

MSSTATS

Meena Choi, Olga Vitek

College of Science

College of Computer and Information Science



Northeastern University

OUTLINE

- iPRG2015: detection of diff. abundance
 - Community effort with label-free shotgun proteomics
- MSstats
 - Normalization, statistical modeling, inference
 - Evaluation
- Extensions to MSstats
 - Assay characterization, longitudinal monitoring

ABRF IPRG STUDY 2015

Detection of differentially abundant proteins in controlled mixture

Name	Origin	Molecular Weight	Samples			
			1	2	3	4
A Ovalbumin	Chicken Egg White	45KD	65	55	15	2
B Myoglobin	Equine Heart	17KD	55	15	2	65
C Phosphorylase b	Rabbit Muscle	97KD	15	2	65	55
D Beta-Galactosidase	Escherichia Coli	116KD	2	65	55	15
E Bovine Serum Albumin	Bovine Serum	66KD	11	0.6	10	500
F Carbonic Anhydrase	Bovine Erythrocytes	29KD	10	500	11	0.6

Spiked into a constant background: tryptic digests of S. cerevisiae

Choi et al., “ABRF Proteome Informatics Research Group (iPRG) 2015 Study: Detection of differentially abundant proteins in label-free quantitative LC-MS/MS experiments”, *Journal of Proteome Research*, 2017.

EXPERIMENTAL PROCEDURES

- **Background**

- ◆ 200ng of tryptic digests of *S. cerevisiae*

- **Spectral acquisition**

- ◆ *Three technical replicates per sample*
- ◆ *Randomized order*
- ◆ Separation:
 - ◆ Thermo nLC 1000 system
 - ◆ 110-min linear gradient
- ◆ Spectral acquisition:
 - ◆ *DDA profile mode in Orbitrap*
 - ◆ Resolution 70,000 for MS and 17,500 for MS/MS
 - ◆ MSI scan range 300-1650 m/z
 - ◆ Normalized collision energy 27%
 - ◆ Singly charged ions excluded

SPECTRAL PROCESSING

- **MS/MS identification**

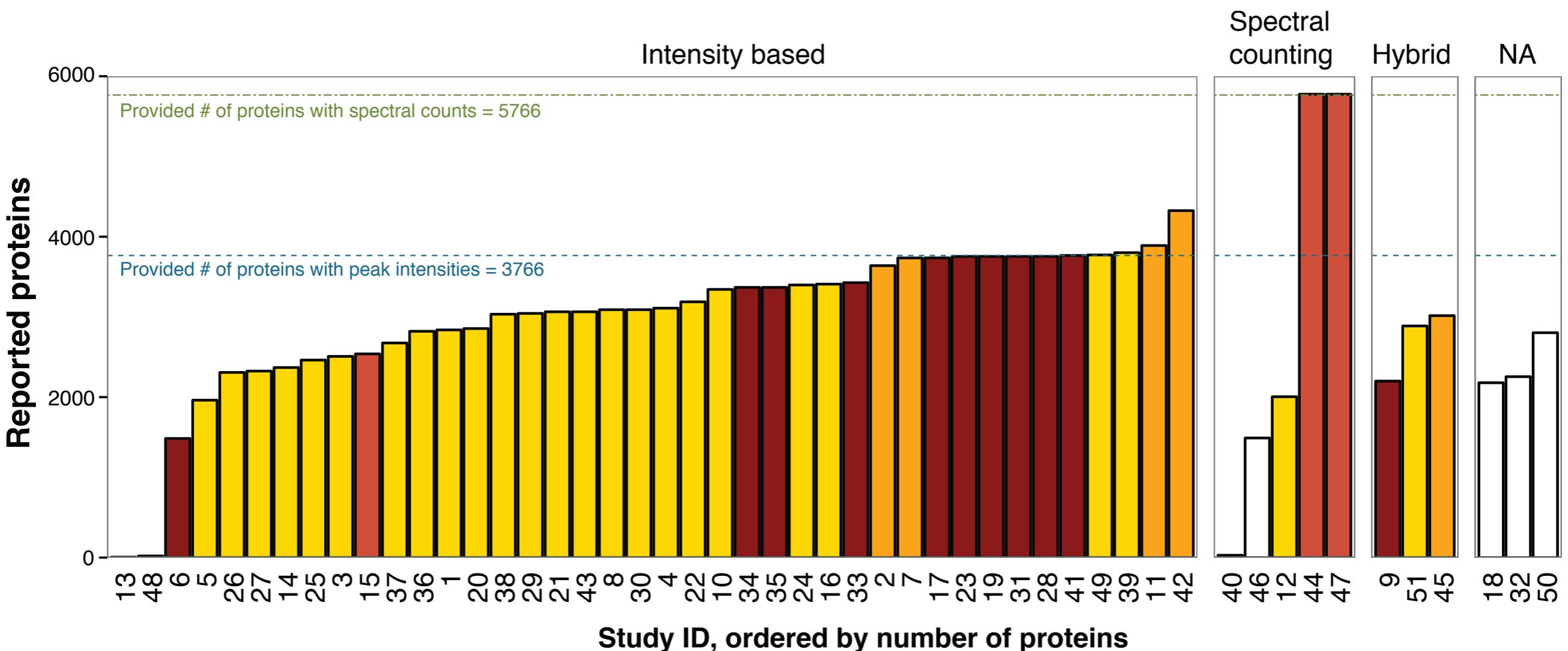
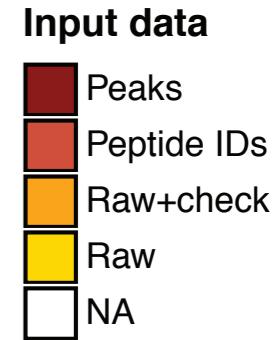
- ◆ Database search: OMSSA7, MS-GF+, Comet
- ◆ Q-value: target-decoy.
 - ◆ No filtering!
- ◆ 5,766 proteins, ~26,242 PSMs/run
- ◆ 48% of proteins had 1 or 2 PSMs
- ◆ 29% of proteins had 15 or more PSMs

- **MSI quantification with Skyline**

- ◆ Original processing
 - ◆ 3,766 proteins, 29,575 features
 - ◆ median 5 features/protein
- ◆ Post-study processing
 - ◆ 3,027 proteins, 34,783 features
 - ◆ median 7 features/protein

DIVERSE SUBMISSIONS

*INPUT, PROTEIN NUMBER,
AND CHOICE OF QUANTIFICATION*

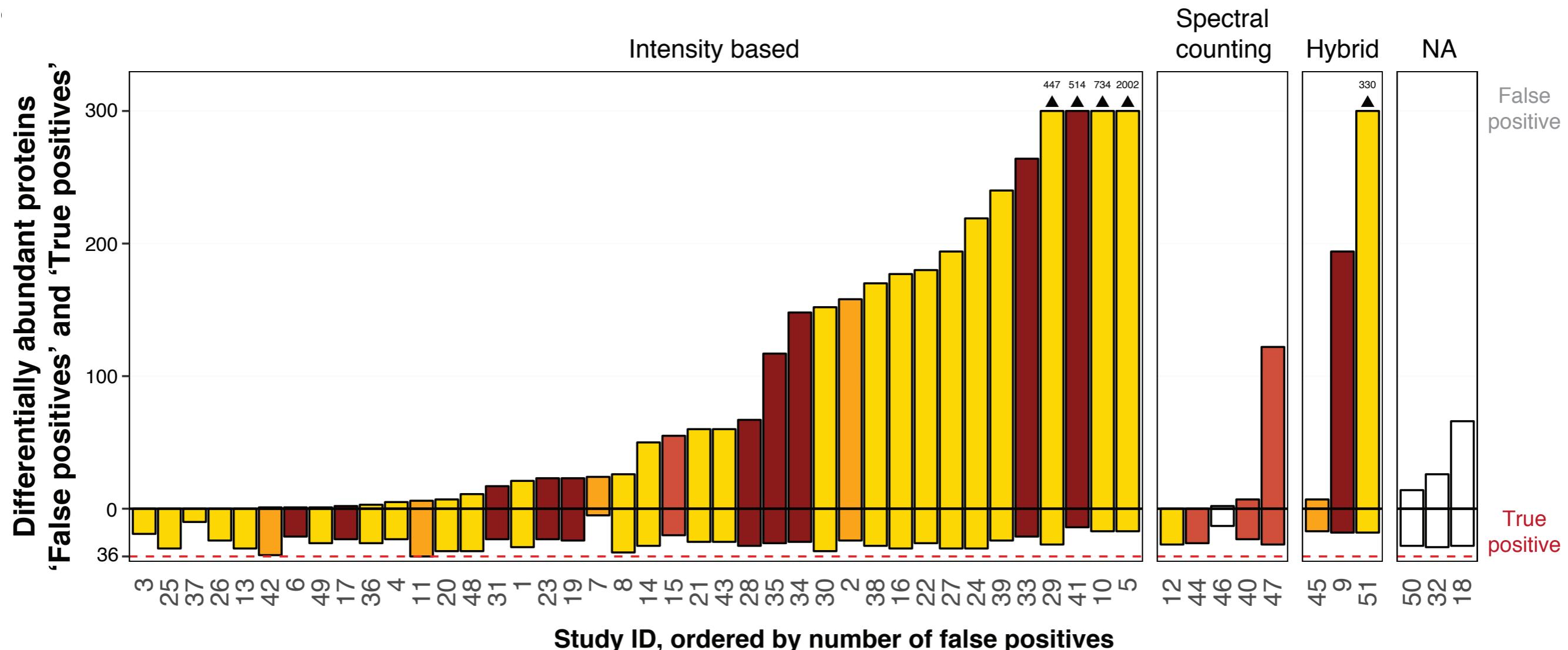


DIVERSE SUBMISSIONS

ACCURACY OF DETECTING DIFFERENTIAL ABUNDANCE

Input data

- Peaks
- Peptide IDs
- Raw+check
- Raw
- NA

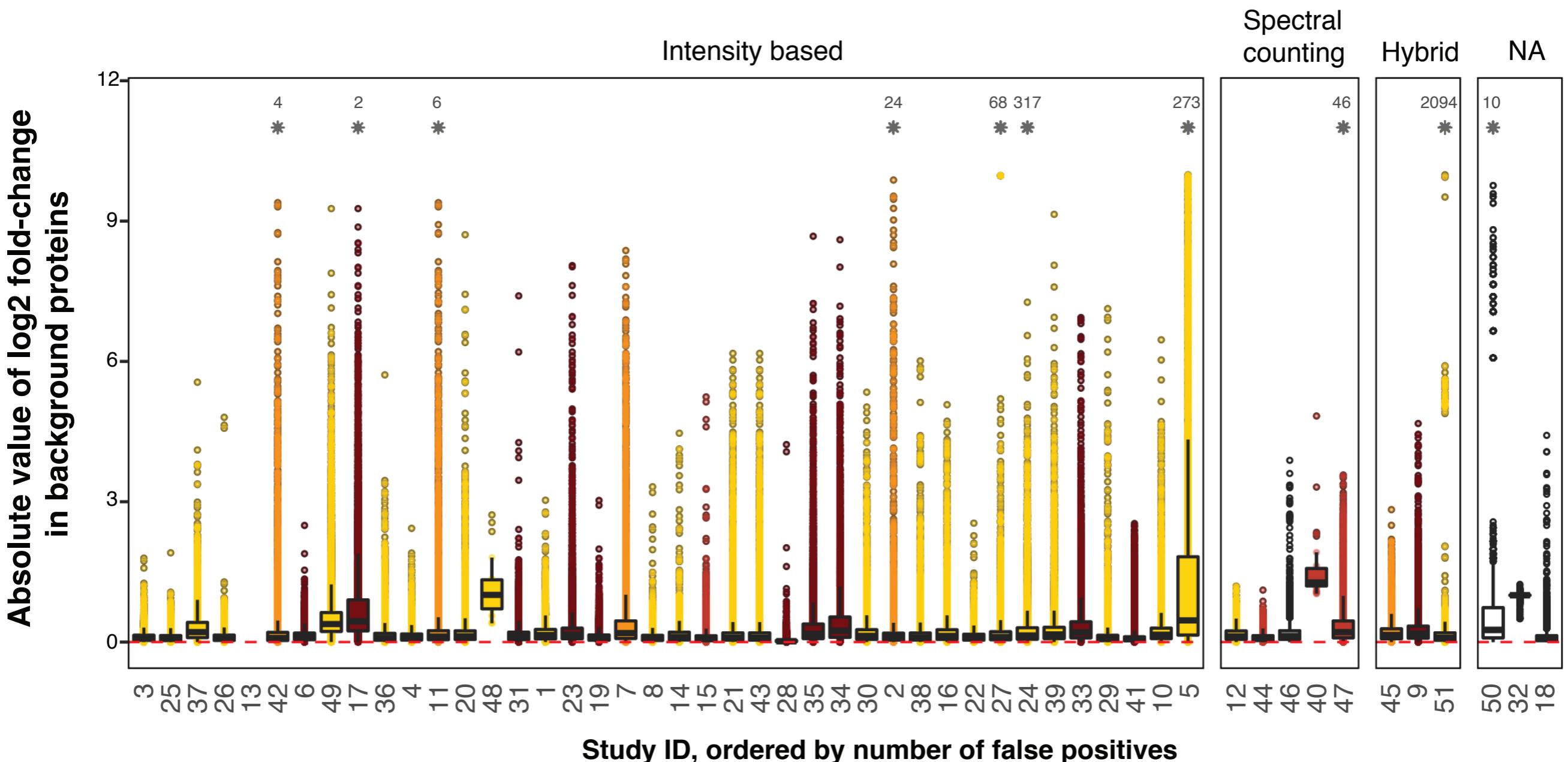


DIVERSE SUBMISSIONS

ACCURACY OF ESTIMATING FOLD CHANGE

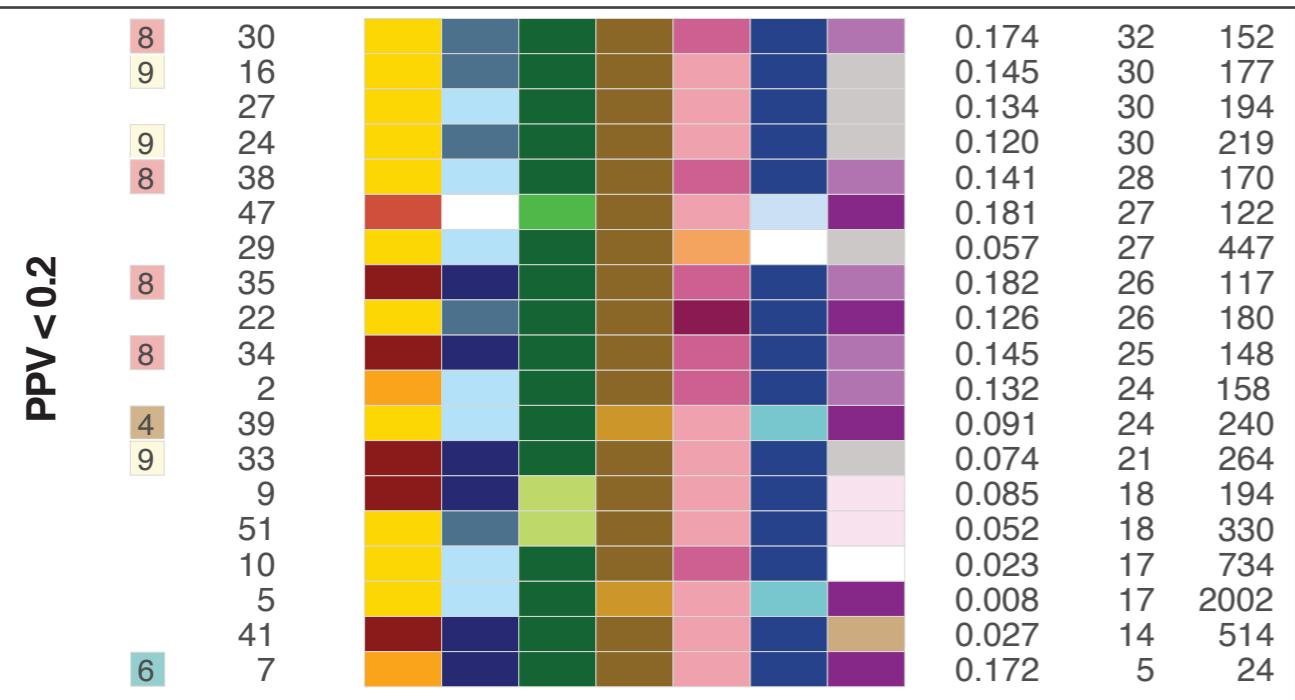
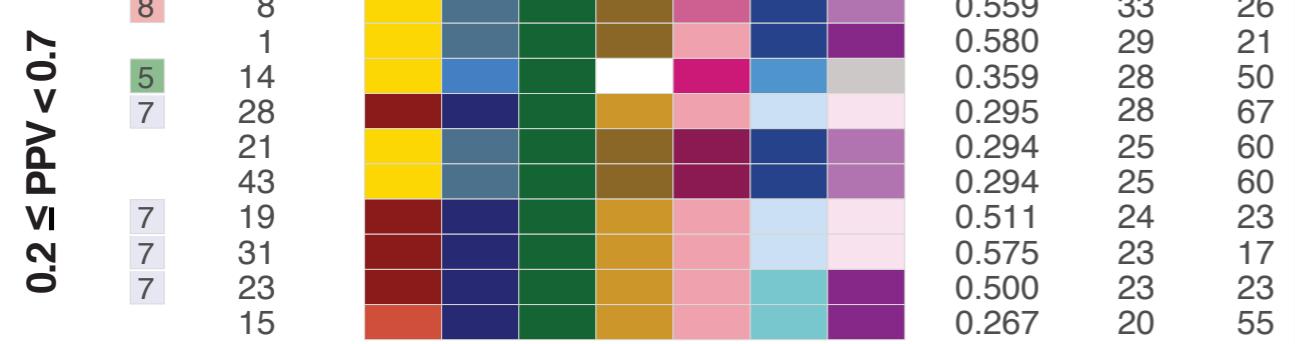
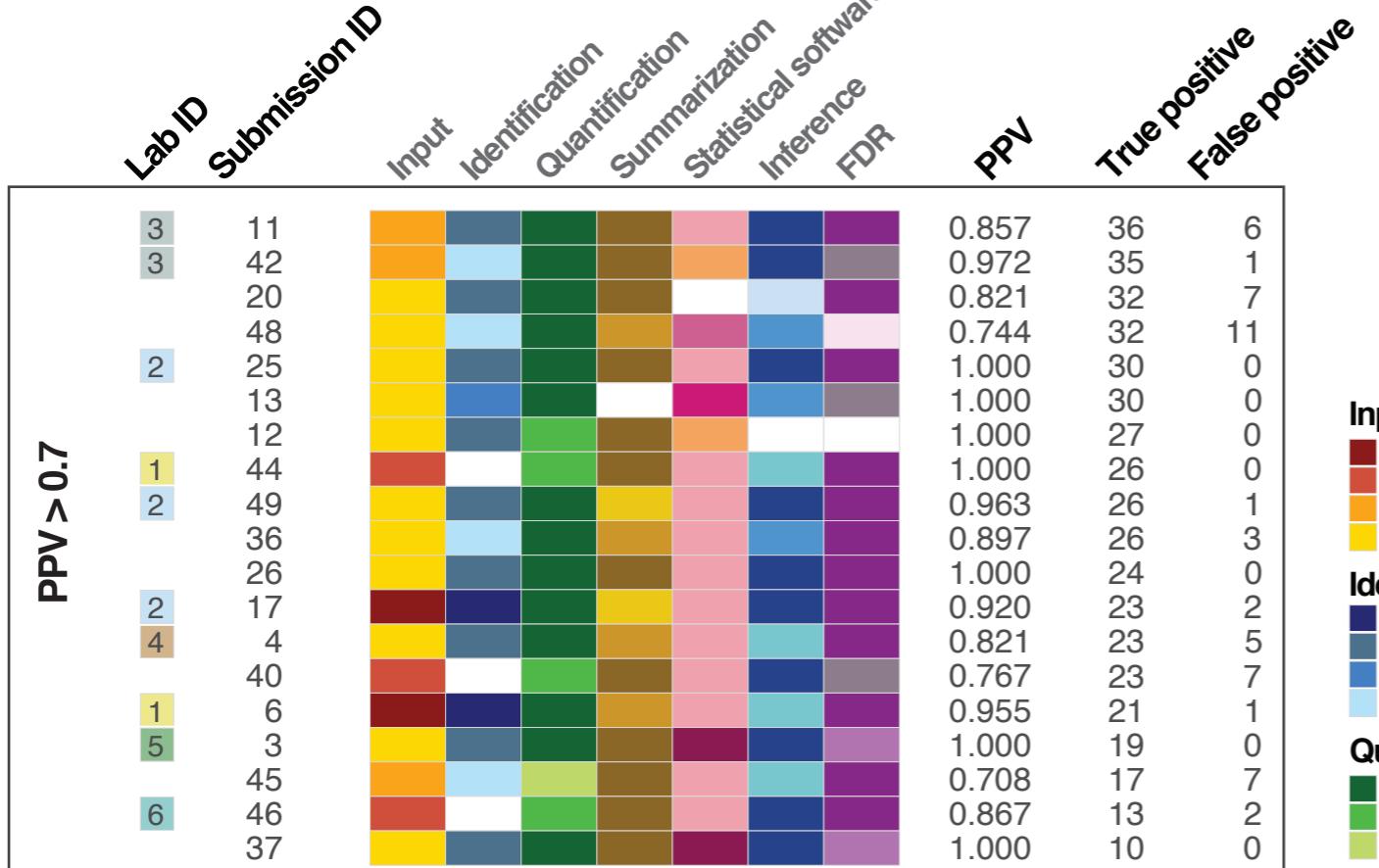
Input data

