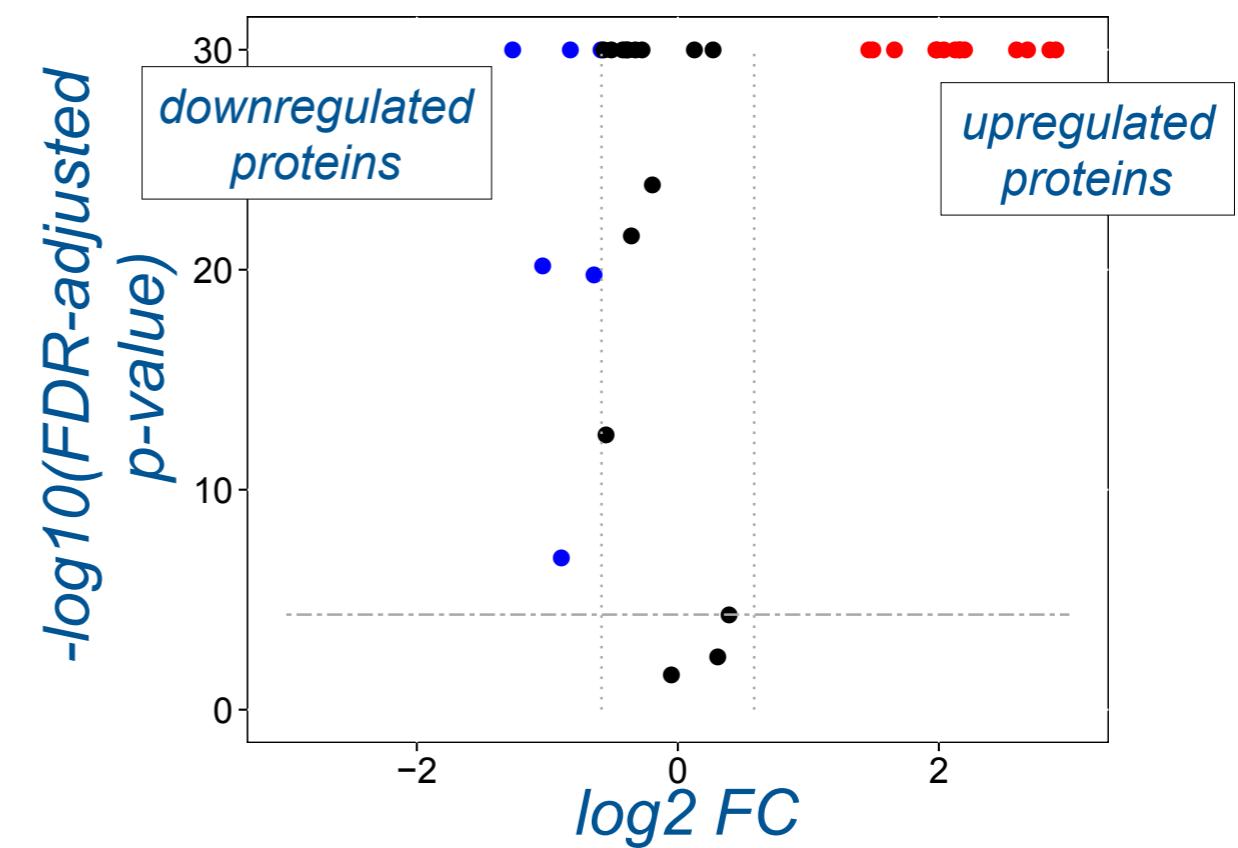
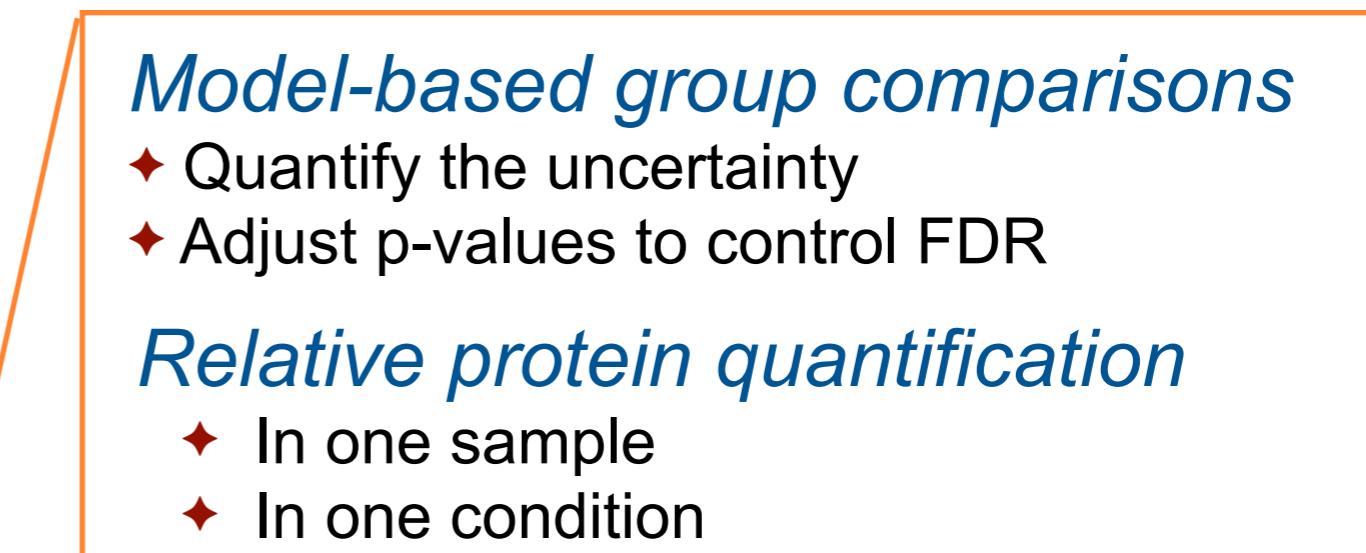
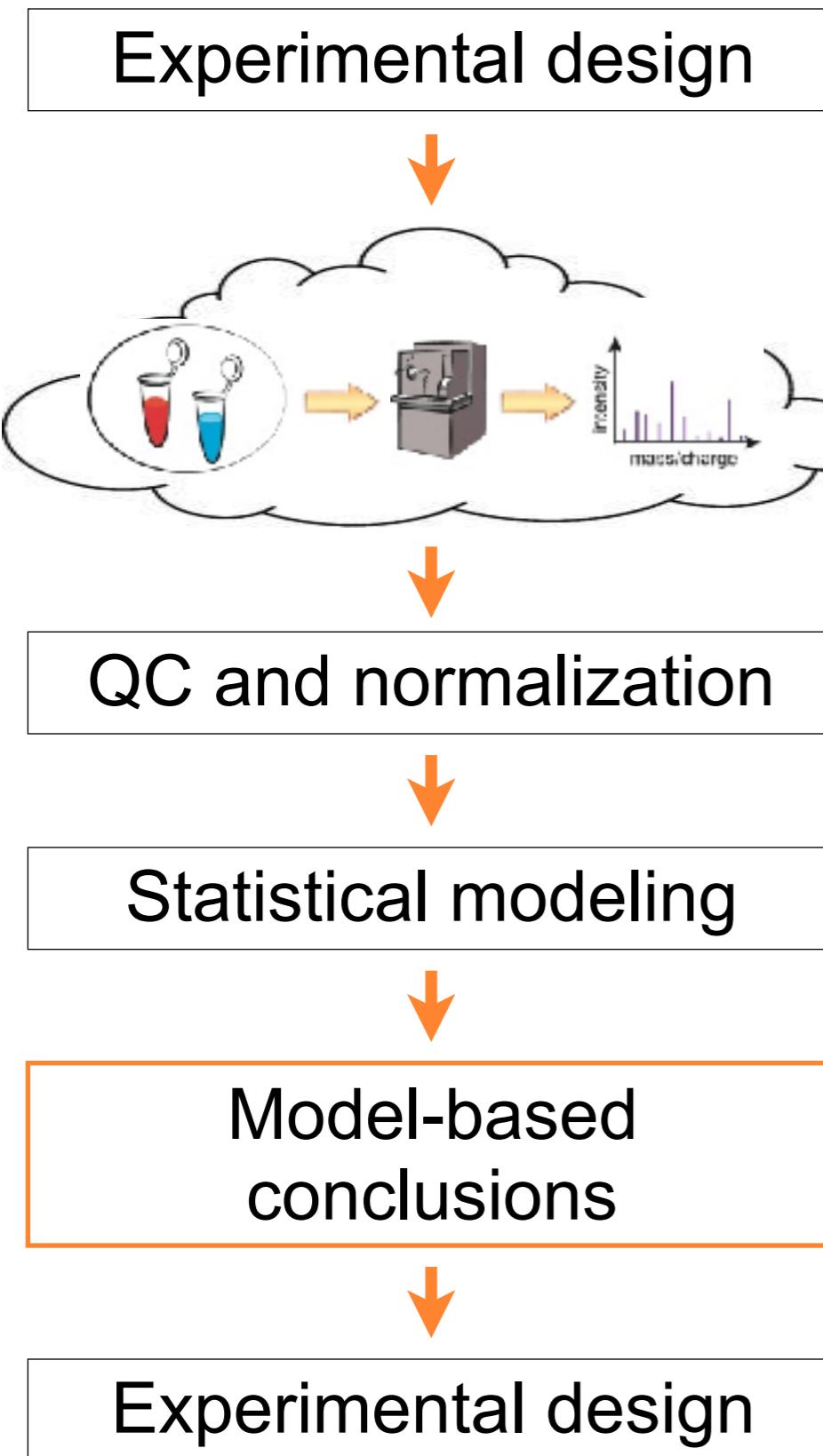


Condition1-Condition2 : True fold change=7.5
EstimatedFC Adj.pvalue

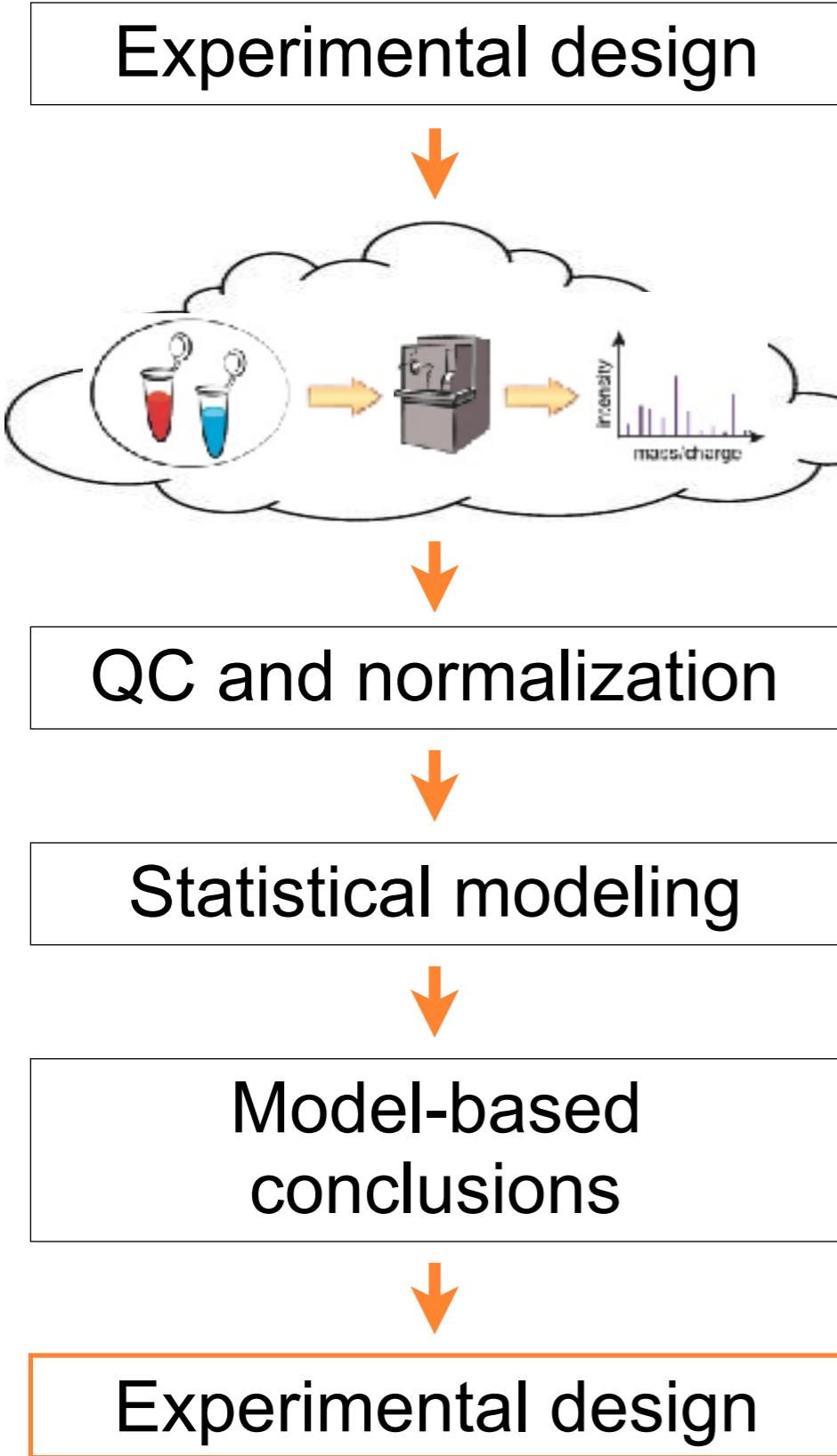
Peptide ions	
Proposed	
TMP	

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

A TYPICAL ANALYSIS WORKFLOW