- Peaks
- Peptide IDs
- Raw+check
- Raw
- NA



SUMMARY OF SUBMISSIONS

USER EXPERTISE IS KEY



Input

- Peaks
- Peptide ids
- Raw+check
- Raw

Identification

- Skyline
- MaxQuant
- Progenesis
- Others

Quantification

- Feature intensity
- Spectral counting
- Hybrid

Summarization

- Protein summarization / Protein-level inference
- Peptide summarization / Protein-level inference
- Peptide summarization / Peptide-level inference

Statistical software

- Persus
- Progenesis QI
- Others
- R, Excel, MatLab, Python
- In-house scripts

Inference

- t-test / SAM's t test
- ANOVA
- Linear (mixed-effects) model
- Others

FDR

- Benjamini Hochberg
- Permutation FDR
- Others
- Manual validation
- FC cutoff
- No adjustment

No information

USER EXPERTISE IS KEY

PPV > 0.7	Lab ID	Submission ID	Process Flow						PPV	True positive	False positive
			Input	Identification	Quantification	Summarization	Statistical software	Inference			
3	11								0.857	36	6
3	42								0.972	35	1
20	20								0.821	32	7
48	48								0.744	32	11
25	25								1.000	30	0
13	13								1.000	30	0
12	12								1.000	27	0
44	44								1.000	26	0
2	49								0.963	26	1
36	36								0.897	26	3
26	26								1.000	24	0
17	17								0.920	23	2
4	4								0.821	23	5
40	40								0.767	23	7
6	6								0.955	21	1
5	3								1.000	19	0
45	45								0.708	17	7
6	46								0.867	13	2
37	37								1.000	10	0

Positive predictive value =

true differentially abundant proteins

claimed differentially abundant proteins

Good

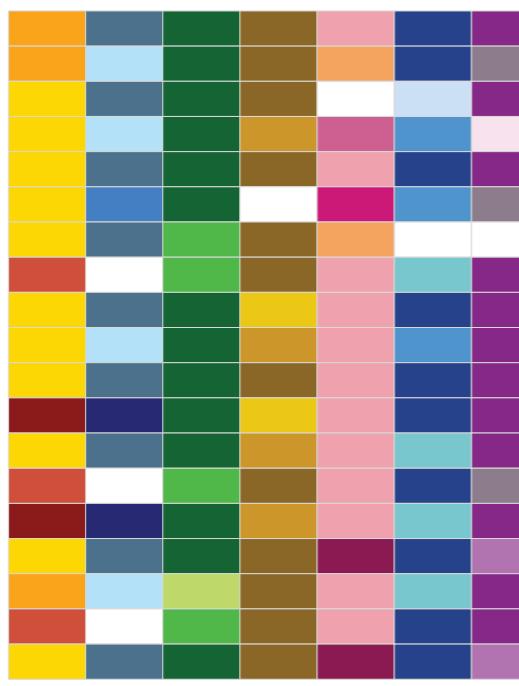
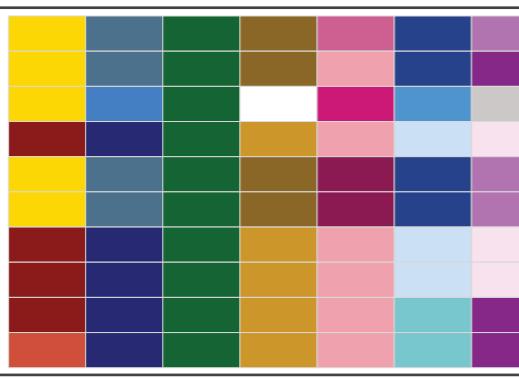
0.2 ≤ PPV < 0.7	Lab ID	Submission ID	Process Flow						PPV	True positive	False positive
			Input	Identification	Quantification	Summarization	Statistical software	Inference			
8	8								0.559	33	26
1	1								0.580	29	21
5	14								0.359	28	50
7	28								0.295	28	67
	21								0.294	25	60
	43								0.294	25	60
7	19								0.511	24	23
7	31								0.575	23	17
7	23								0.500	23	23
	15								0.267	20	55

Bad

PPV < 0.2	Lab ID	Submission ID	Process Flow						PPV	True positive	False positive
			Input	Identification	Quantification	Summarization	Statistical software	Inference			
8	30								0.174	32	152
9	16								0.145	30	177
	27								0.134	30	194
9	24								0.120	30	219
8	38								0.141	28	170
	47								0.181	27	122
	29								0.057	27	447
8	35								0.182	26	117
	35								0.126	26	180
8	22								0.145	25	148
	34								0.132	24	158
2	2								0.091	24	240
4	39								0.074	21	264
9	33								0.085	18	194
	51								0.052	18	330
10	10								0.023	17	734
5	5								0.008	17	2002
6	41								0.027	14	514
	7								0.172	5	24

Very bad

USER EXPERTISE IS KEY

	Lab ID	Submission ID	Input	Identification	Quantification	Summarization	Statistical software	Inference	FDR	PPV	True positive	False positive
PPV > 0.7 	3	11	0.857	36	6							
	3	42	0.972	35	1							
	20	0.821	32	7								
	48	0.744	32	11								
2	25	1.000	30	0								
	13	1.000	30	0								
	12	1.000	27	0								
1	44	1.000	26	0								
2	49	0.963	26	1								
	36	0.897	26	3								
	26	1.000	24	0								
2	17	0.920	23	2								
4	4	0.821	23	5								
	40	0.767	23	7								
1	6	0.955	21	1								
5	3	1.000	19	0								
	45	0.708	17	7								
6	46	0.867	13	2								
	37	1.000	10	0								
0.2 ≤ PPV < 0.7 	8	8	0.559	33	26							
	1	0.580	29	21								
5	14	0.359	28	50								
7	28	0.295	28	67								
	21	0.294	25	60								
	43	0.294	25	60								
7	19	0.511	24	23								
7	31	0.575	23	17								
7	23	0.500	23	23								
	15	0.267	20	55								
PPV < 0.2 	8	30	0.174	32	152							
9	16	0.145	30	177								
	27	0.134	30	194								
9	24	0.120	30	219								
8	38	0.141	28	170								
	47	0.181	27	122								
	29	0.057	27	447								
8	35	0.182	26	117								
	22	0.126	26	180								
8	34	0.145	25	148								
	2	0.132	24	158								
4	39	0.091	24	240								
9	33	0.074	21	264								
	9	0.085	18	194								
	51	0.052	18	330								
	10	0.023	17	734								
	5	0.008	17	2002								
6	41	0.027	14	514								
	7	0.172	5	24								

MaxQuant and Perseus

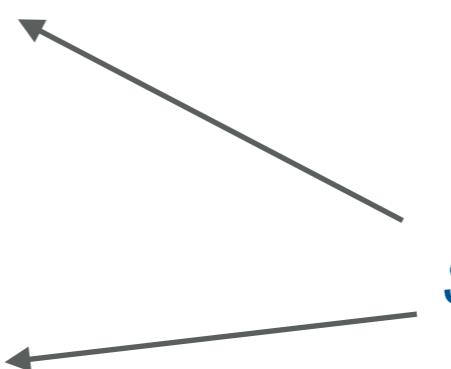
USER EXPERTISE IS KEY

Lab ID	Submission ID	Input	Identification	Quantification	Summarization	Statistical software	Inference	FDR	PPV	True positive	False positive
3	11								0.857	36	6
3	42								0.972	35	1
20									0.821	32	7
48									0.744	32	11
2	25								1.000	30	0
13									1.000	30	0
12									1.000	27	0
1	44								1.000	26	0
2	49								0.963	26	1
36									0.897	26	3
2	26								1.000	24	0
4	17								0.920	23	2
4	4								0.821	23	5
40									0.767	23	7
1	6								0.955	21	1
5	3								1.000	19	0
45									0.708	17	7
6	46								0.867	13	2
	37								1.000	10	0

PPV > 0.7	PPV ≤ 0.7	PPV < 0.2
8	8	30
1	1	16
5	14	27
7	28	24
	21	38
	43	47
7	19	29
7	31	35
7	23	35
	15	22

PPV < 0.2	PPV > 0.7	PPV ≤ 0.7
8	30	8
9	16	16
	27	17
9	24	24
8	38	27
	47	27
8	29	38
8	35	47
8	35	47
8	22	35
8	34	35
2	34	34
4	2	34
9	39	39
9	33	33
	9	33
51		51
10		51
5		51
41		41
6	7	7

Skyline and linear modeling in R



USER EXPERTISE IS KEY

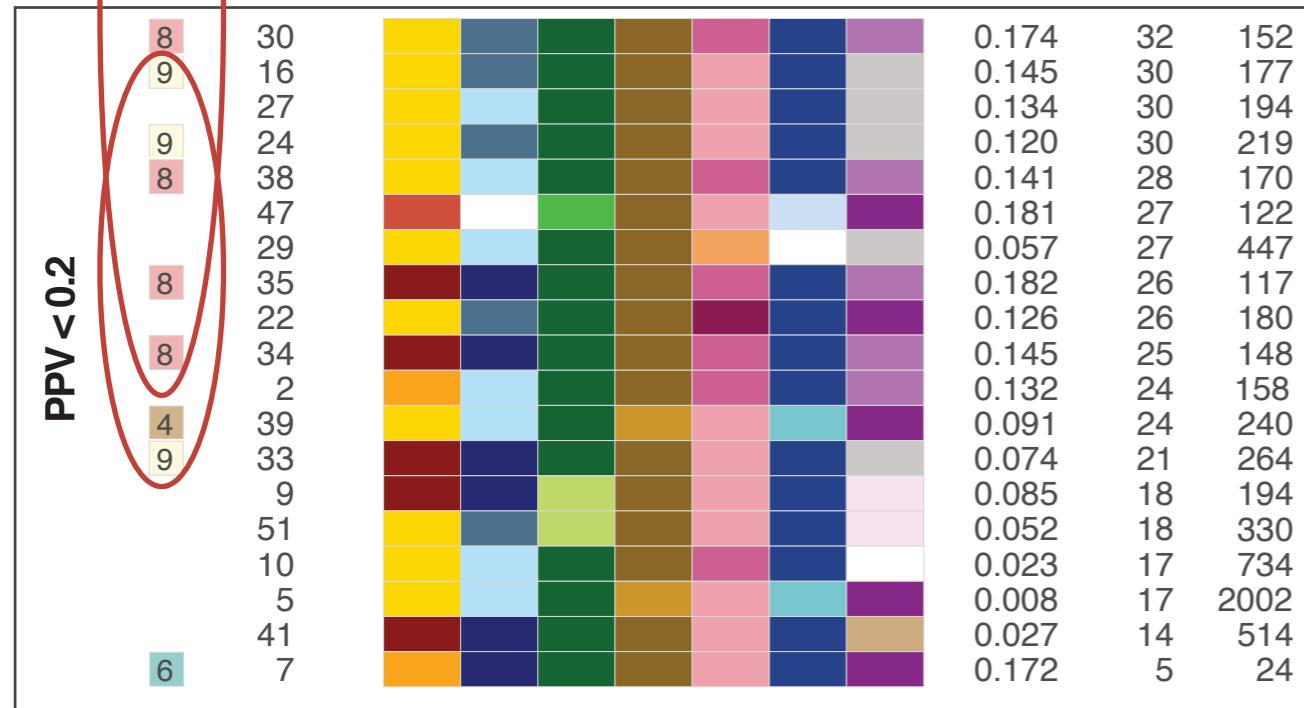
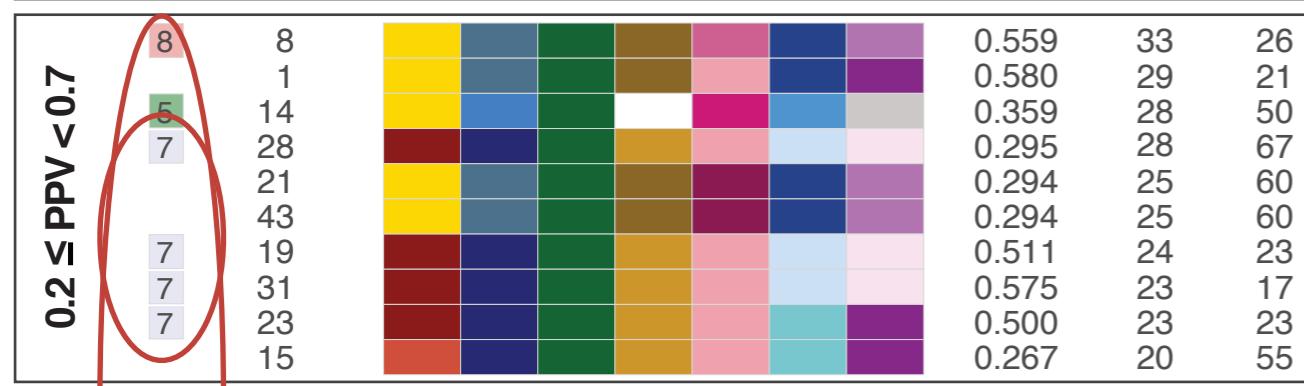
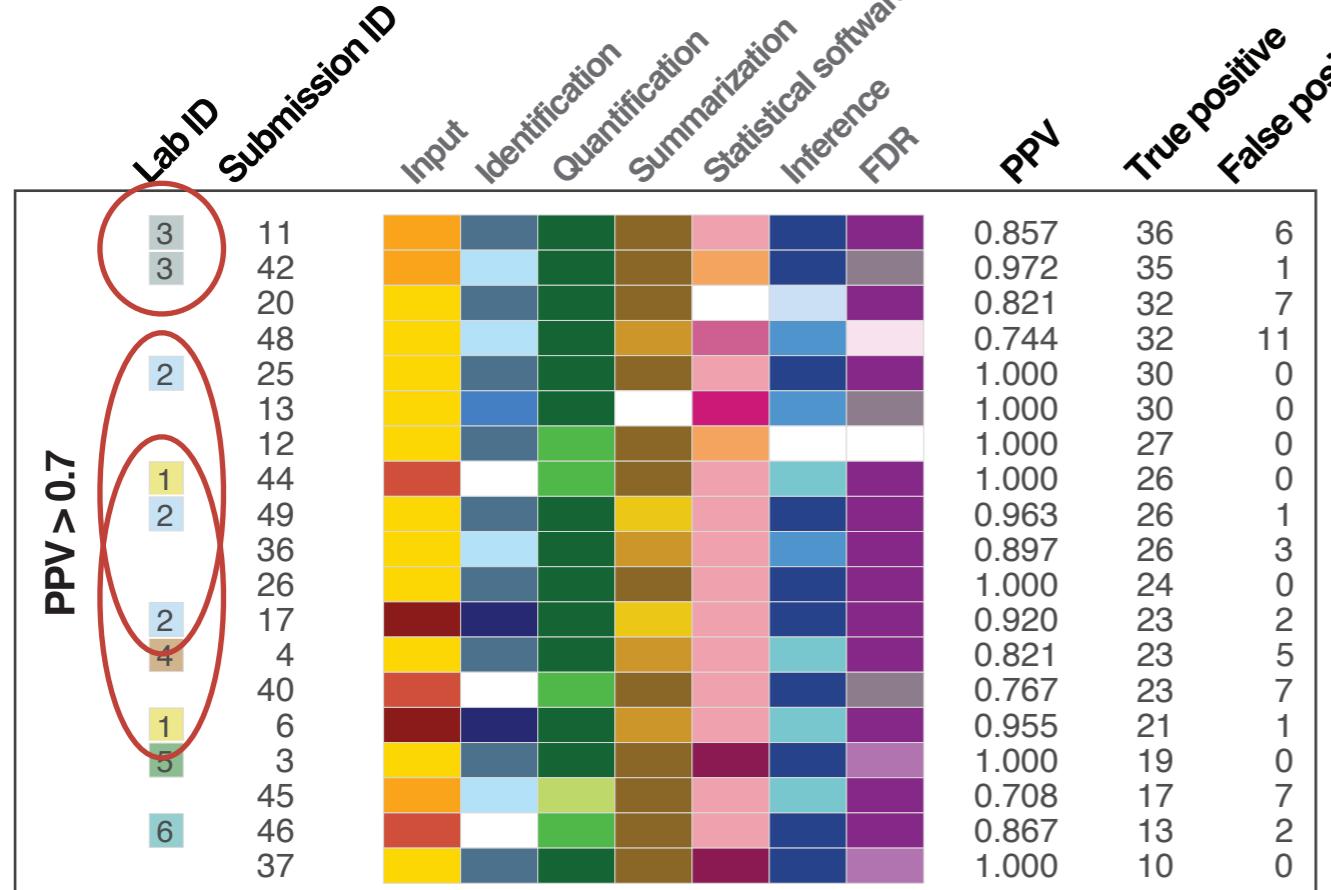
Lab ID	Submission ID	Input	Identification	Quantification	Summarization	Statistical software	Inference	FDR	PPV	True positive	False positive
3	11								0.857	36	6
3	42								0.972	35	1
20	20								0.821	32	7
48	48								0.744	32	11
2	25								1.000	30	0
13	13								1.000	30	0
12	12								1.000	27	0
1	44								1.000	26	0
2	49								0.963	26	1
36	36								0.897	26	3
2	26								1.000	24	0
4	17								0.920	23	2
4	4								0.821	23	5
40	40								0.767	23	7
1	6								0.955	21	1
5	3								1.000	19	0
6	45								0.708	17	7
6	46								0.867	13	2
	37								1.000	10	0

PPV > 0.7	PPV ≤ 0.7	PPV < 0.2
8	8	30
1	1	16
5	14	27
7	28	24
	21	38
	43	47
7	19	29
7	31	35
7	23	35
	15	22

PPV < 0.2	PPV > 0.7	PPV ≤ 0.7
8	30	8
9	16	16
	27	27
9	24	24
8	38	38
	47	47
8	29	29
8	35	35
8	35	35
8	22	22
8	34	34
2	2	2
4	39	39
9	33	33
	9	9
51	51	51
10	10	10
5	5	5
41	41	41
6	7	7

Compared peak intensity vs spectral counts

USER EXPERTISE IS KEY



Input

- Peaks
- Peptide ids
- Raw+check
- Raw

Identification

- Skyline
- MaxQuant
- Progenesis
- Others

Quantification

- Feature intensity
- Spectral counting
- Hybrid

Summarization

- Protein summarization / Protein-level inference
- Peptide summarization / Protein-level inference
- Peptide summarization / Peptide-level inference

Statistical software

- Persus
- Progenesis QI
- Others
- R, Excel, MatLab, Python
- In-house scripts

Inference

- t-test / SAM's t test
- ANOVA
- Linear (mixed-effects) model
- Others

FDR

- Benjamini Hochberg
- Permutation FDR
- Others
- Manual validation
- FC cutoff
- No adjustment

No information

Article

[!\[\]\(e5d4c1253f90f386527cfb2278e2ccef_img.jpg\) Previous Article](#)[!\[\]\(2c3352433bff267ed8ae00945ed009eb_img.jpg\) Next Article](#)[Table of Contents](#)

ABRF Proteome Informatics Research Group (iPRG) 2015 Study: Detection of Differentially Abundant Proteins in Label-Free Quantitative LC–MS/MS Experiments

Meena Choi^{#†} , Zeynep F. Eren-Dogu^{#‡}, Christopher Colangelo[§], John Cottrell[¶], Michael R. Hoopmann[⊥], Eugene A. Kapp[¶], Sangtae Kim[®], Henry Lam[□], Thomas A. Neubert^{*}, Magnus Palmblad[○], Brett S. Phinney^{*}, Susan T. Weintraub[△], Brendan MacLean[▲], and Olga Vitek^{#†} 

[#] Northeastern University, Boston, Massachusetts 02115, United States

[†] Mugla Sitki Koçman University, 48000 Mugla, Turkey

[‡] Primary Ion, LLC, Old Lyme, Connecticut 06371, United States

[¶] Matrix Science Ltd., London W1U 7GB, U.K.

[○] Institute for Systems Biology, Seattle, Washington 98109, United States

[○] Walter and Eliza Hall Institute of Medical Research, Melbourne 3052, Australia

[○] Pacific Northwest National Laboratory, Richland, Washington 99354, United States

[□] Department of Chemical and Biomolecular Engineering and Division of Biomedical Engineering, The Hong Kong University of Science and Technology, Clear Water Bay, Hong Kong

[■] Skirball Institute and Department of Biochemistry and Molecular Pharmacology, New York University School of Medicine, New York, New York 10016, United States

[○] Center for Proteomics and Metabolomics, Leiden University Medical Center, 2333 ZA Leiden, The Netherlands

[●] University of California at Davis, Davis, California 95616, United States

[△] University of Texas Health Science Center at San Antonio, San Antonio, Texas 78229, United States

[▲] University of Washington, Seattle, Washington 98105, United States

J. Proteome Res., 2017, 16 (2), pp 945–957

DOI: 10.1021/acs.jproteome.6b00881

Publication Date (Web): December 19, 2016

Copyright © 2016 American Chemical Society

*E-mail: o.vitek@neu.edu. Tel: 617-373-2194.



Article Options

 ACS ActiveView PDF

Hi-Res Print, Annotate, Reference QuickView

[Abstract](#)

[Supporting Info](#)

 PDF (3135 KB)

[Figures](#)

 PDF w/ Links (886 KB)

[References](#)

 Full Text HTML

 Add to ACS ChemWorx

 Add to Favorites

 Download Citation

 Email a Colleague

 Order Reprints

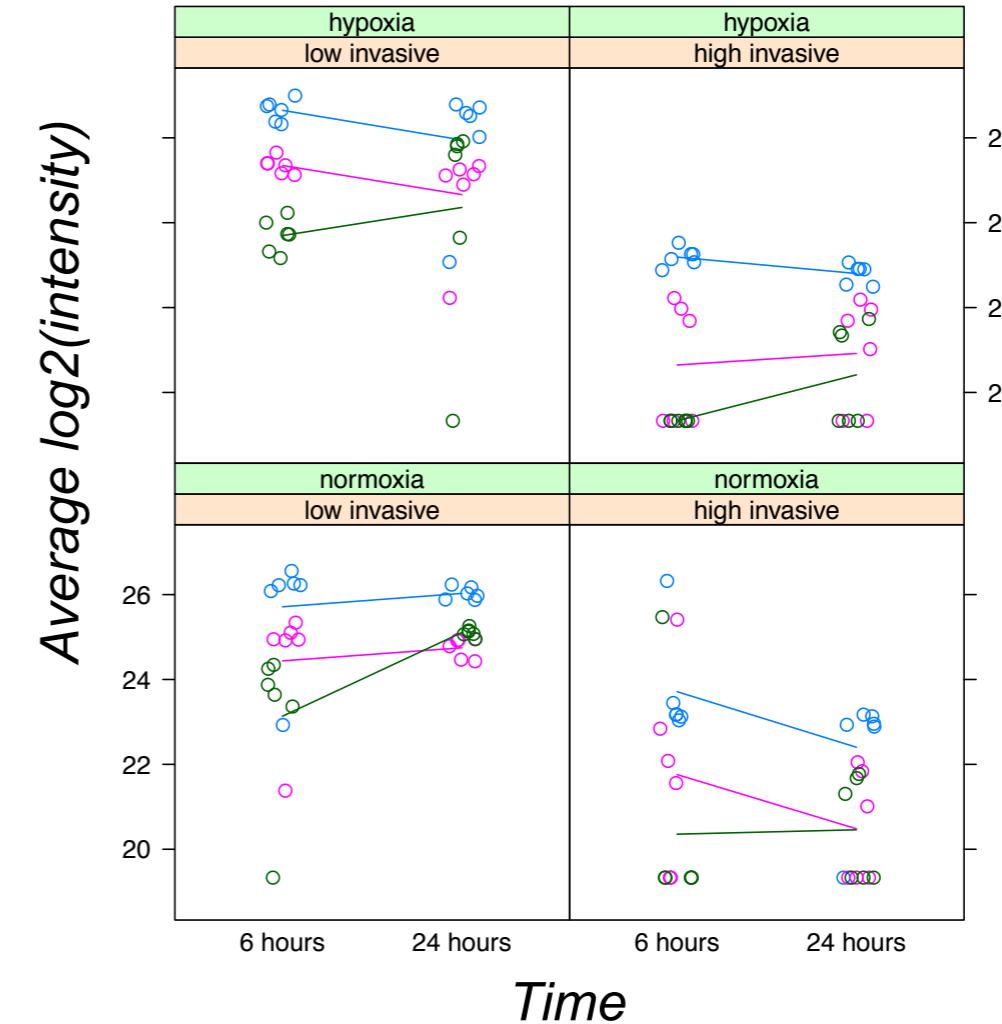
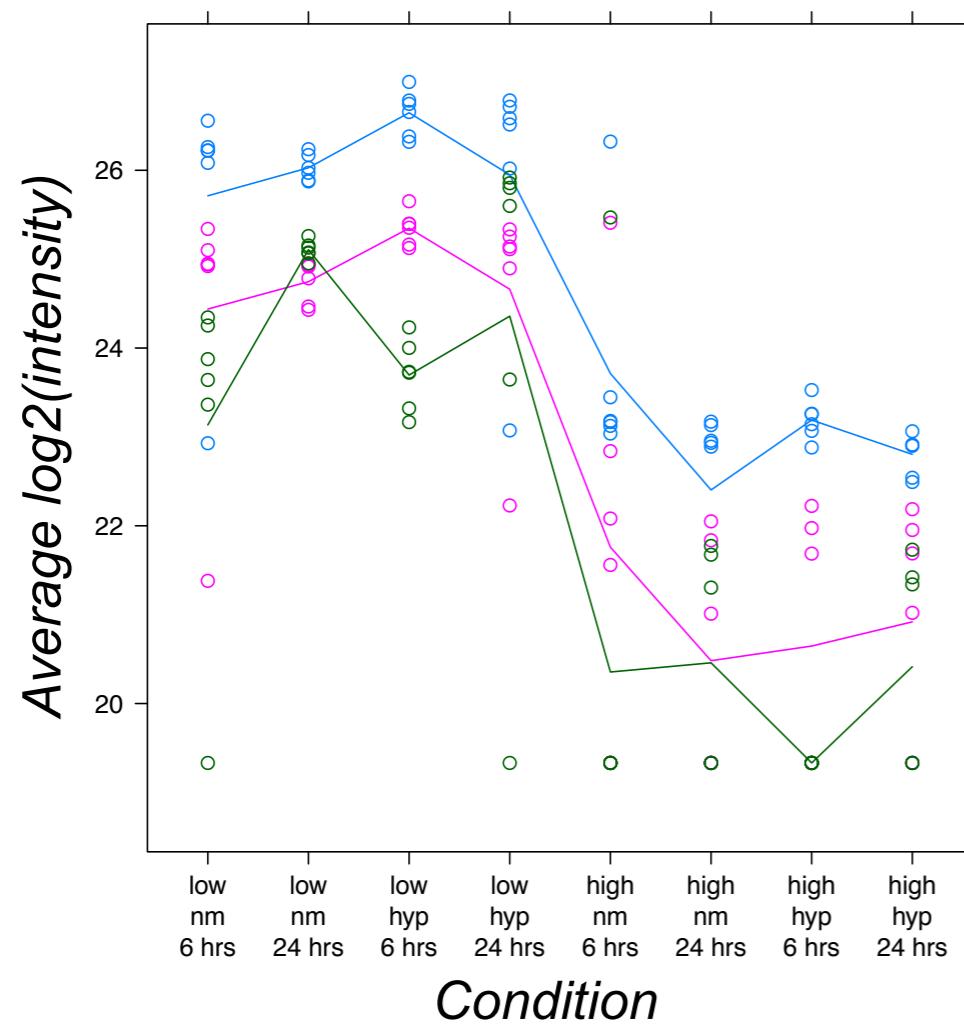
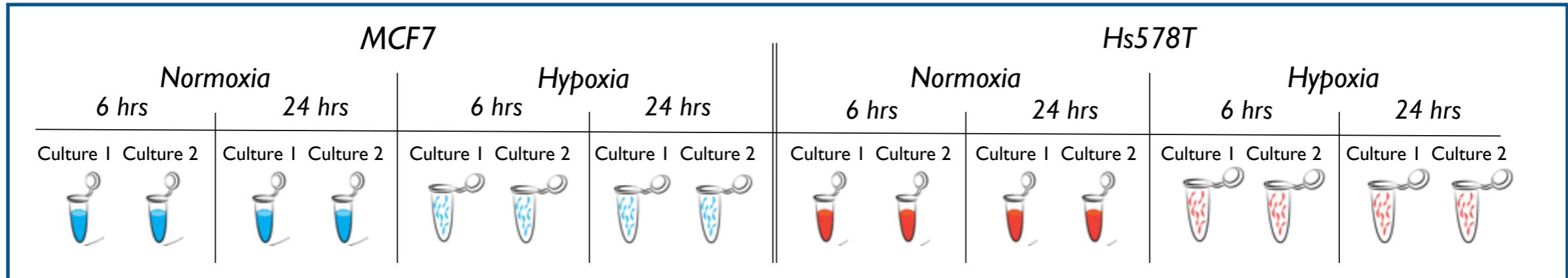
 Rights & Permissions

OUTLINE

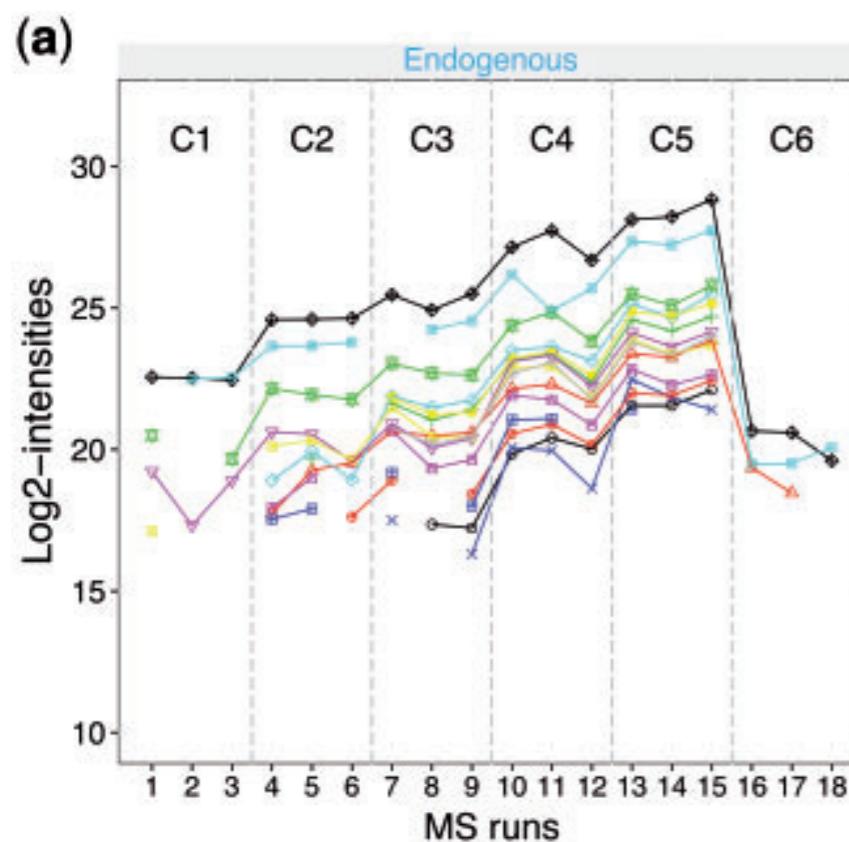
- iPRG2015: detection of diff. abundance
 - Community effort with label-free shotgun proteomics
- MSstats
 - Normalization, statistical modeling, inference
 - Evaluation
- Extensions to MSstats
 - Assay characterization, longitudinal monitoring

EXAMPLE: A LABEL-FREE EXPERIMENT

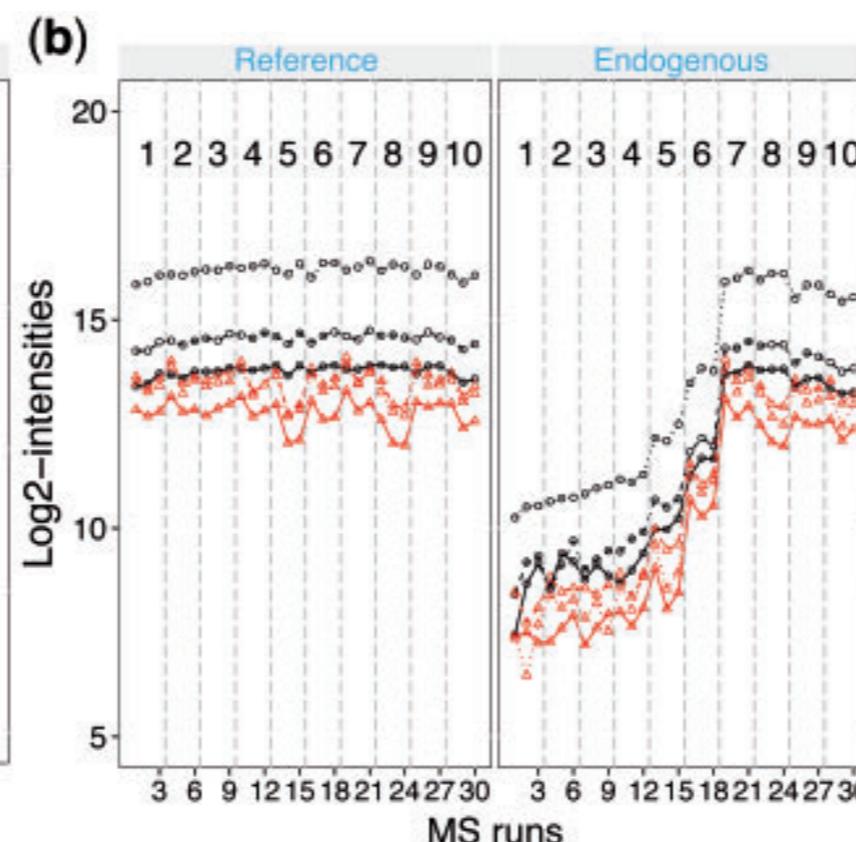
Question: which proteins change in abundance?