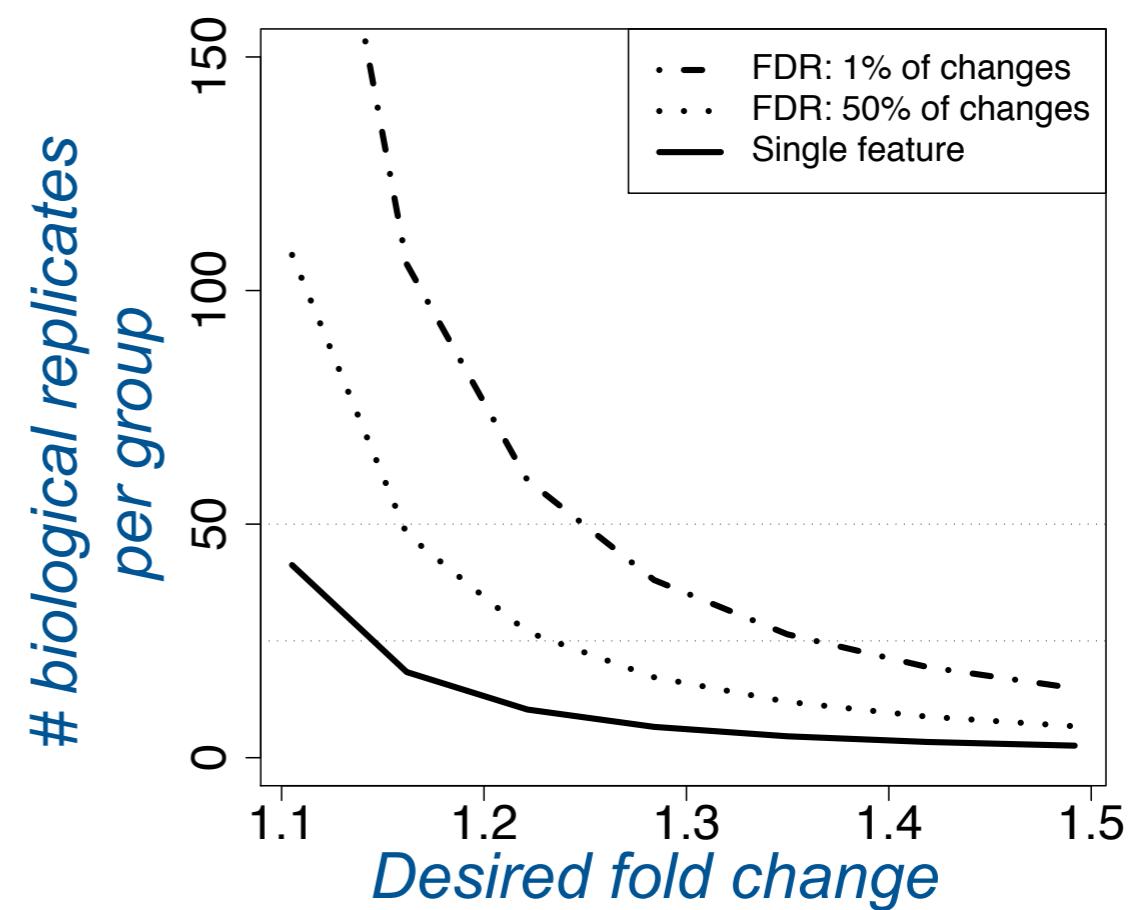


A TYPICAL ANALYSIS WORKFLOW



Use the dataset to improve:

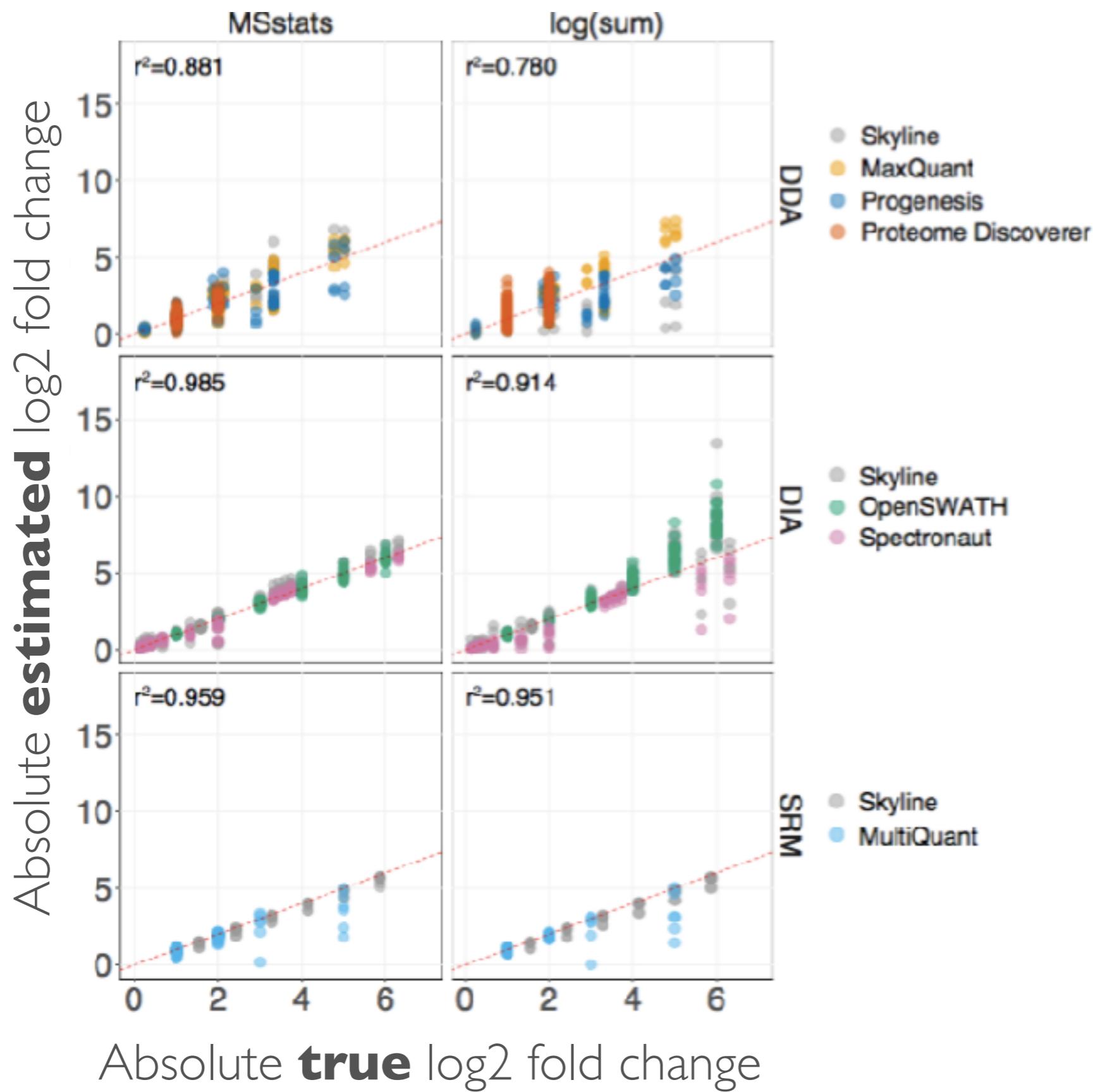
- Subject selection: matching
- Resource allocation: blocking
- Calculation of sample size