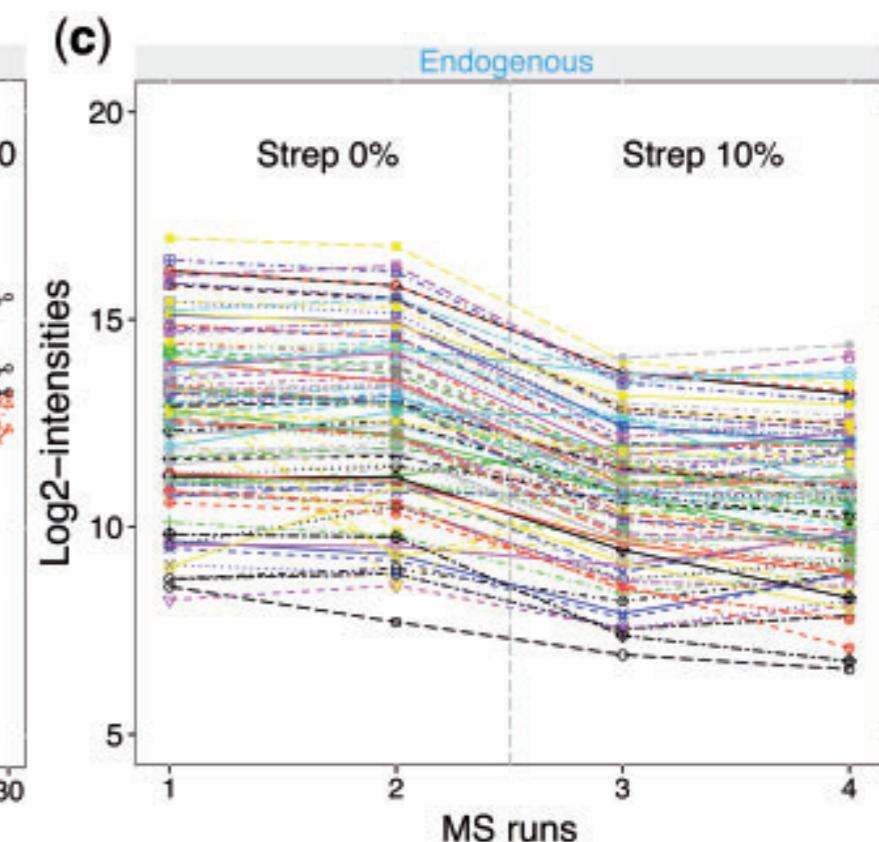
CHROMATOGRAPHY-BASED QUANTIFICATION YIELDS DATA WITH VARYING PROPERTIES



Data-dependent acquisition (DDA)



Selected reaction monitoring with labeled reference peptides (SRM)



Data-independent acquisition (DIA)