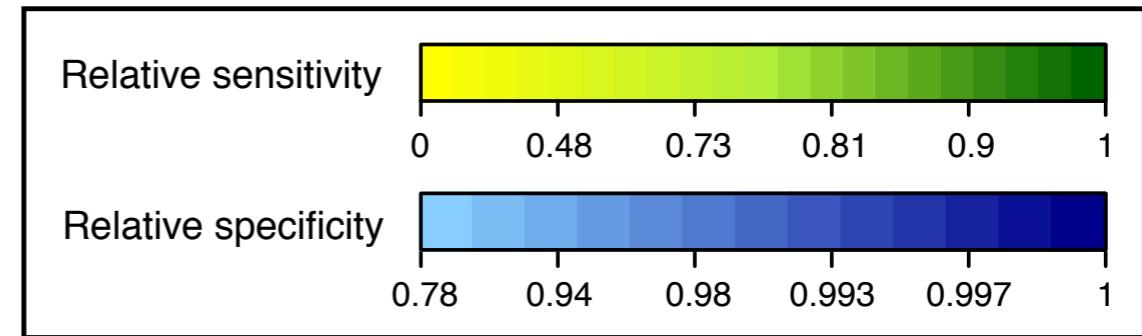
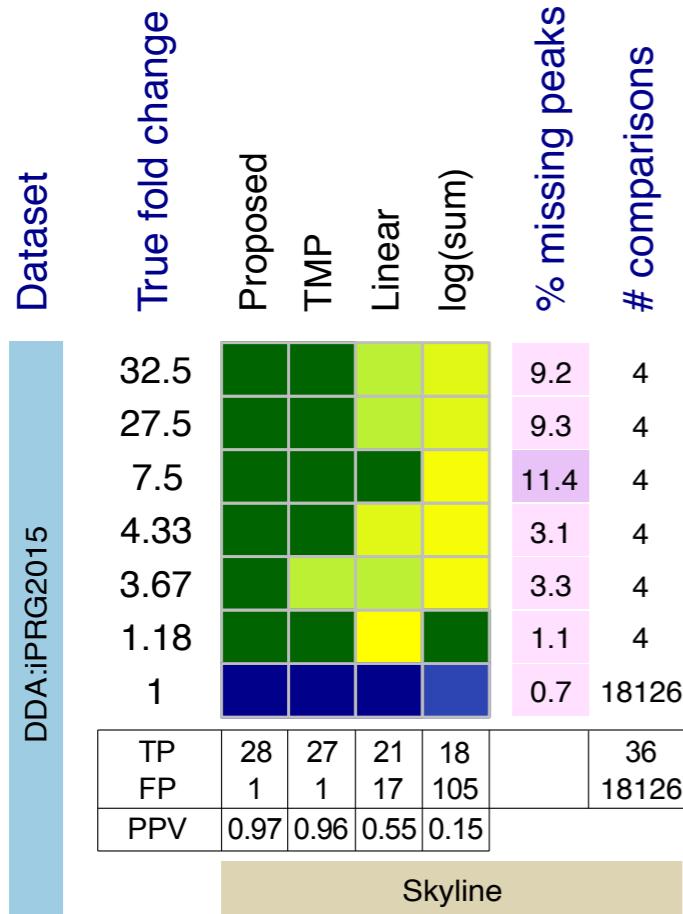
OUTLINE

- Methods: background
 - Linear mixed effects models
- Evaluation
 - Spike-in and experimental datasets

ESTIMATION OF LOG-FOLD CHANGE

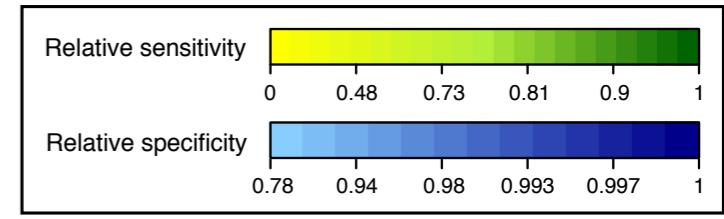
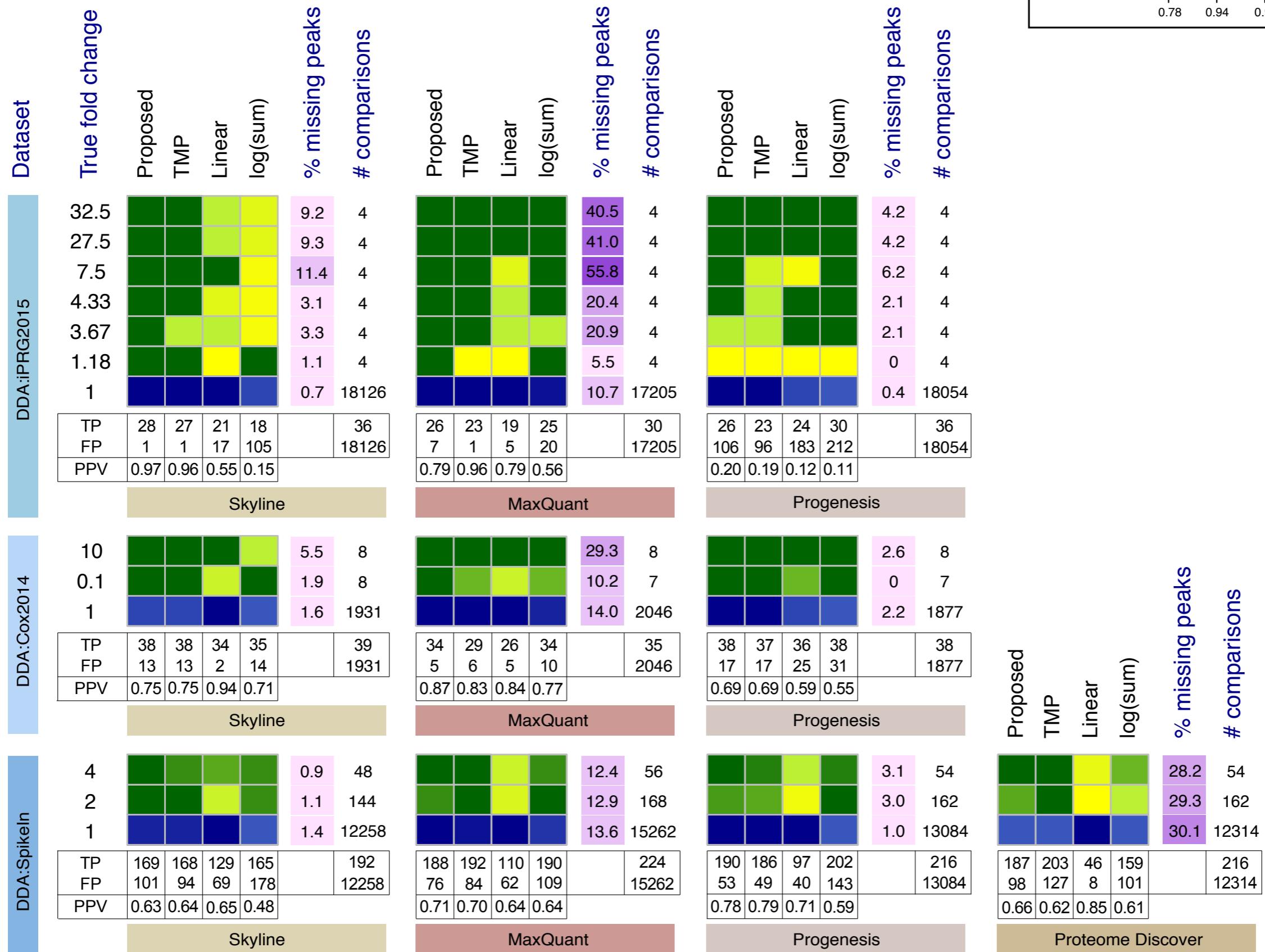