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

SCHEMATIC DATA REPRESENTATION

Repeat for every protein

*Conditions, subjects and runs:
biological aspects of the experiment*



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	...	y	NA	y

*Spectral features:
technological aspects of the
experiment*

Missing values

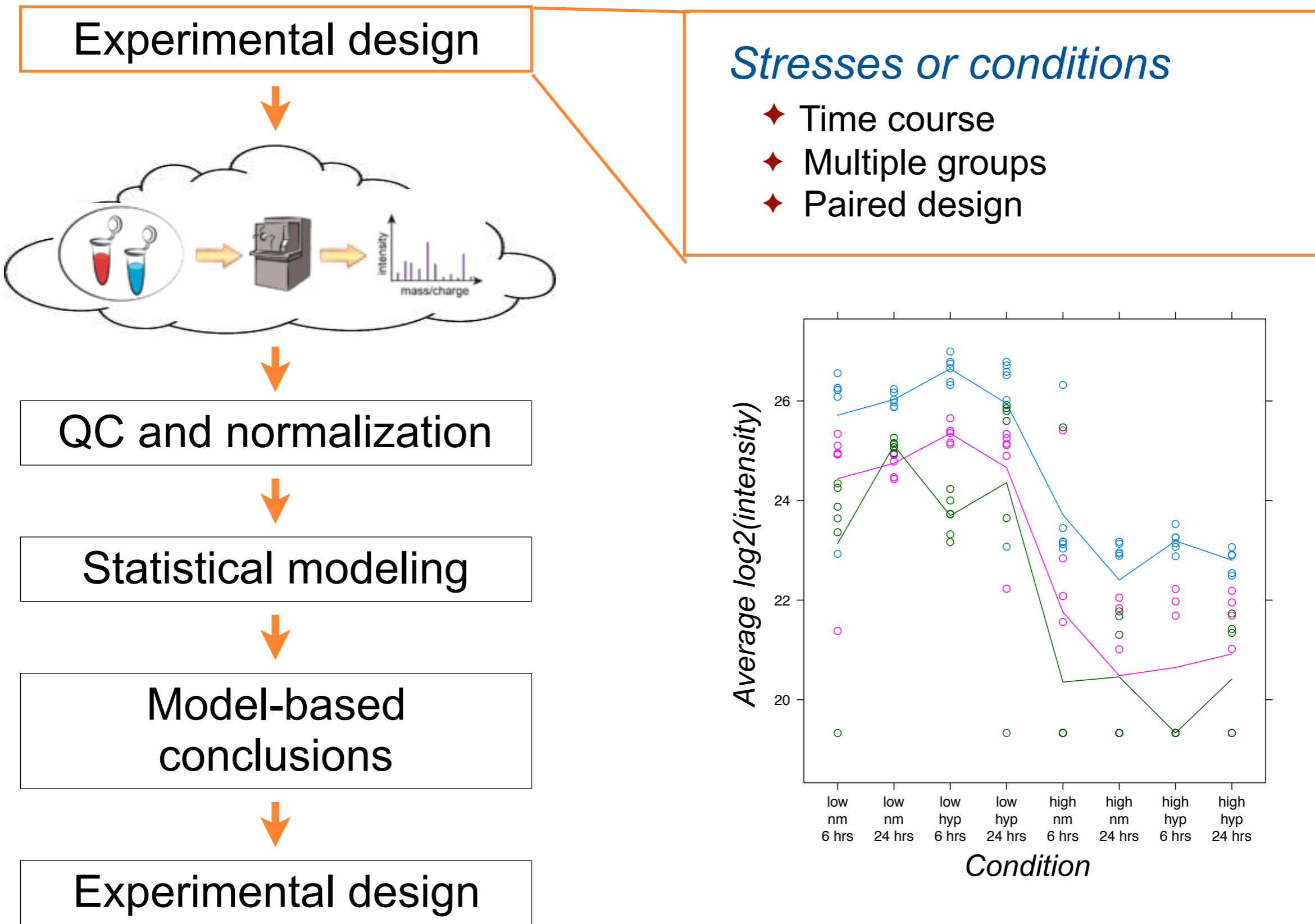
*Log(feature
intensities)*



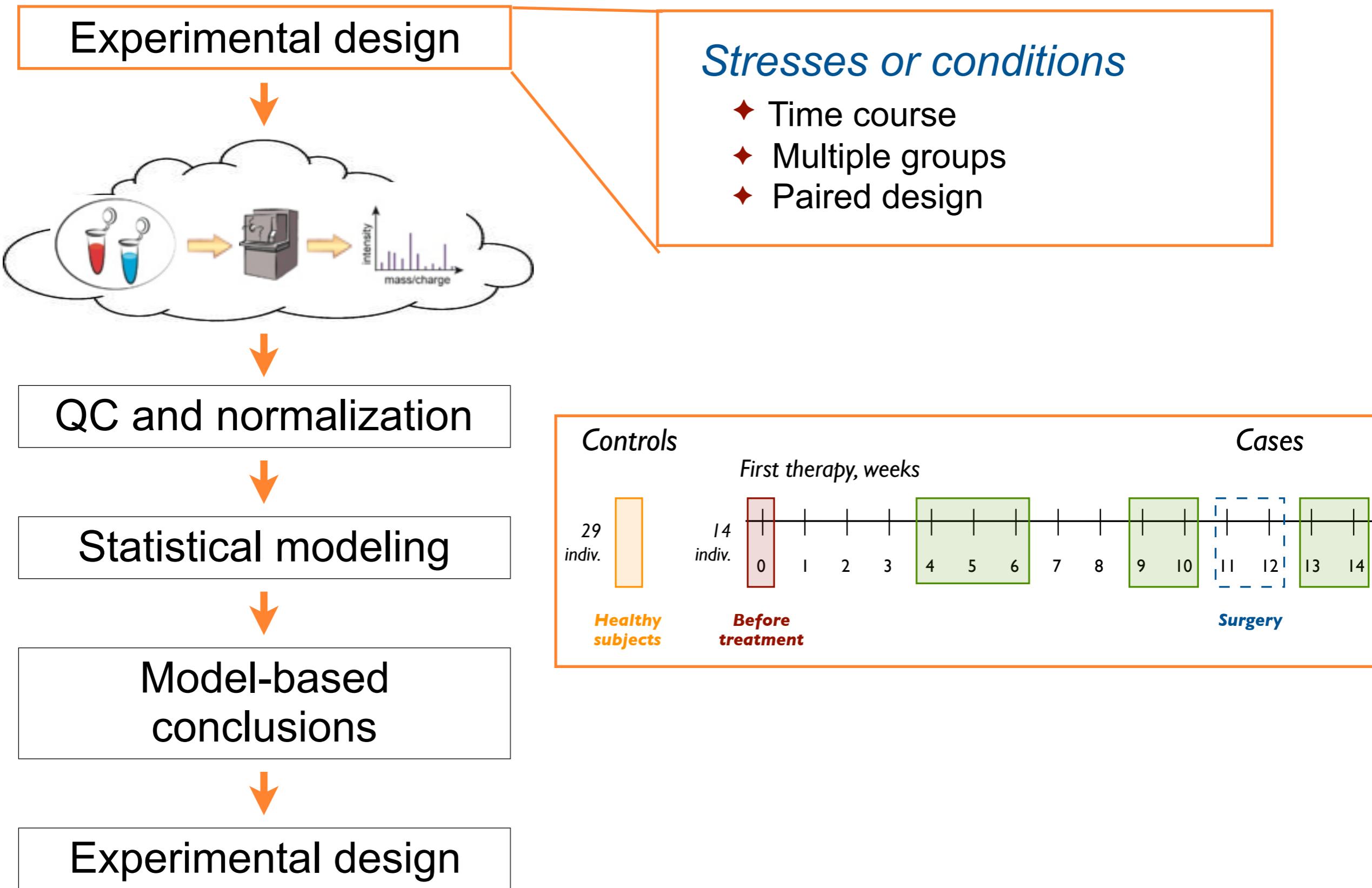
INPUT DATA REPRESENTATION

	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

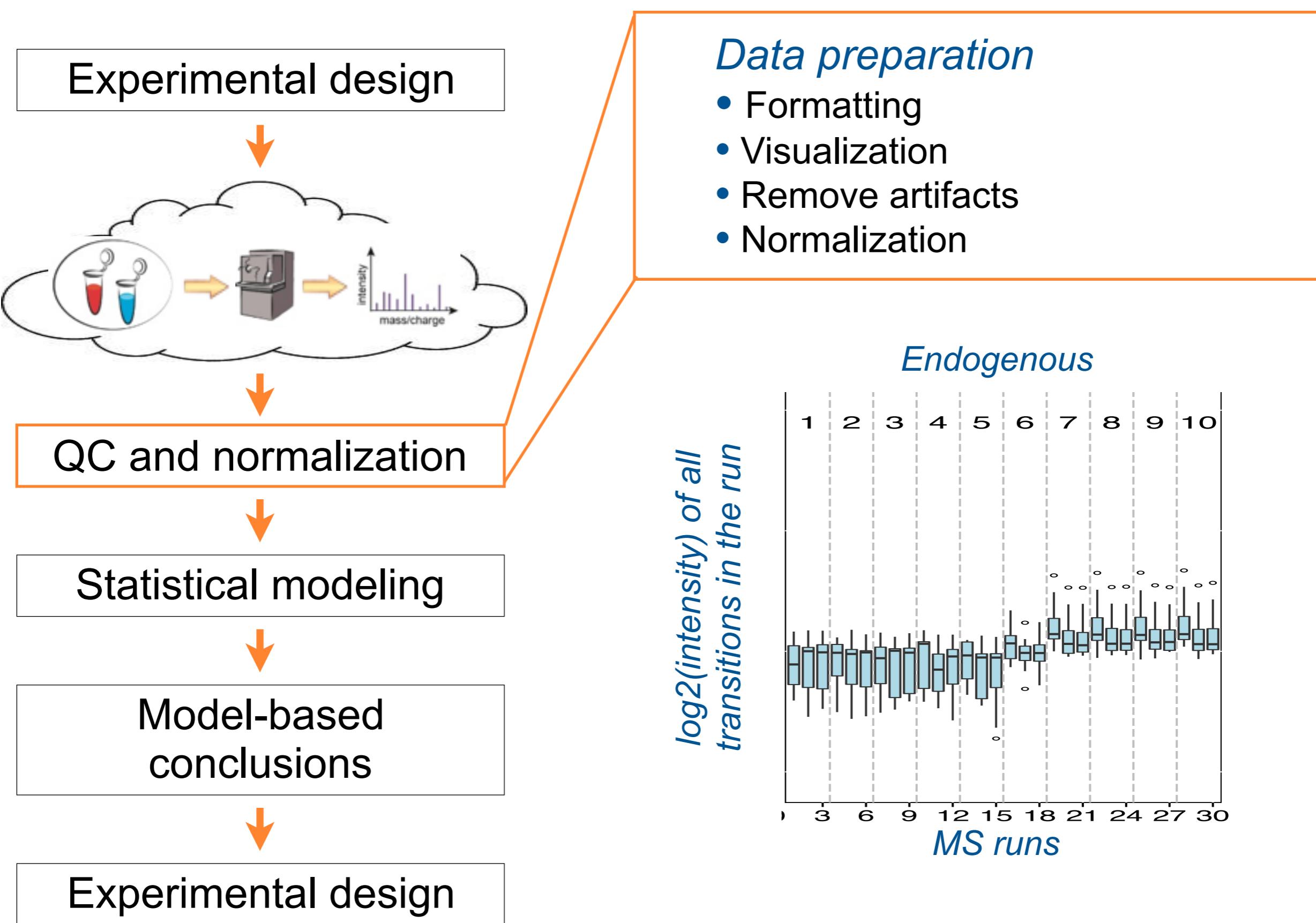
A TYPICAL ANALYSIS WORKFLOW



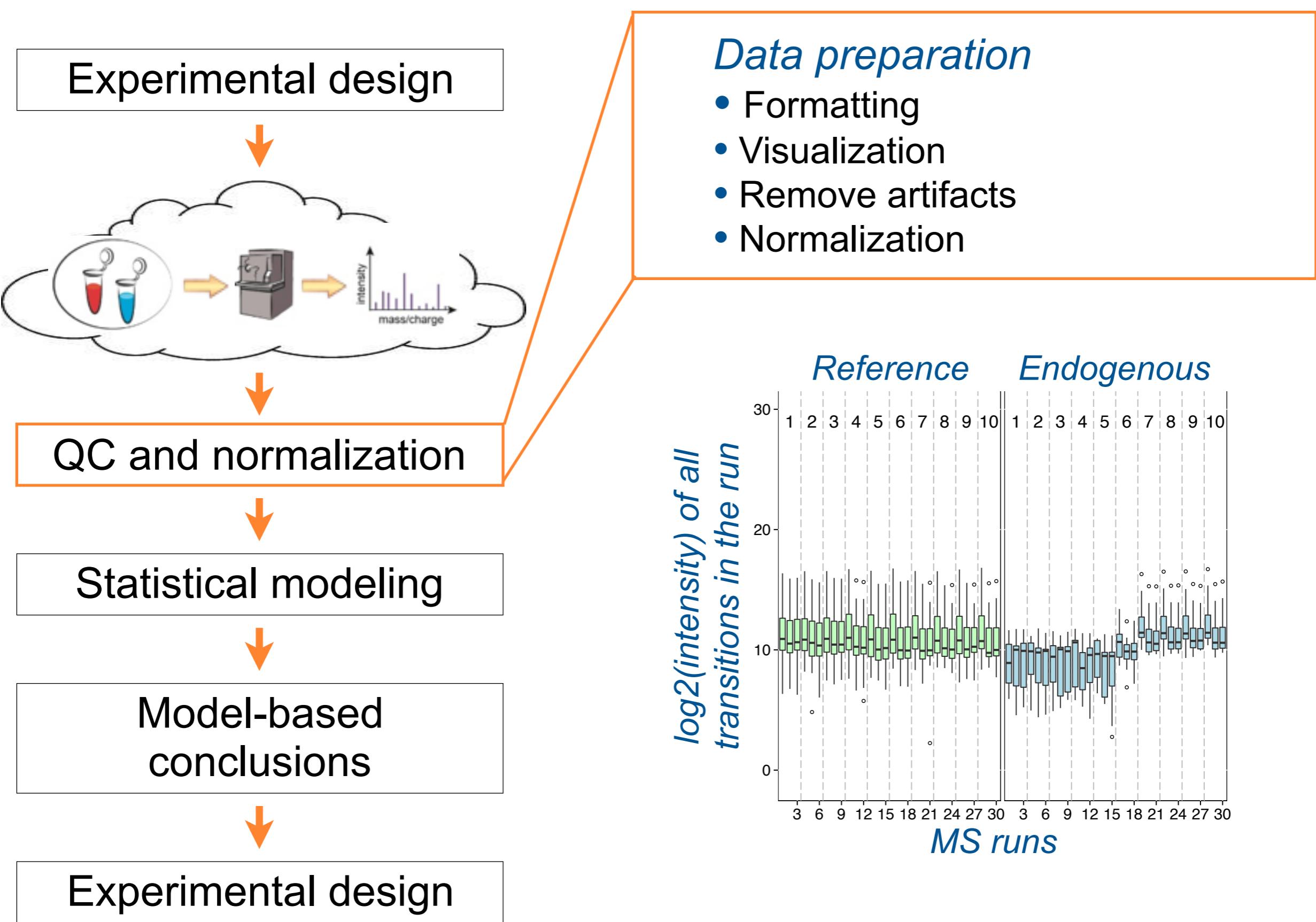
A TYPICAL ANALYSIS WORKFLOW



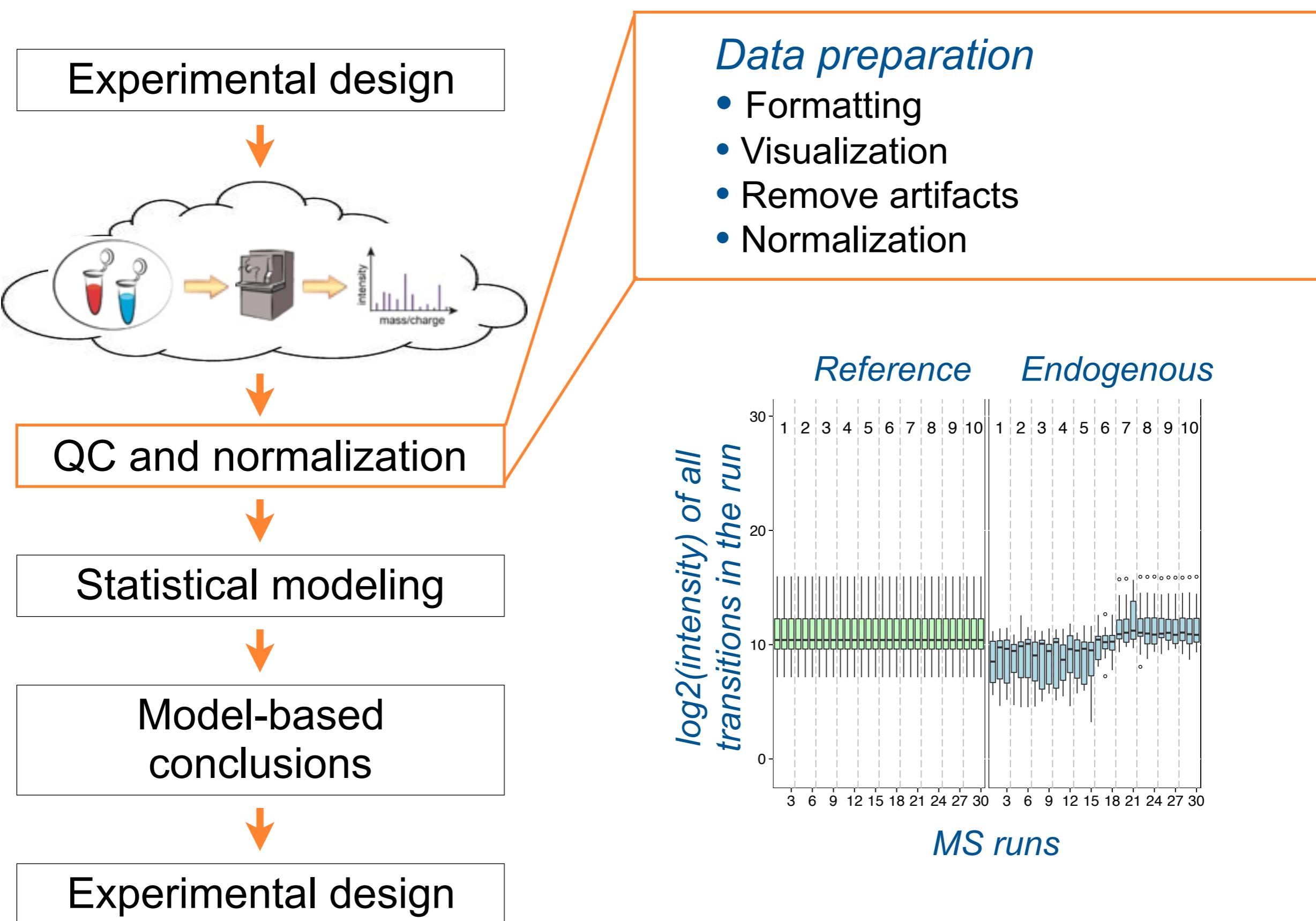
A TYPICAL ANALYSIS WORKFLOW



A TYPICAL ANALYSIS WORKFLOW

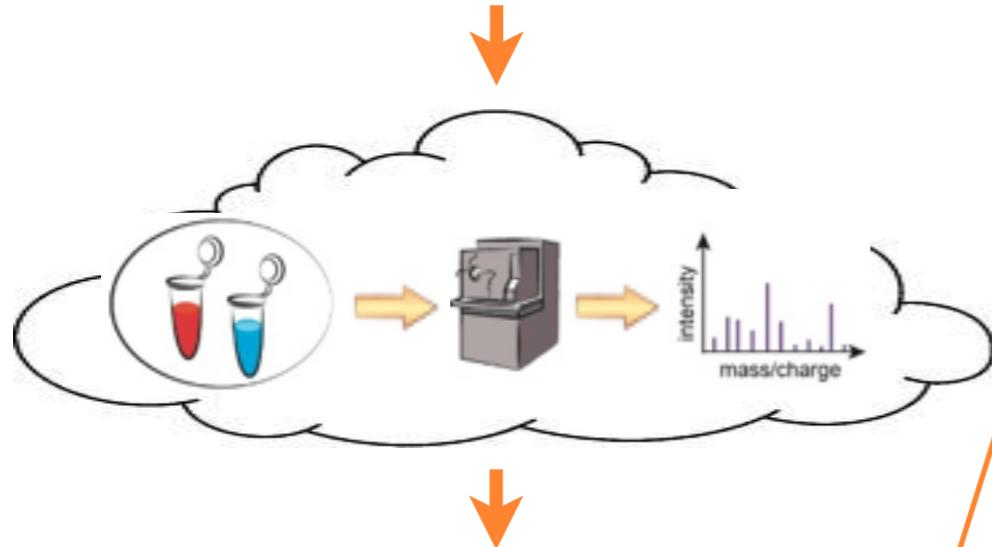


A TYPICAL ANALYSIS WORKFLOW



A TYPICAL ANALYSIS WORKFLOW

Experimental design

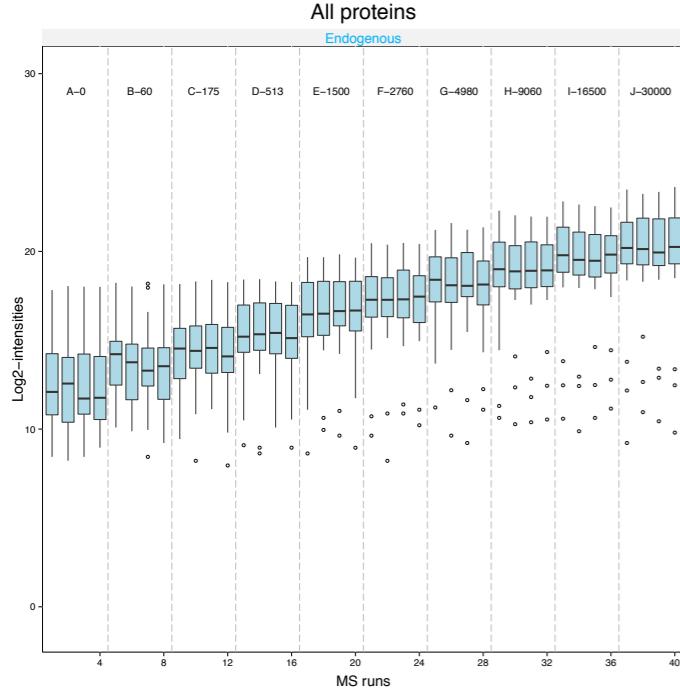


Data preparation

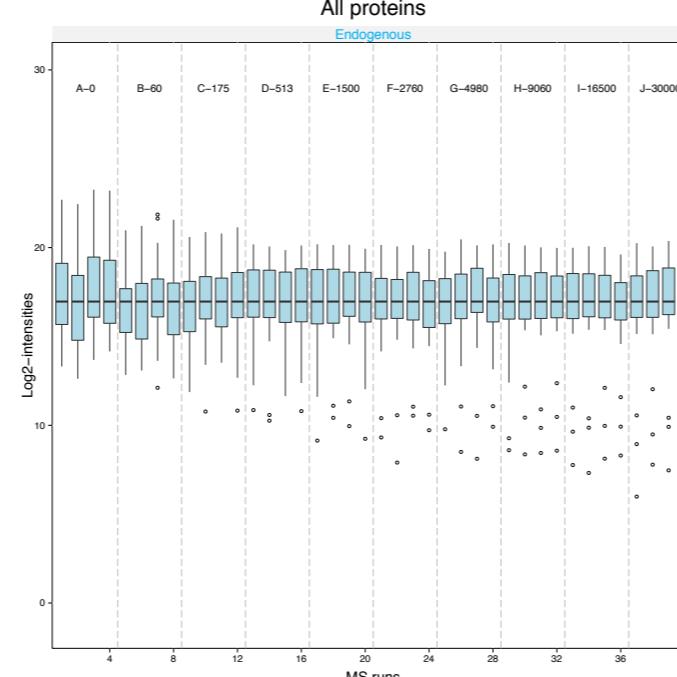
- Formatting
- Visualization
- Remove artifacts
- Normalization

QC and normalization

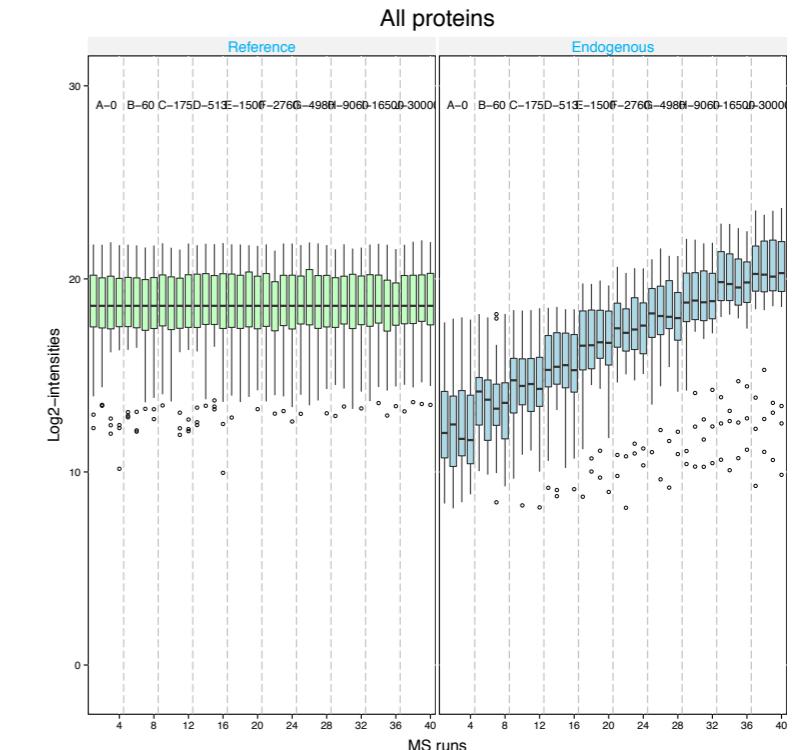
No normalization



Equalize medians normalization



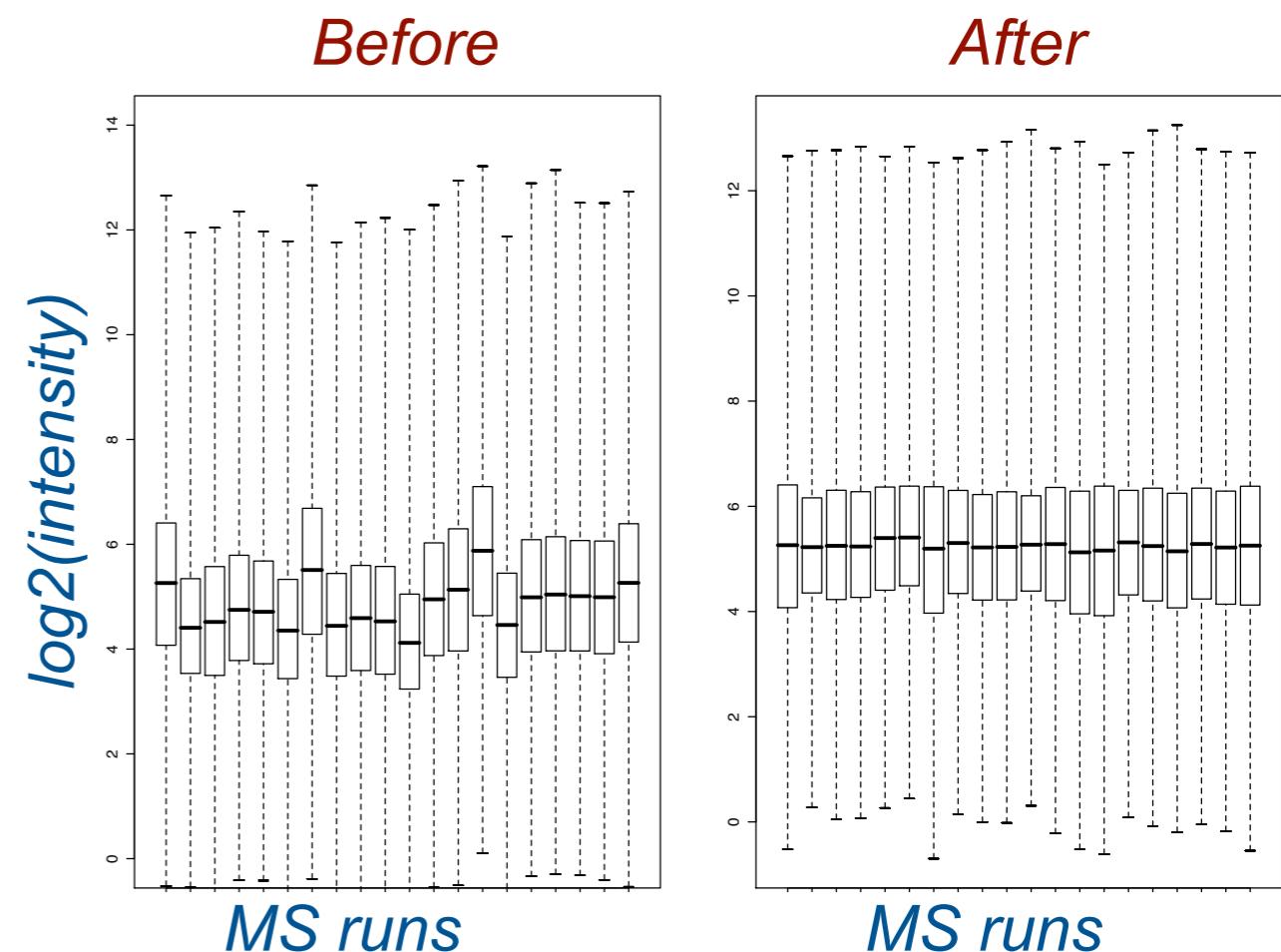
Equalize medians normalization



STANDARD-BASED NORMALIZATION

Assumes constant abundance of features from a standard

- Algorithm
 - Subtract median[log(feature int)] of the standard
 - Add back the median of the medians
- Comments
 - + Standard does not have biological variation
 - + Independent of type/number of features
 - Only accounts for deviations that occur after adding the standard
 - Standards can be noisy (unequal spiked abundance; overlapped peaks)
- Best practice
 - Use one standard for normalization, another for verification

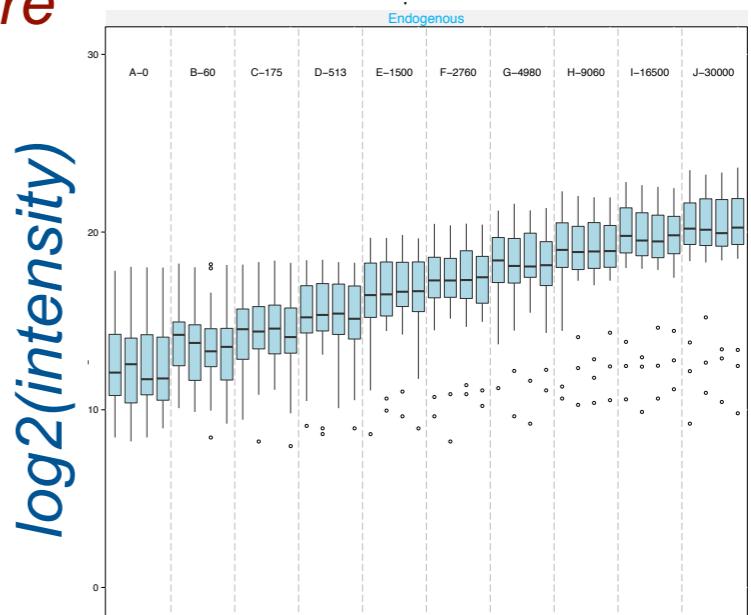


MEDIAN NORMALIZATION

Assumes constant abundance of median of $\log(\text{int})$ of endogenous features

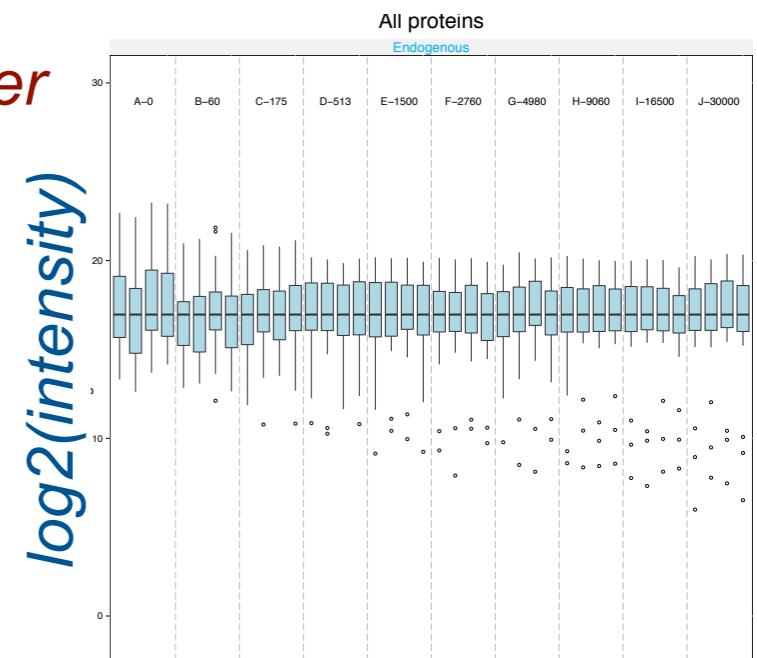
- Algorithm
 - Subtract median[$\log(\text{feature int})$] of all endogenous features
 - Add back the median of the medians
- Comments
 - + More stable than a single standard
 - + Accounts for all data processing steps
 - Assumes that the majority of endogenous proteins are not affected by the conditions
- Best practice
 - Use in discovery experiments with many features

Before



MS runs

After

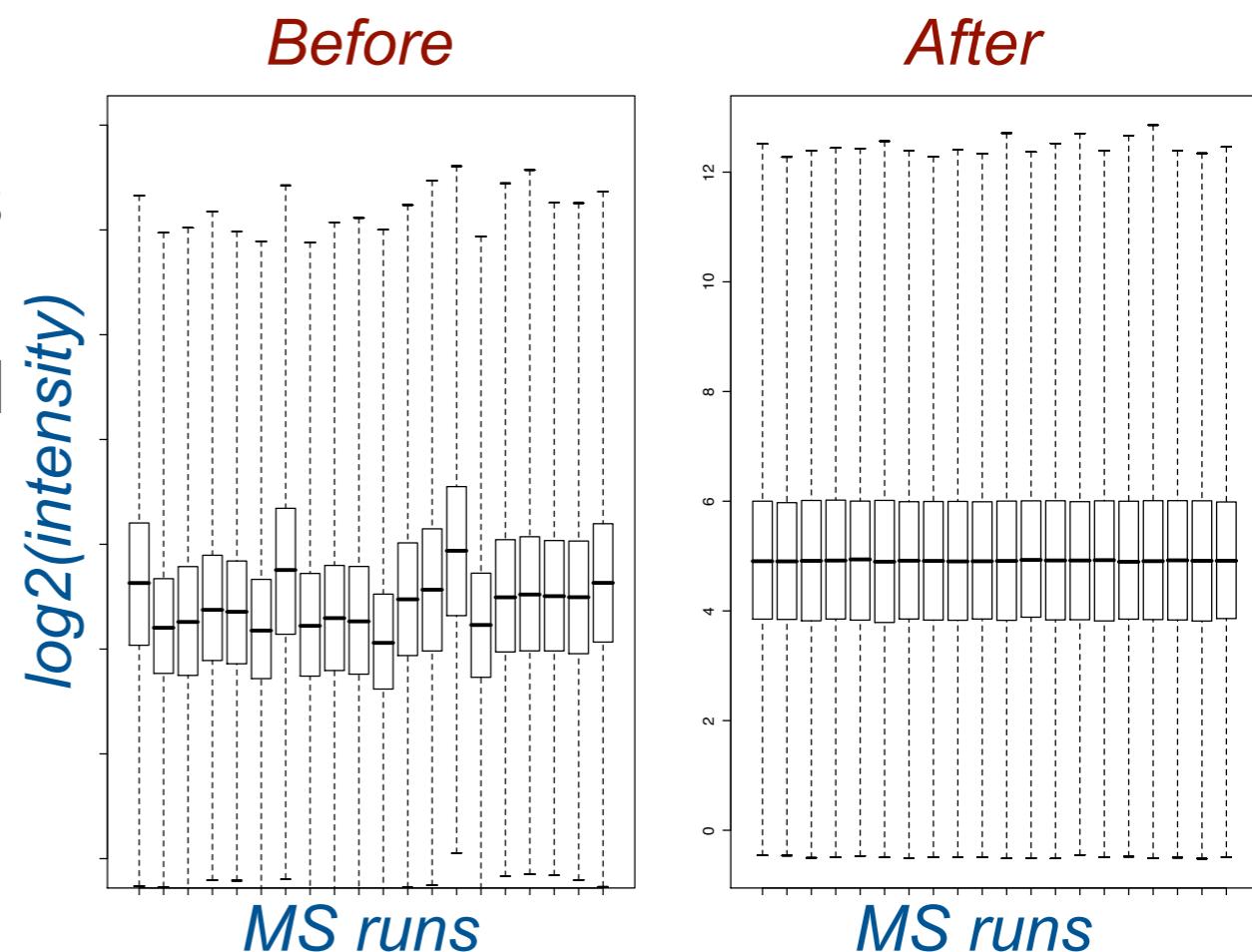


MS runs

QUANTILE NORMALIZATION

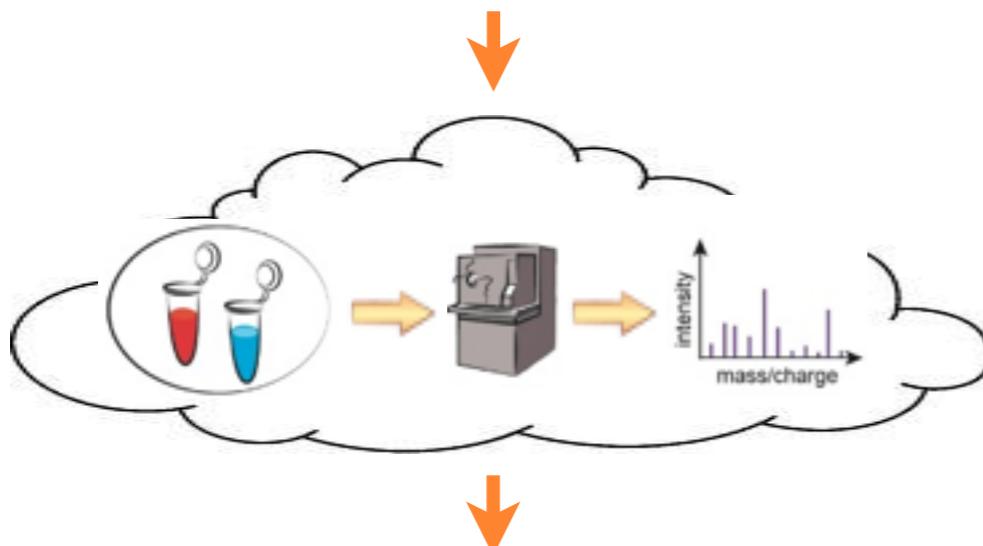
Assumes constant abundance of all quantiles of $\log(\text{int})$ of endogenous features

- Algorithm
 - Order $\log(\text{int})$ in each run
 - Calculate average of each quantile across runs
 - Substitute each $\log(\text{int})$ with the average
 - Re-arrange $\log(\text{int})$ in original order
- Comments
 - + More stable than a single standard
 - + Accounts for all data processing steps
 - Assumes that the majority of endogenous proteins are not affected by the conditions
 - Often very aggressive
- Best practice
 - Use to normalize multiple standards



A TYPICAL ANALYSIS WORKFLOW

Experimental design



Summarize all protein features in a statistical model

- Systematic variation
- Random variation

Verify the assumptions!