TESTING: DDA

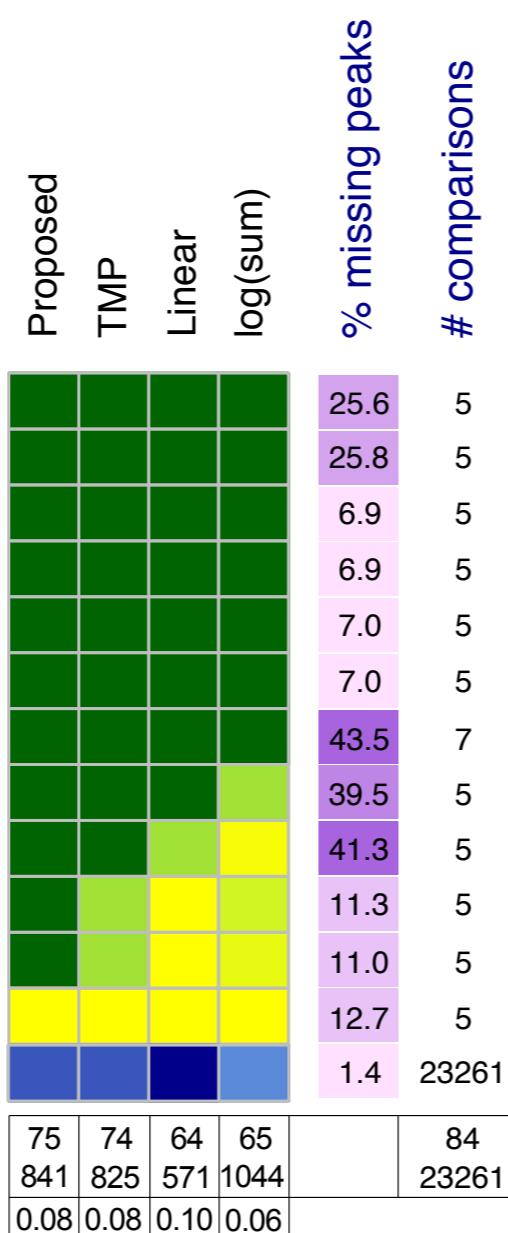
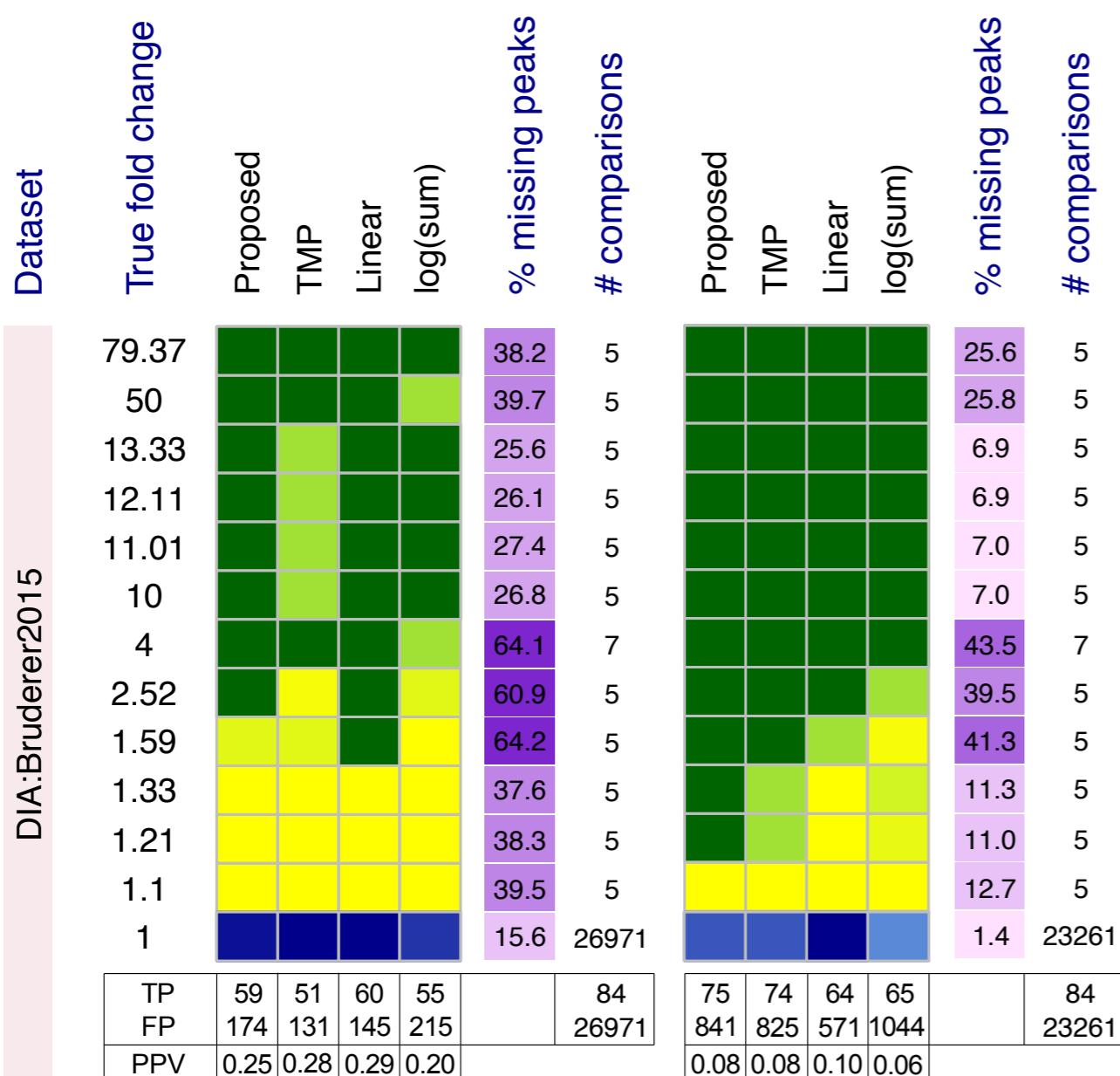