QC and normalization

Statistical modeling

Model-based conclusions

Experimental design

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

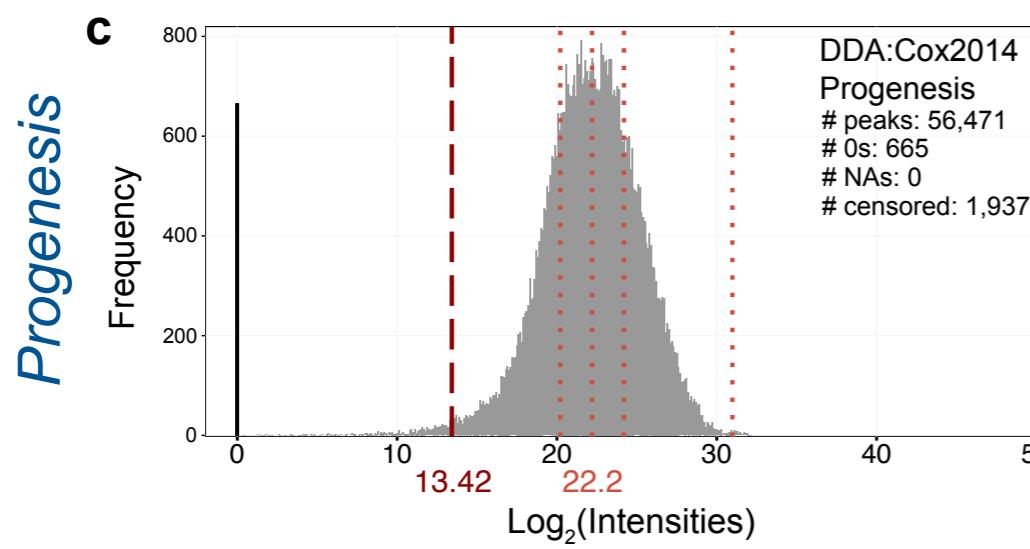
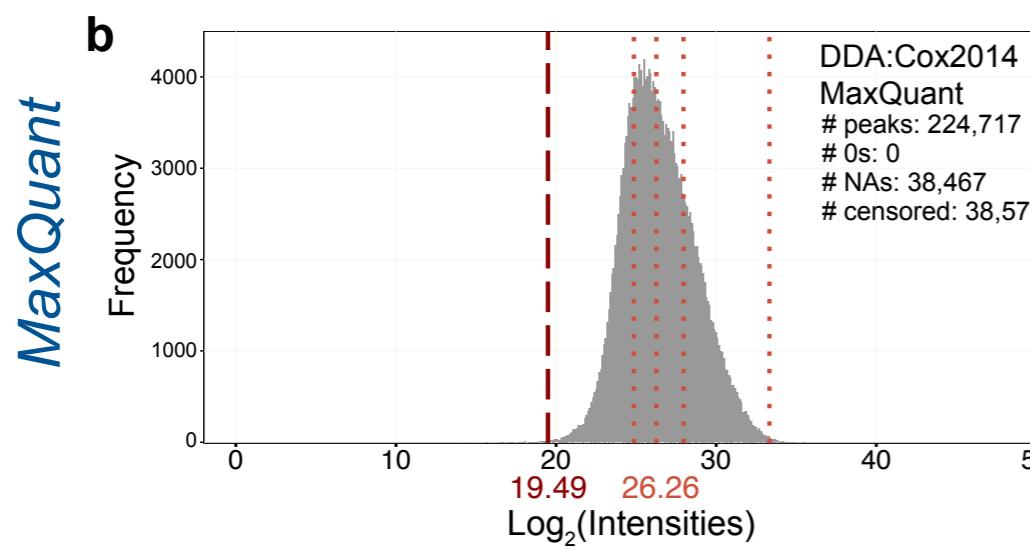
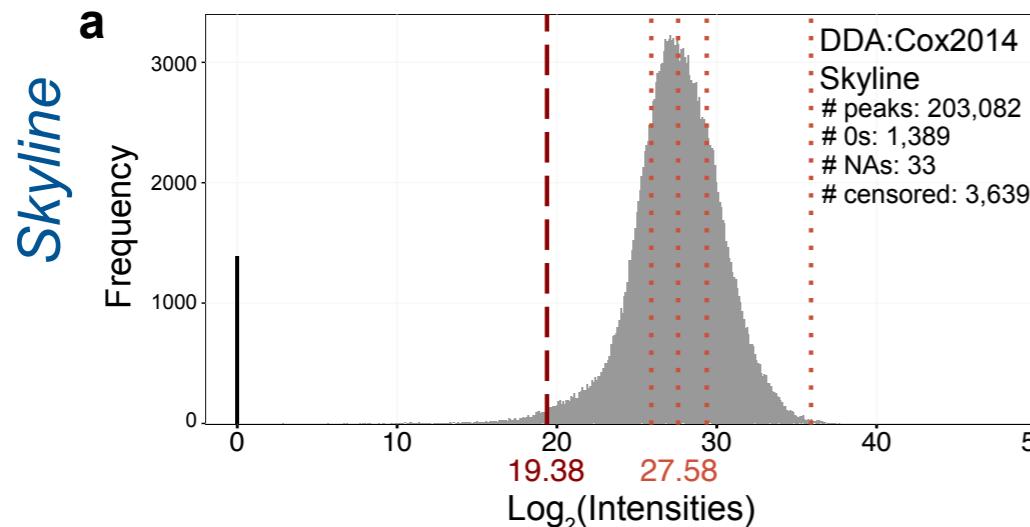
$$\text{Subject(Condition)}_{j(i)} \stackrel{\text{iid}}{\sim} \mathcal{N}(\mathbf{0}, \sigma_{\text{Subject}}^2)$$

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

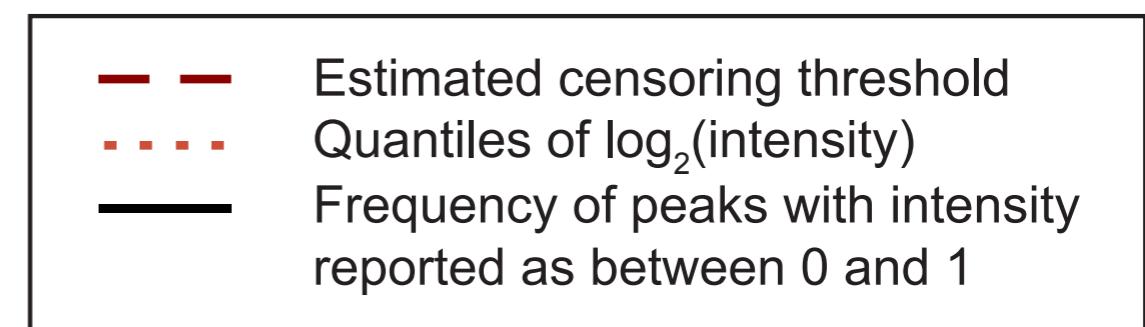
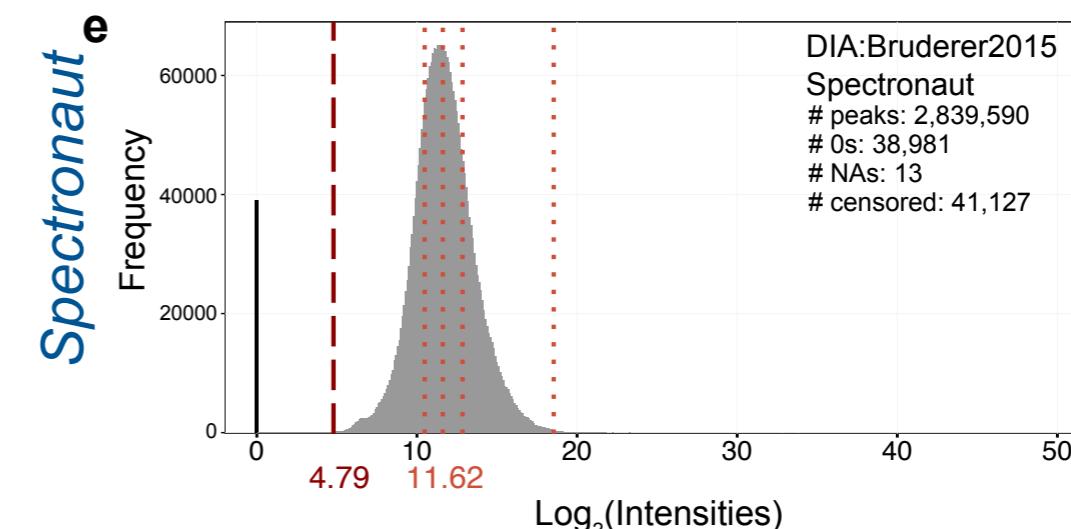
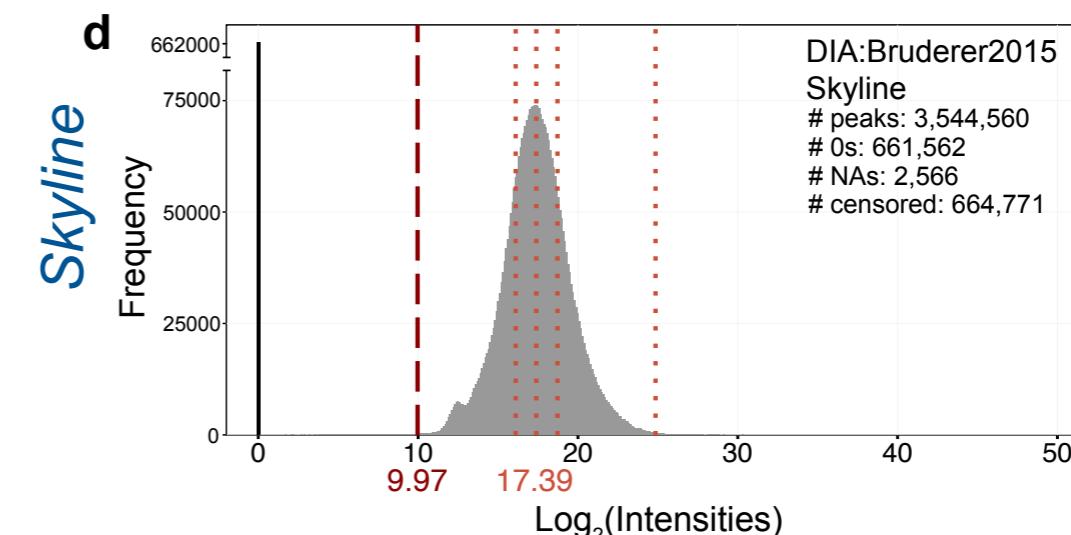
$$\text{Run} \times \text{Feature}_{ijkl} = \epsilon_{ijkl} \stackrel{\text{iid}}{\sim} \mathcal{N}(\mathbf{0}, \sigma_{\epsilon}^2)$$

PROPERTIES OF PEAK INTENSITIES VARY BETWEEN DATA PROCESSING TOOLS

DDA: Cox 2014



DIA: Bruderer 2015



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	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	



B 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	
...
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



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}



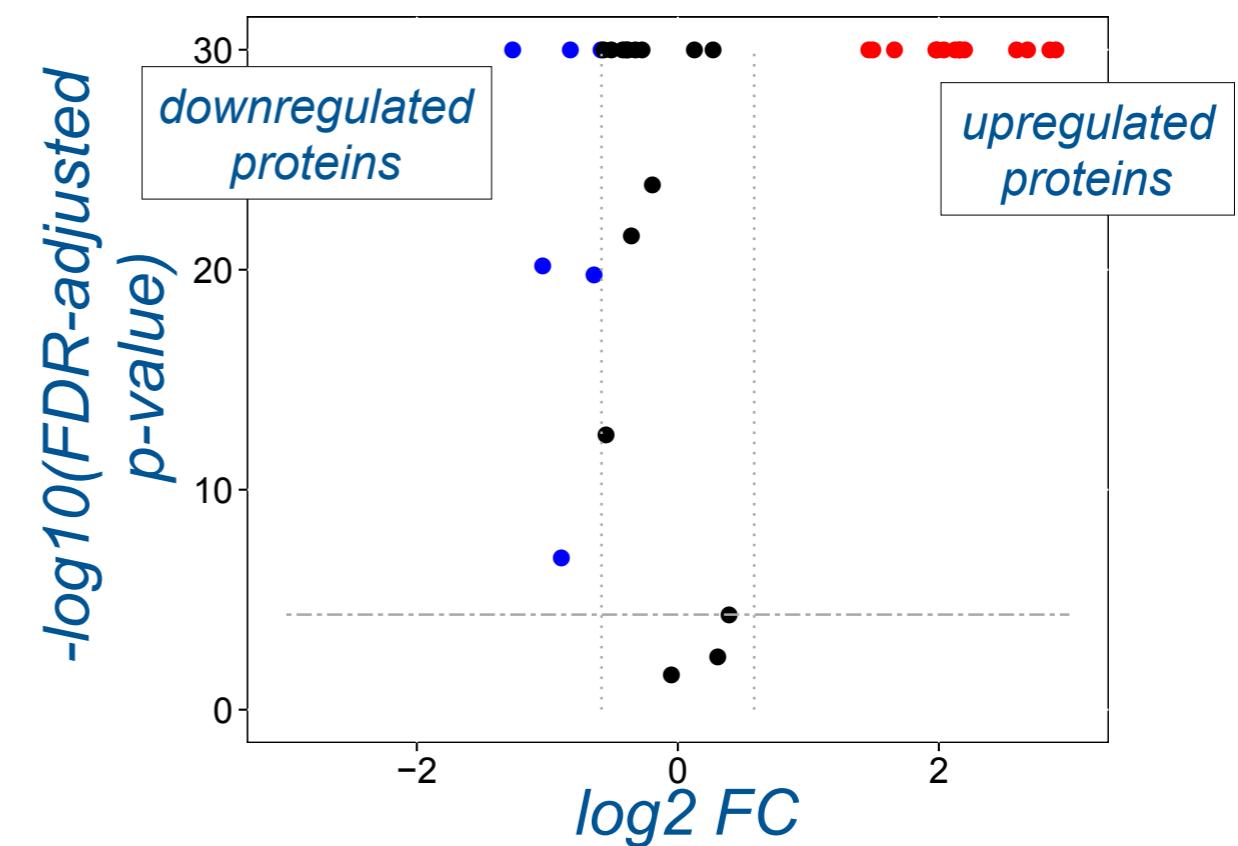
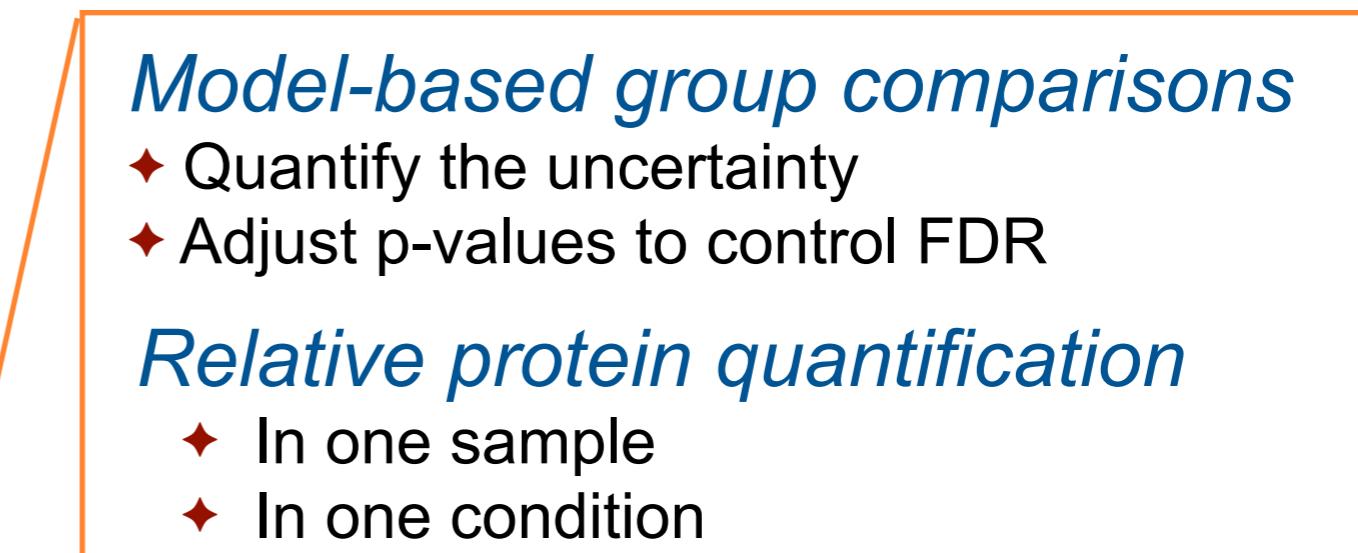
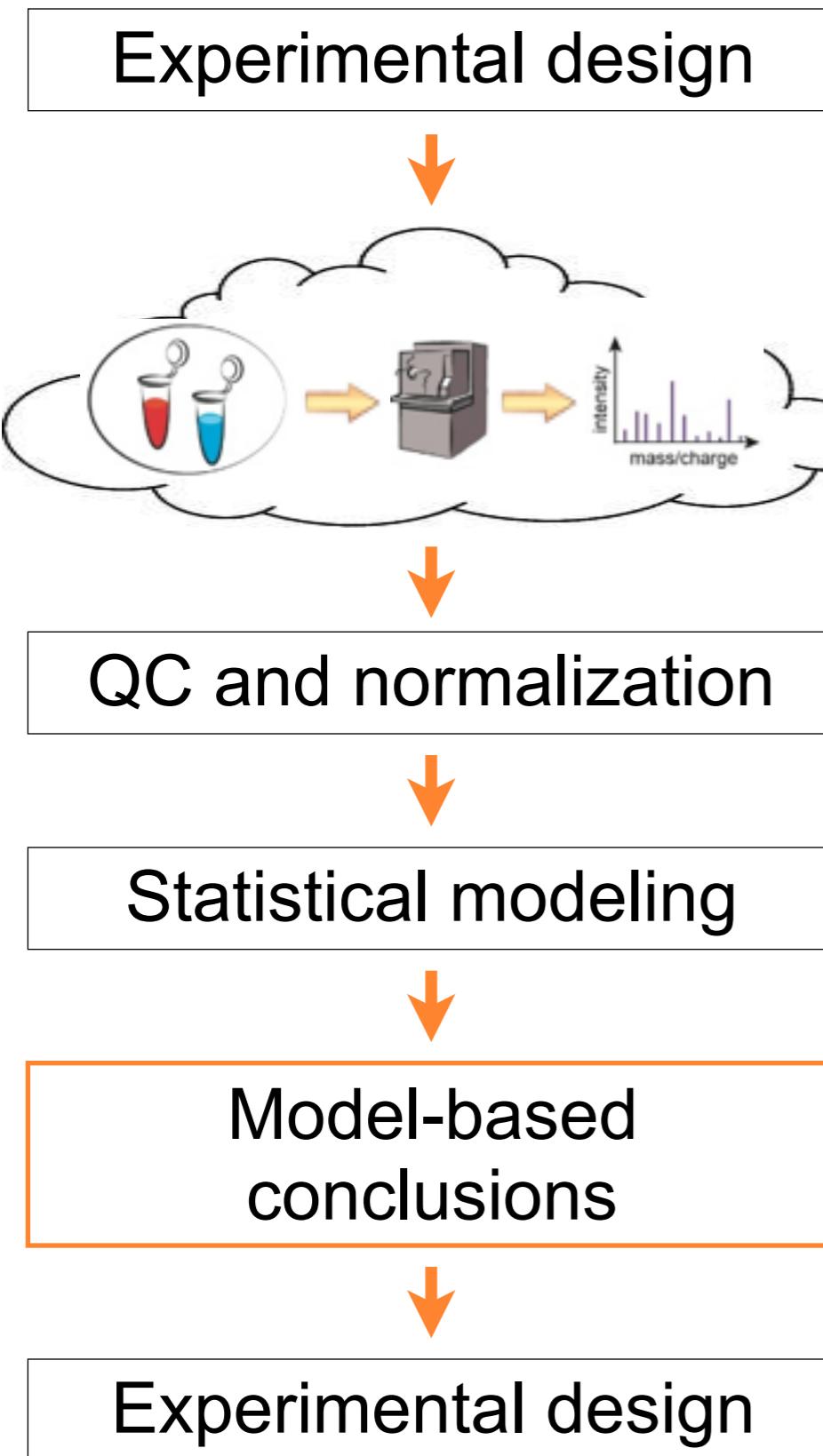
C Step 2 : Model-based inference by whole plot

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

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

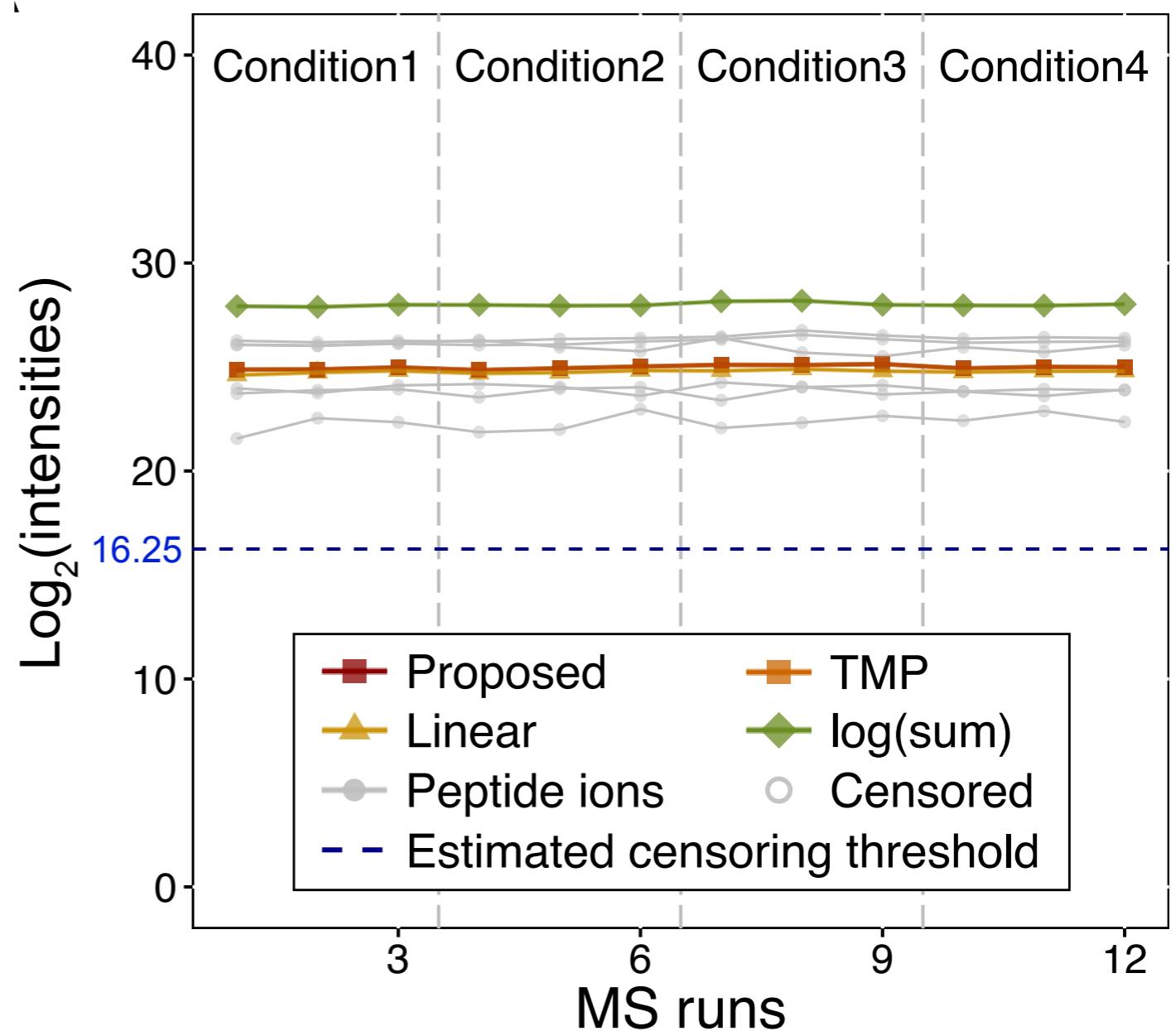
Condition ₁									...	Condition _I											
Subject ₁			Subject ₂			...	Subject _J			...	Subject _{(I-1)J+1}			Subject _{(I-1)+2}			...	Subject _{IJ}			
Run ₁	Run ₂	Run ₃	Run ₄	Run ₅	Run ₆	...	Run _{JK-2}	Run _{JK-1}	Run _{JK}	...	Run _{(I-1)JK+1}	Run _{(I-1)JK+2}	Run _{(I-1)JK+3}	Run _{(I-1)JK+4}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

A TYPICAL ANALYSIS WORKFLOW



ROBUSTNESS TO OUTLIERS

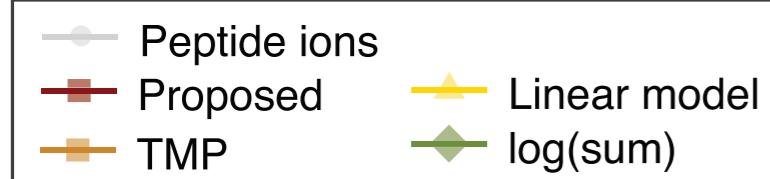
Methods perform similarly with high quality data



Condition2-Condition1 : True fold change=1

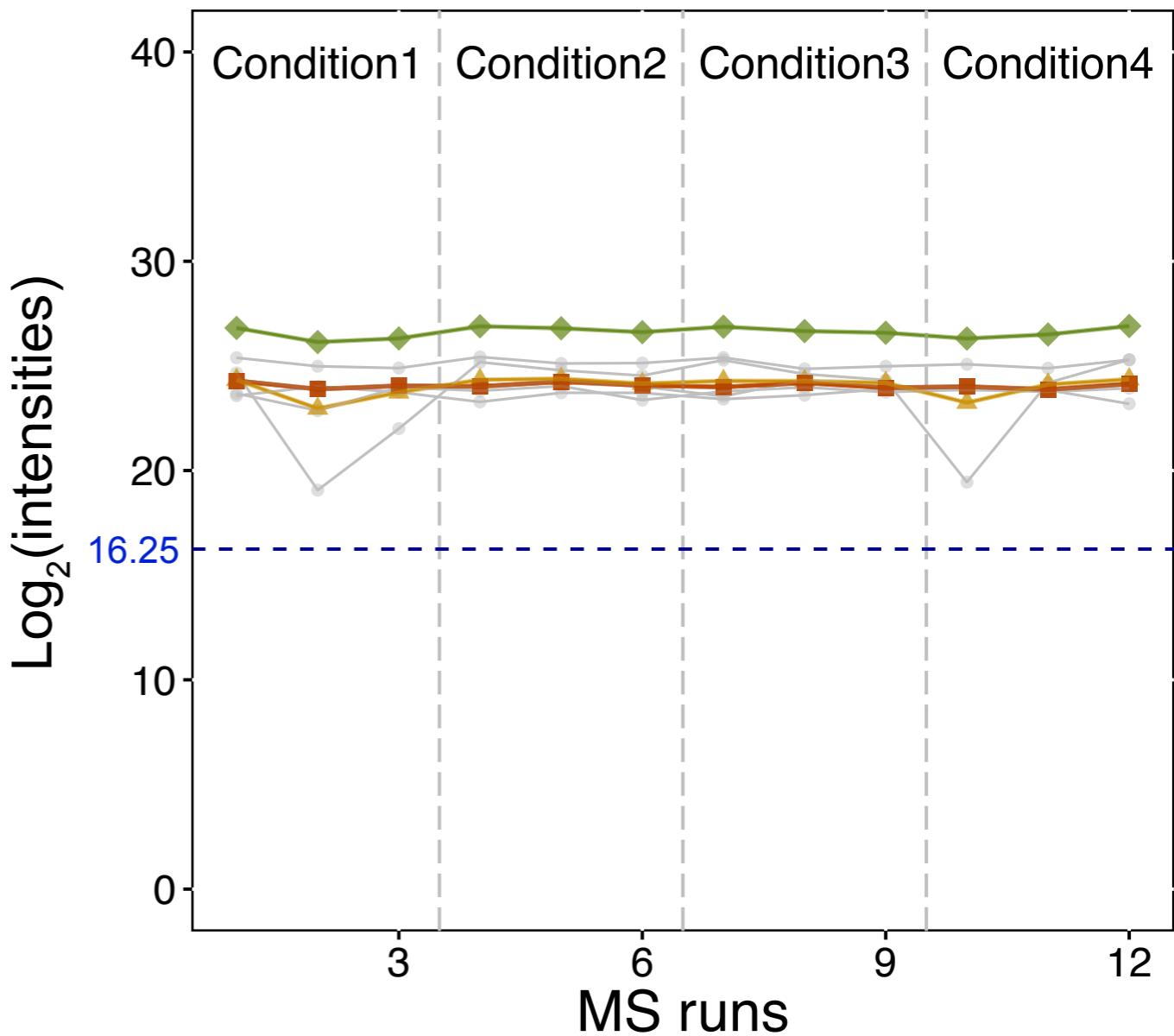
EstimatedFC Adj.pvalue

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

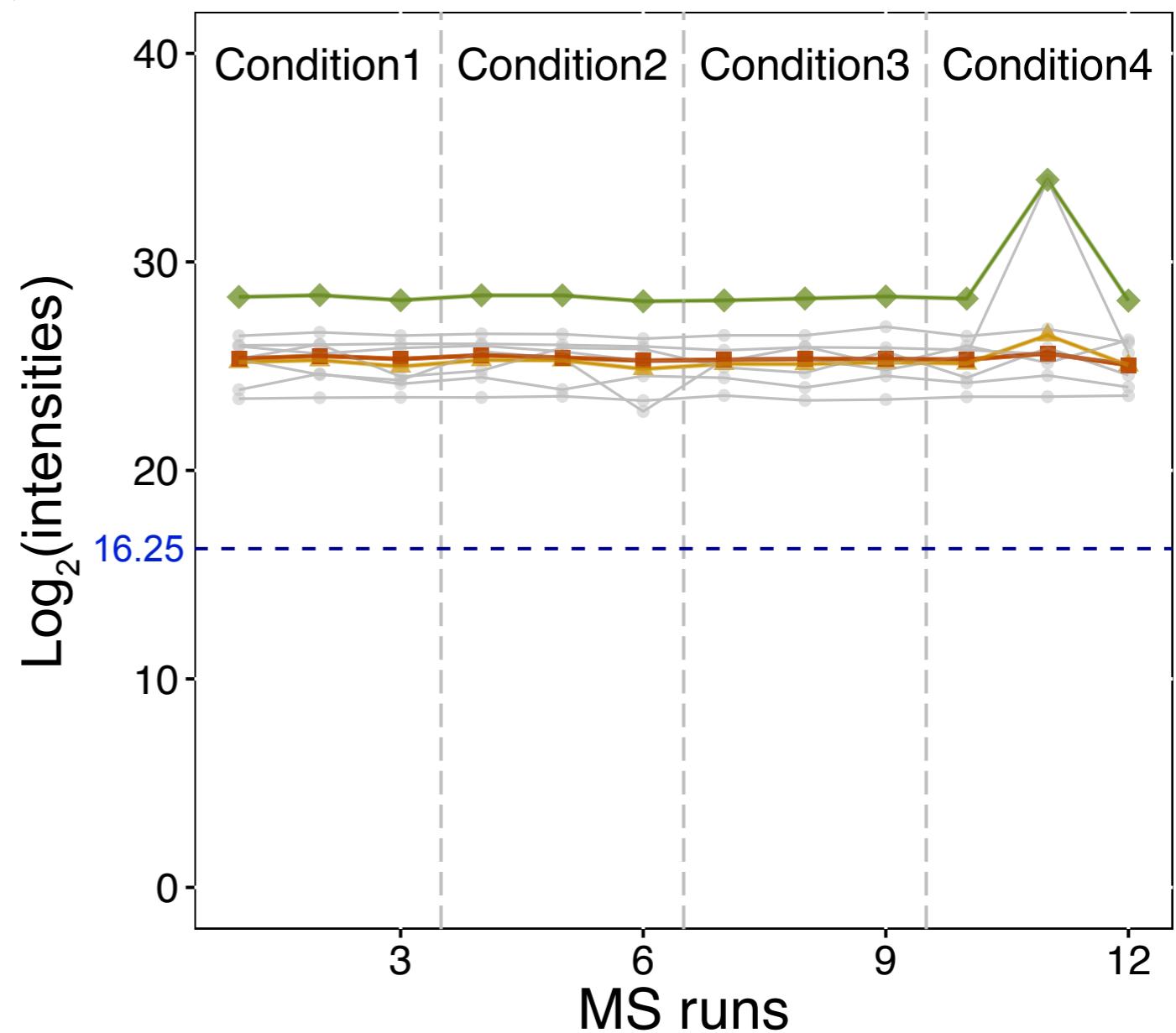
EstimatedFC Adj.pvalue

Peptide ions	
Proposed	Linear model
TMP	log(sum)

	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)*



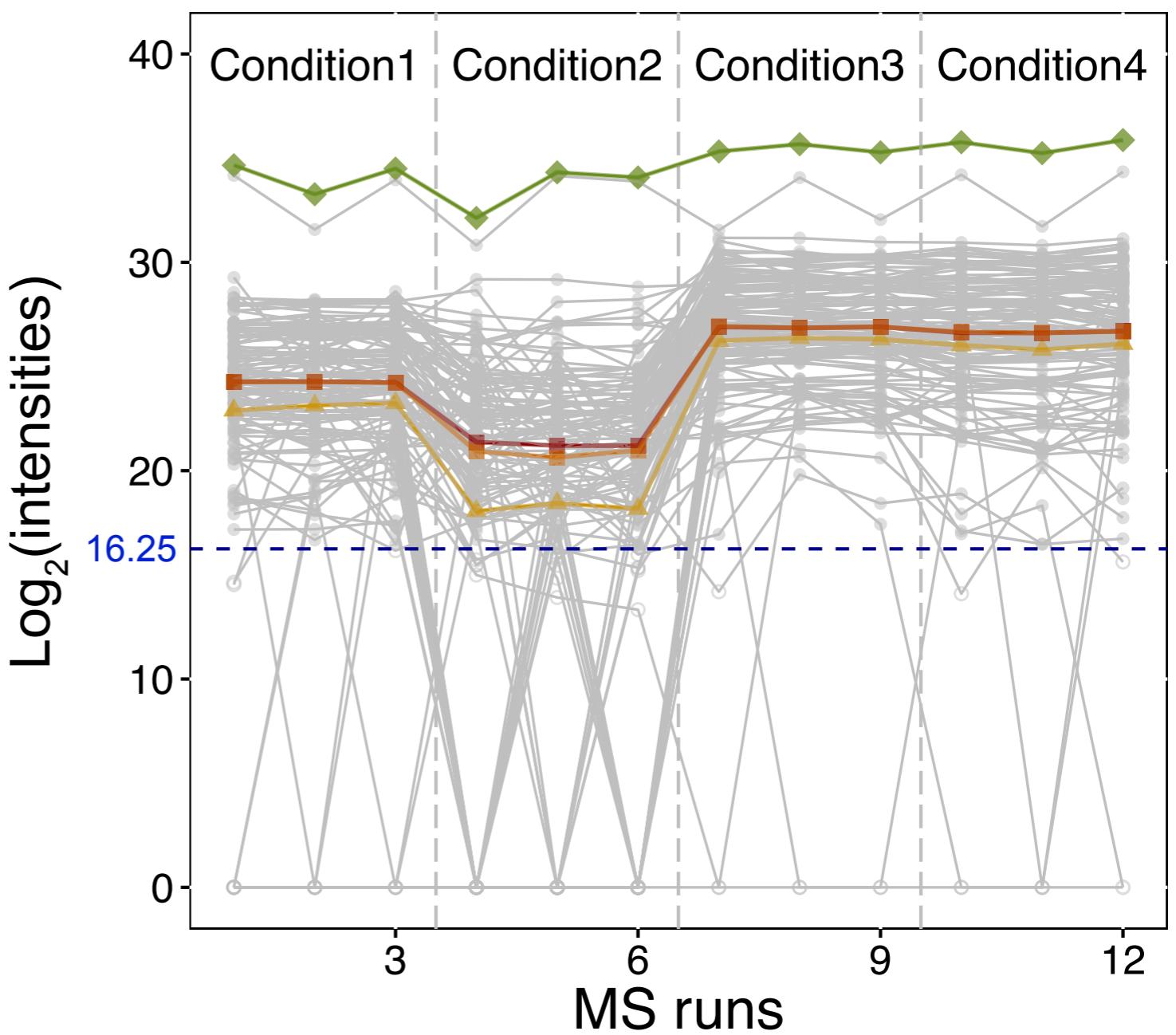
Condition4-Condition1 : True fold change=1
EstimatedFC Adj.pvalue