TESTING: DDA

a

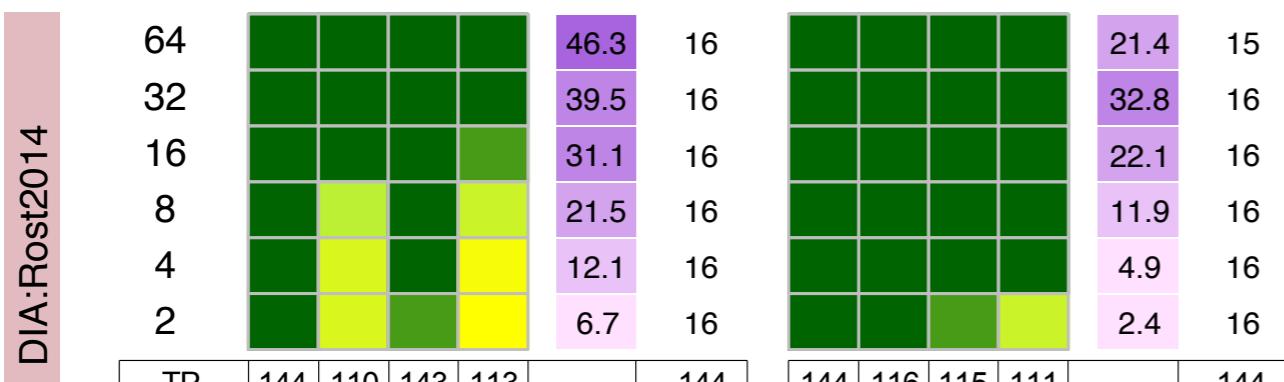


TESTING: DIA



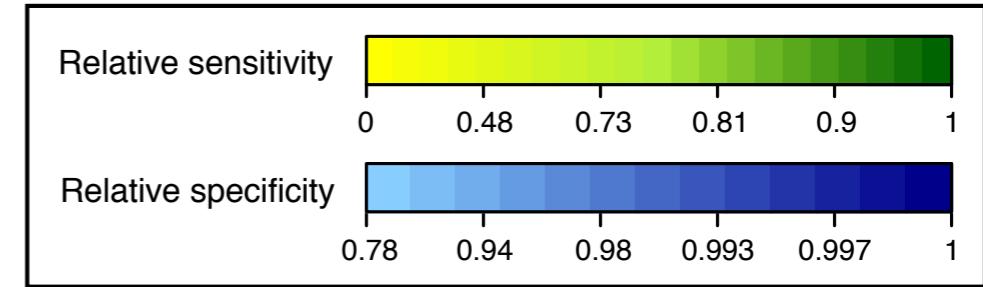
Skyline

Spectronaut



Skyline

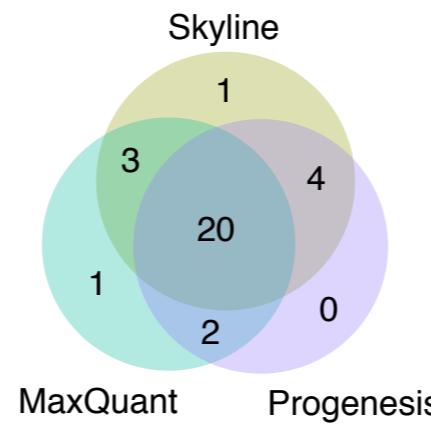
OpenSWATH



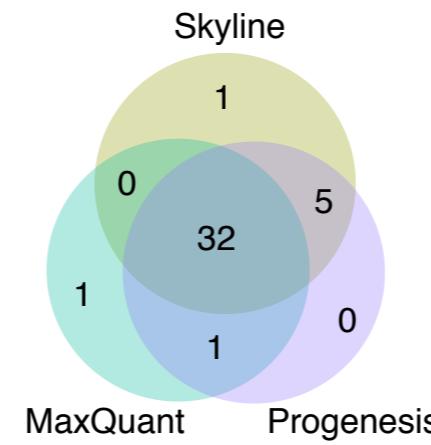
BETTER STATISTICAL METHODS ENHANCE REPRODUCIBLE RESEARCH

Better agreement
in #differentially
abundant proteins
between tools

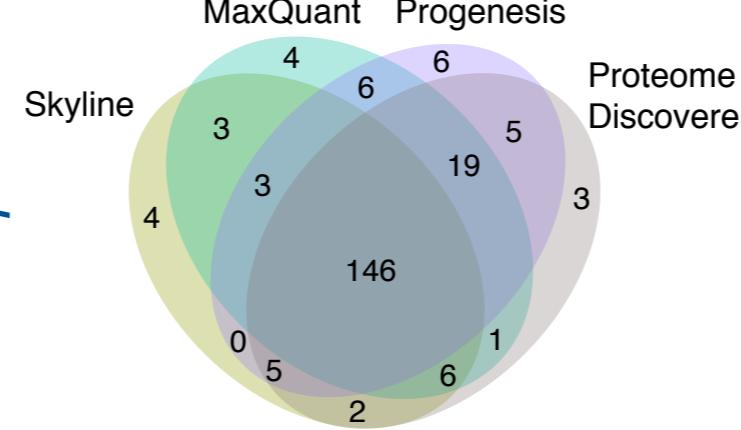
DDA: iPRG2015



DDA: Cox 2014

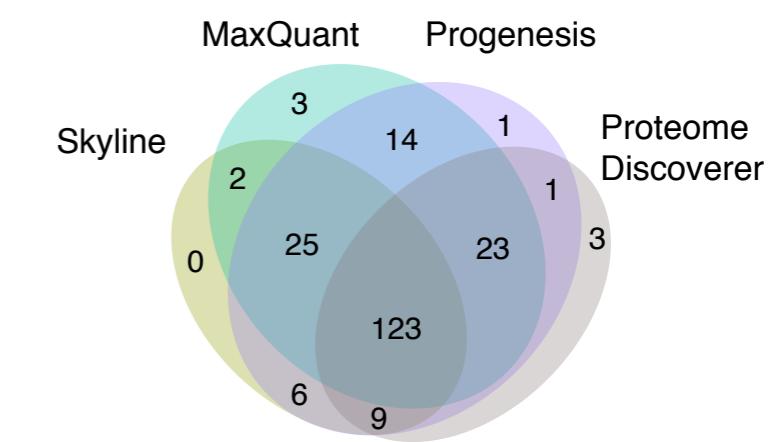
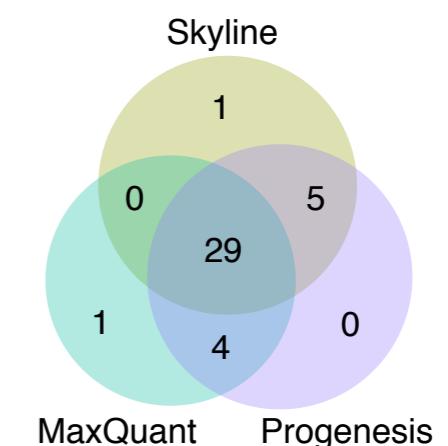
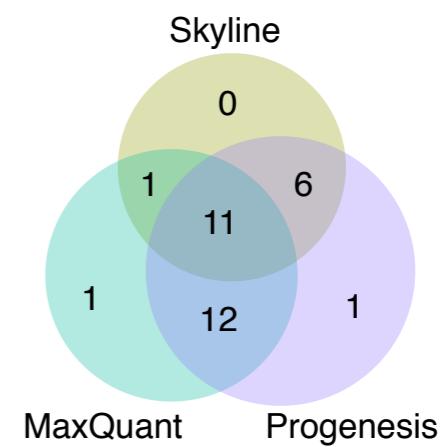


DDA: Spike-in



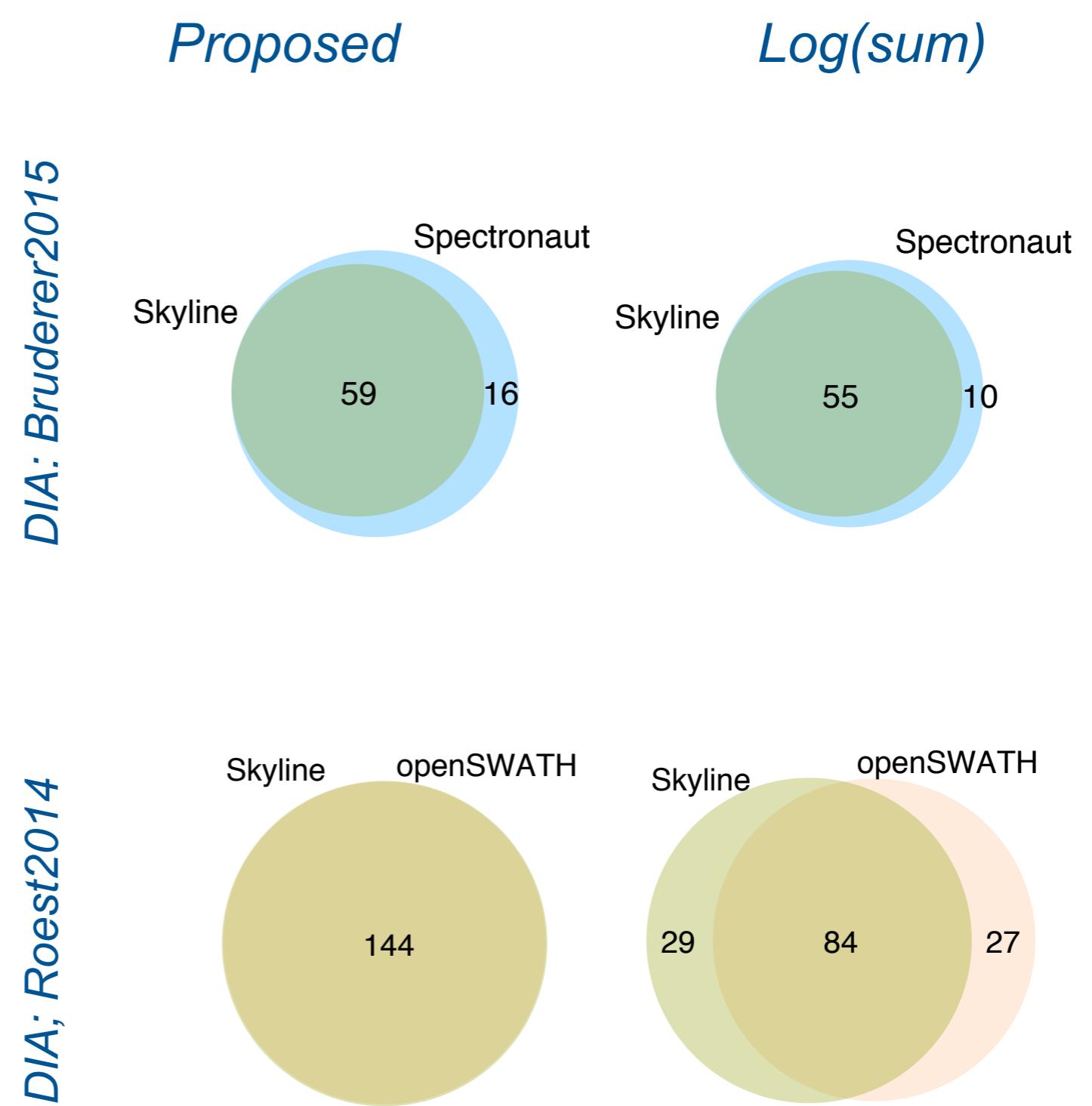
Proposed

Log(sum)

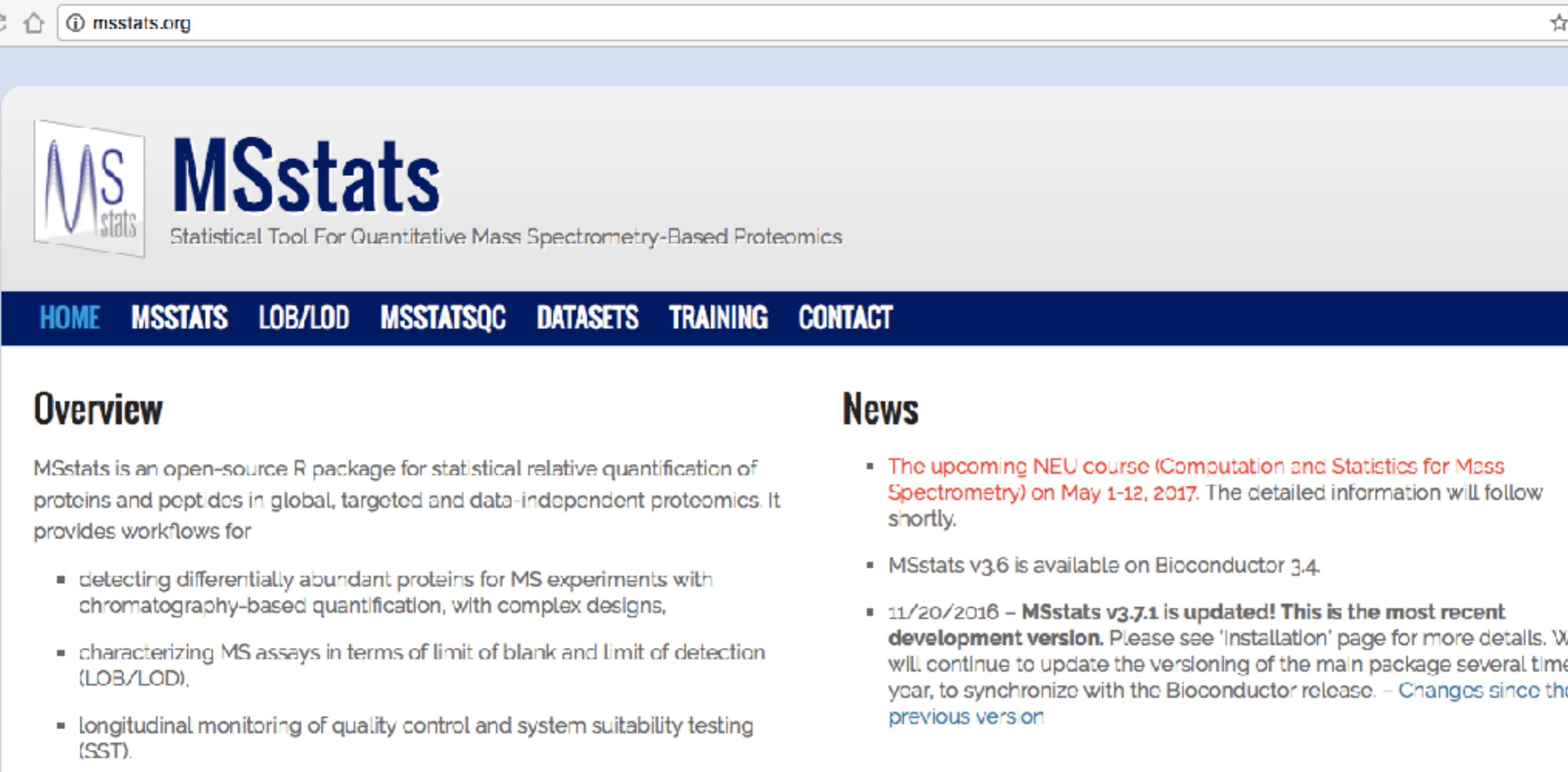


BETTER STATISTICAL METHODS ENHANCE REPRODUCIBLE RESEARCH

Better agreement
in #differentially
abundant proteins
between tools



MSSTATS IS OPEN-SOURCE, R-BASED AND PUBLICLY AVAILABLE



The screenshot shows the MSstats website. At the top, there's a header bar with icons for search, home, and a link to msstats.org. Below the header is the MSstats logo, which consists of a stylized 'M' and 'S' icon followed by the word 'stats'. The main title 'MSstats' is in large blue letters, with the subtitle 'Statistical Tool For Quantitative Mass Spectrometry-Based Proteomics' in smaller text below it. A navigation menu bar at the bottom of the page includes links for HOME, MSSTATS, LOB/LOD, MSSTATSQC, DATASETS, TRAINING, and CONTACT.

Overview

MSstats is an open-source R package for statistical relative quantification of proteins and peptides in global, targeted and data-independent proteomics. It provides workflows for

- detecting differentially abundant proteins for MS experiments with chromatography-based quantification, with complex designs,
- characterizing MS assays in terms of limit of blank and limit of detection (LOB/LOD),
- longitudinal monitoring of quality control and system suitability testing (SST).

News

- The upcoming NEU course (Computation and Statistics for Mass Spectrometry) on May 1-12, 2017. The detailed information will follow shortly.
- MSstats v3.6 is available on Bioconductor 3.4.
- 11/20/2016 – **MSstats v3.7.1 is updated! This is the most recent development version.** Please see 'Installation' page for more details. We will continue to update the versioning of the main package several times per year, to synchronize with the Bioconductor release. – Changes since the previous version

STATISTICAL METHODS FOR ASSAY CHARACTERIZATION

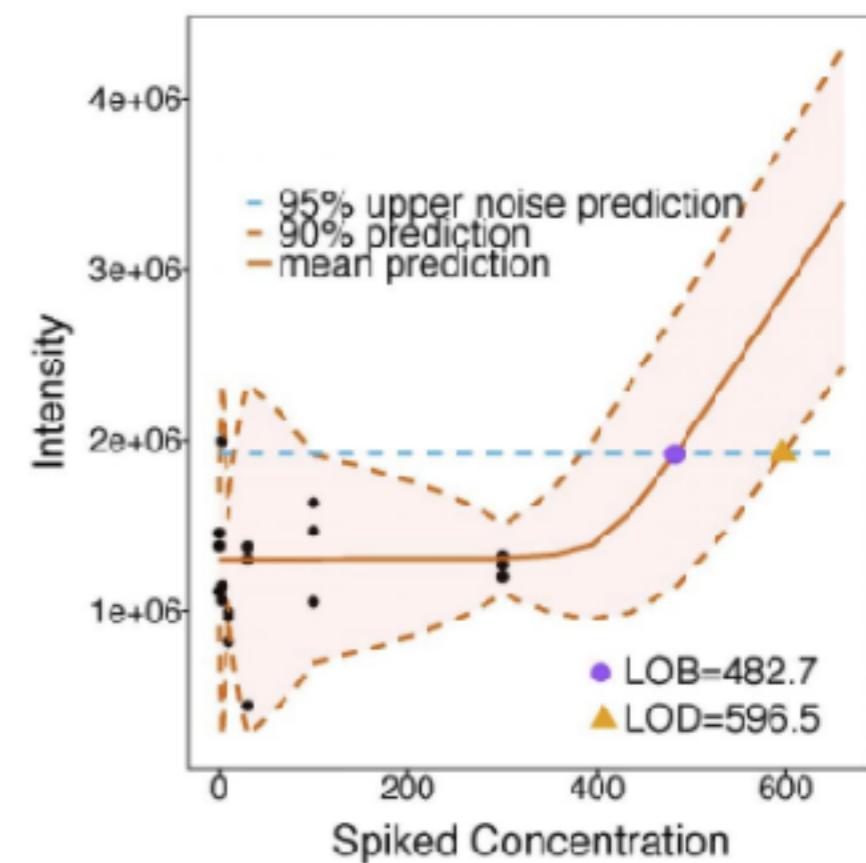


[HOME](#) [MSSTATS](#) [LOB/LOD](#) [MSSTATSQC](#) [DATASETS](#) [TRAINING](#) [CONTACT](#)

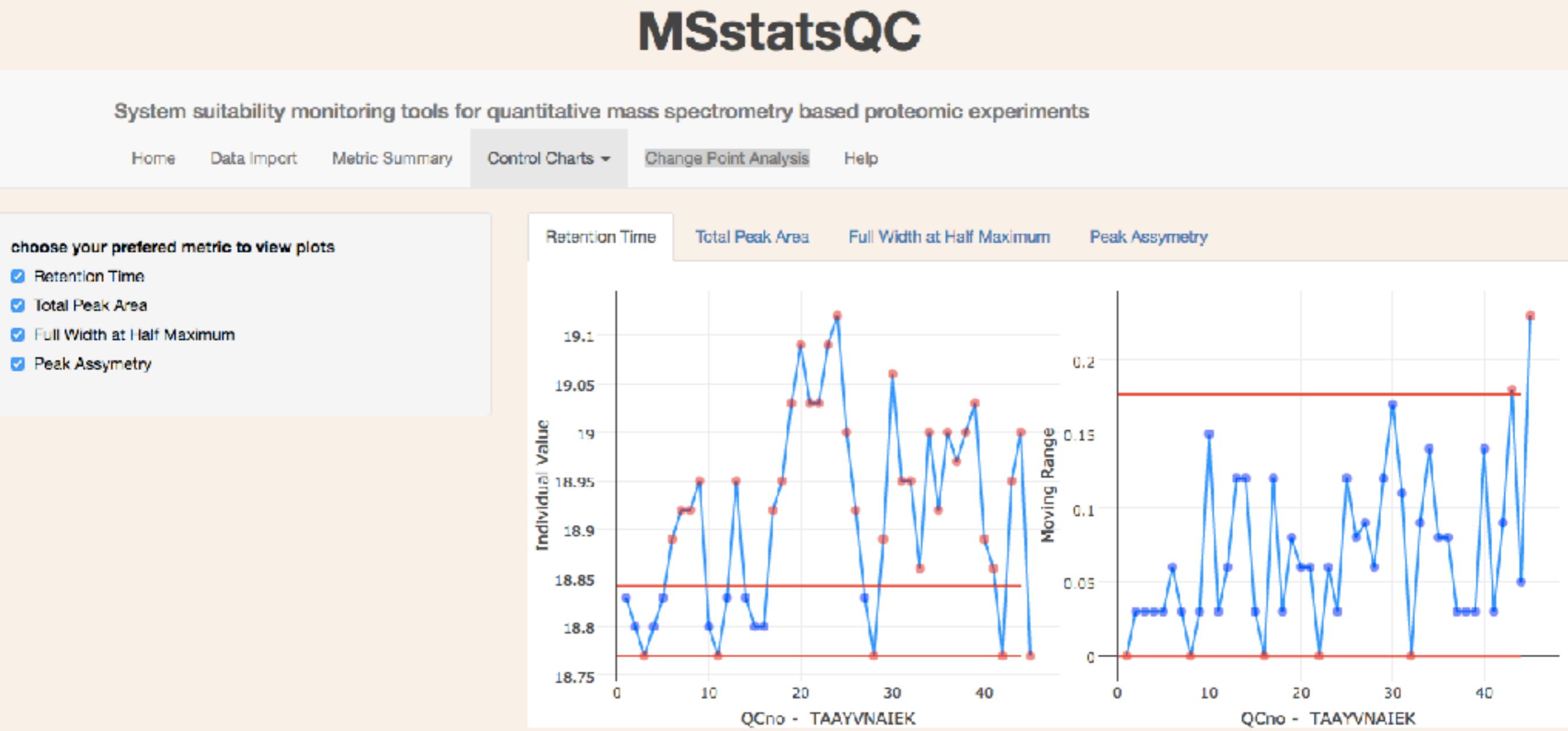
LOB/LOD ANALYSIS

Assay characterization : estimation of limit of blanc (LoB) and limit of detection (LoD)

The need for assay characterization is ubiquitous in quantitative mass spectrometry-based proteomics. Although many assay characteristics exist, the limit of blank (LOB) and limit of detection (LOD) are particularly useful figures of merit. LOB and LOD are determined by repeatedly measuring the peak intensities of peptides in samples with known peptide concentrations, and deriving an Intensity versus concentration response curve. Most commonly, a weighted linear regression is fit to the intensity-concentration response, and LOB and LOD are estimated from the fit. Linear methods, however, inaccurately characterize assays containing a noise threshold at low concentrations, which is a very common situation. We propose a new approach based on non-linear regression that correctly captures the noise threshold. In absence of a noise threshold, the estimates of LOB/LOD obtained with non-linear statistical modeling are identical to those of weighted linear regression. However, in presence of a noise threshold the non-linear model changed the estimates of LOB/LOD by up to 20-40%. It improved the accuracy of the results, and avoided the unduly optimistic estimation of these figures of merit. We implemented the non-linear regression approach in the open-source R-based software MSstats, and advocate its general use for mass spectrometric protein assay characterization.



LONGITUDINAL PROFILING FOR SYSTEM SUITABILITY AND QUALITY CONTROL



ACKNOWLEDGEMENTS

Northeastern University

Kylie Bemis
Meena Choi
Eralp Dogu
Dan Guo
April Harry
Ting Huang
Cyril Galitzine
Robert Ness
Sara Taheri
Tsung-Heng Tsai

ETH Zurich

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

University of Washington

Michael MacCoss
Brendan MacLean
Jarrett Egertson

Biognosis

Lukas Reiter

iPRG 2015

Zeynep Eren-Dogu
Chris Colangelo
John Cottrell
Michael Hopman
Eugene Kapp
Santa Kim
Henry Lam
Tom Neubert
Magnus Palmblad
Brett Phinney
Sue Weintraub,