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

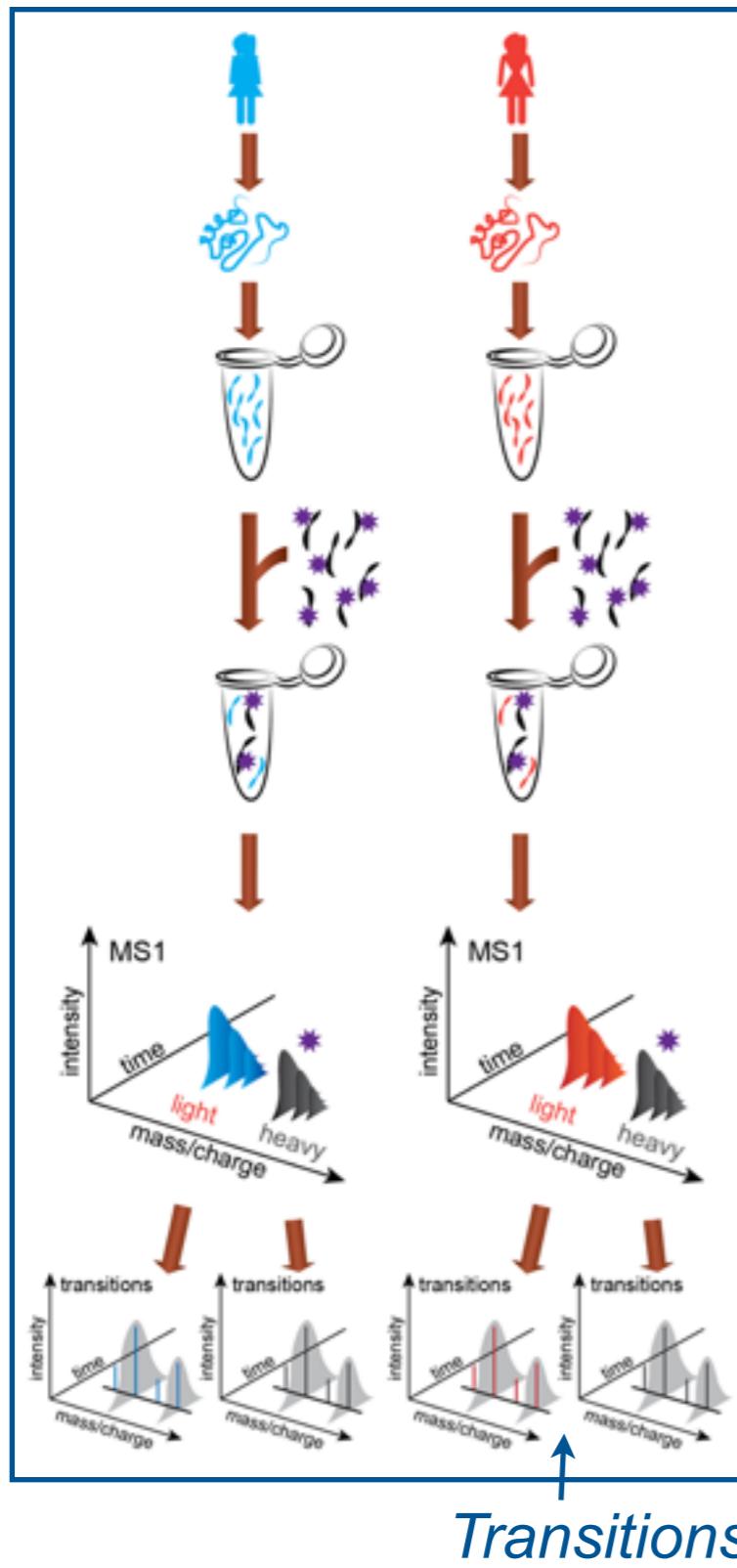
Peptide ions	
Proposed	
TMP	

EstimatedFC	Adj.pvalue	
Proposed	8.015	< 0.001
TMP	10.605	< 0.001
Linear model	29.106	< 0.001
log(sum)	1.552	0.999

EXAMPLE: LABELED REFERENCE PEPTIDES

Labeling + statistical modeling reduces the variation

Label-based SRM workflow



Analysis of heavy/light peak pairs

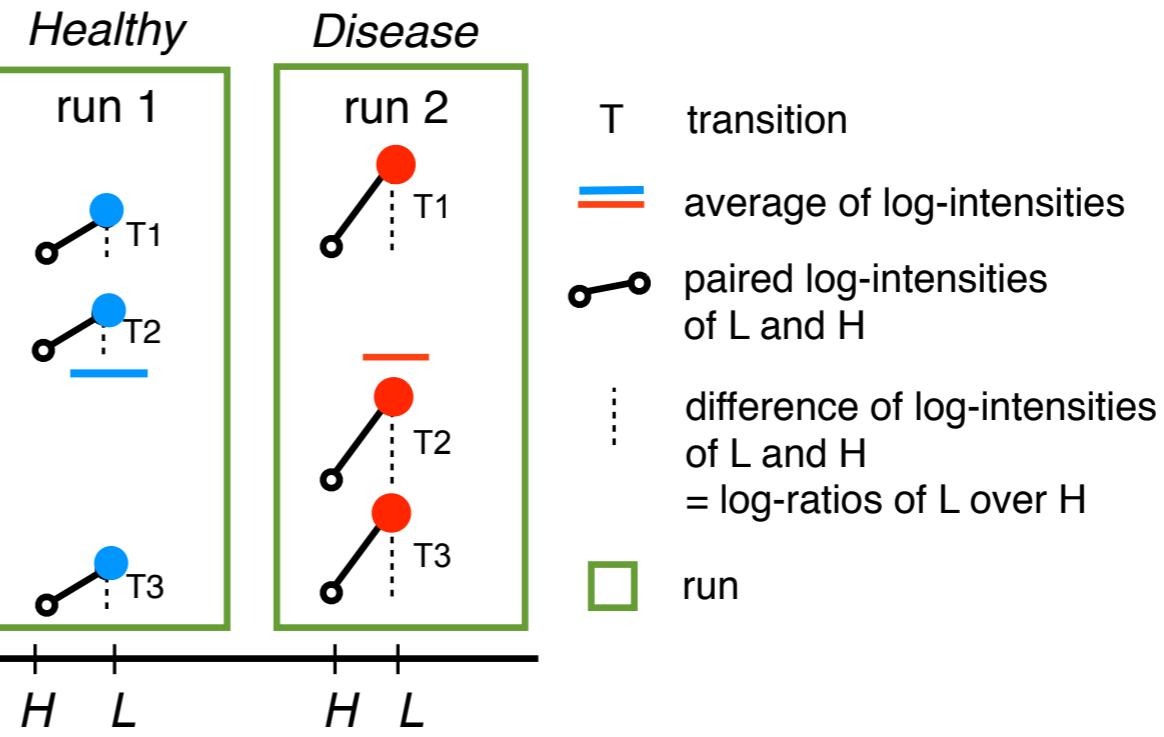


Table of quantified peaks

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

Legend :

- Label: orange oval
- Feature: Transition/Peptide: pink oval
- Group: blue oval
- Run: green oval
- Subject: brown oval

EXTENSION: LABELED REFERENCE PEPTIDES

Whole plot

Subplot		Condition _i									Condition _j										
		Subject ₁			Subject ₂			...			Subject _J			...	Subject _{(I-1)J+1}			Subject _{(I-1)J+2}			
		Run	Run	Run	Run	Run	Run	...	Run	Run	Run	Run	Run	Run	Run	Run	Run	...	Run	Run	Run
Endogenous	Feature ₁	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

	Feature _L	y		y			y	...		y	y			y		y		...	y		y
Reference	Feature ₁	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

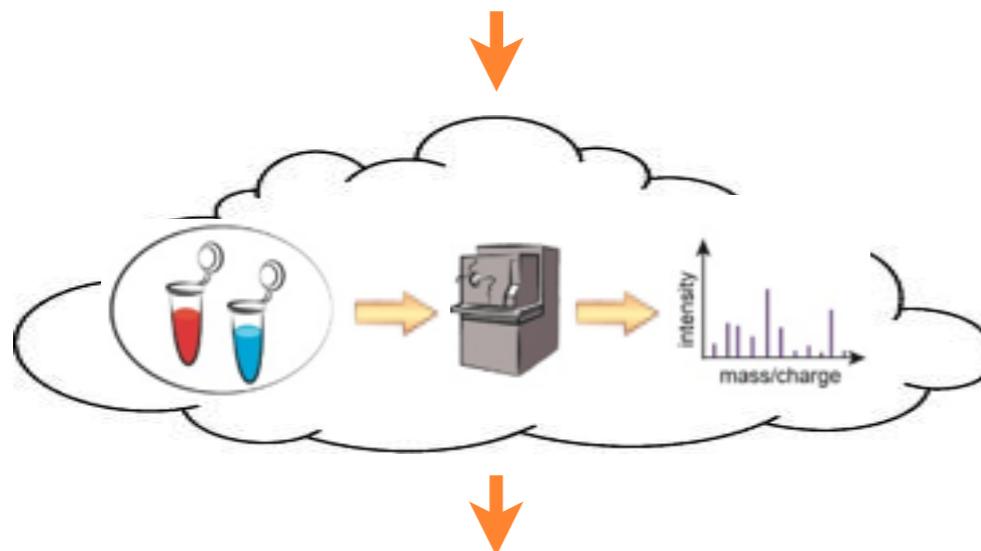
	Feature _L	y	y	y	y	y	y	...	y	y	y	...	y	y	y	y	y	...	y	y	y

Whole plot

$y_{ijklm} = \mu + \text{Condition}_i + \text{Subject}(\text{Condition})_{j(i)} + \frac{\text{Run}_{ijk}}{\text{Whole-plot biological variation}} + \frac{\text{Label}_m + \text{Run} \times \text{Label}_{ijkm}}{\text{Whole-plot technical variation}} + \text{Feature}_l + \frac{\epsilon_{ijklm}}{\text{Subplot error}}$
where $\sum_{i=1}^I \text{Condition}_i = 0, \sum_{j=1}^L \text{Feature}_l = 0, \sum_{m=0}^1 \text{Label}_m = 0$
$\text{Subject}(\text{Condition})_{j(i)} \stackrel{\text{iid}}{\sim} \mathcal{N}(\mathbf{0}, \sigma_{\text{Subject}}^2)$
$\text{Run}_{ijk} = \psi_{ijk} \stackrel{\text{iid}}{\sim} \mathcal{N}(\mathbf{0}, \sigma_{\psi}^2)$
$\epsilon_{ijklm} \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \sigma_{\epsilon}^2)$

A TYPICAL ANALYSIS WORKFLOW

Experimental design



QC and normalization



Statistical modeling



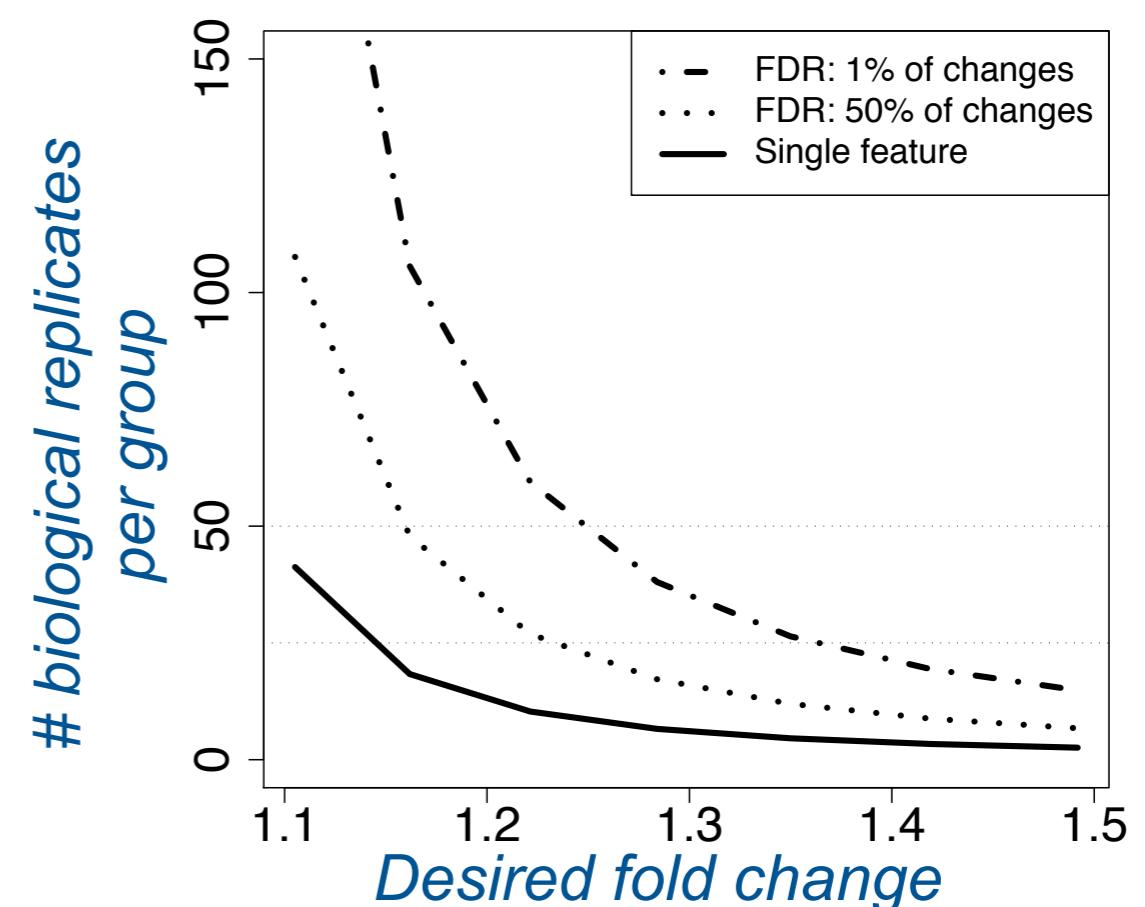
Model-based conclusions



Experimental design

Use the dataset to improve:

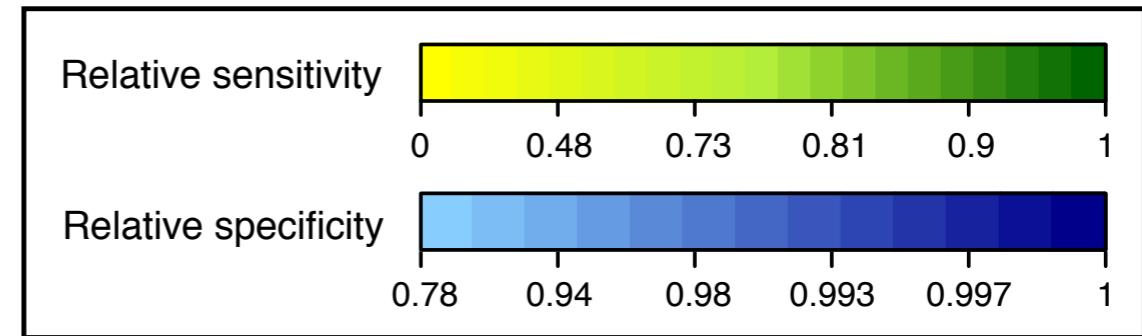
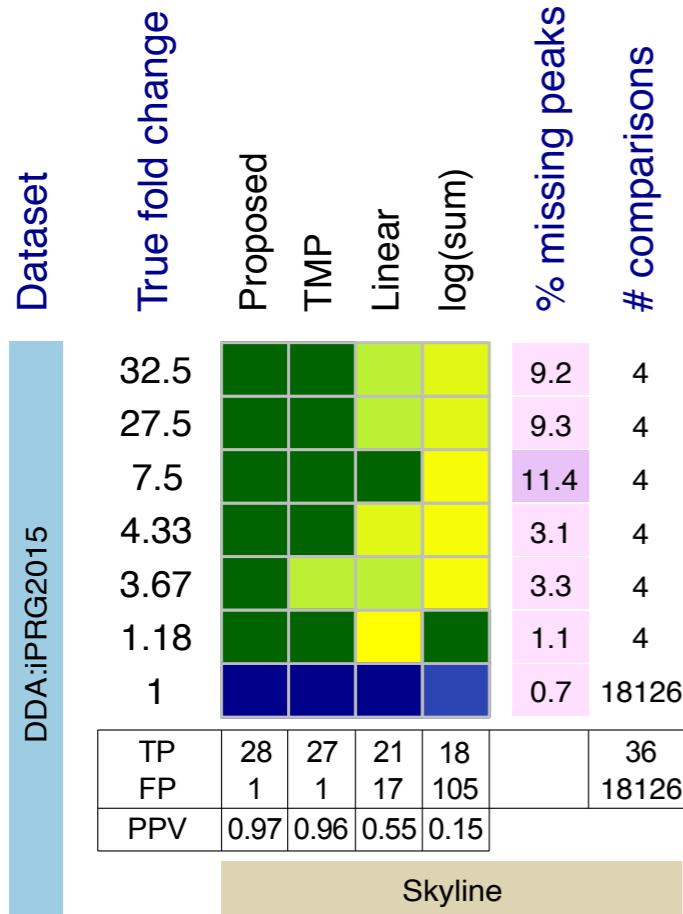
- Subject selection: matching
- Resource allocation: blocking
- Calculation of sample size



OUTLINE

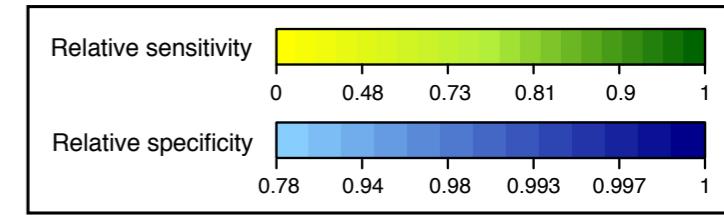
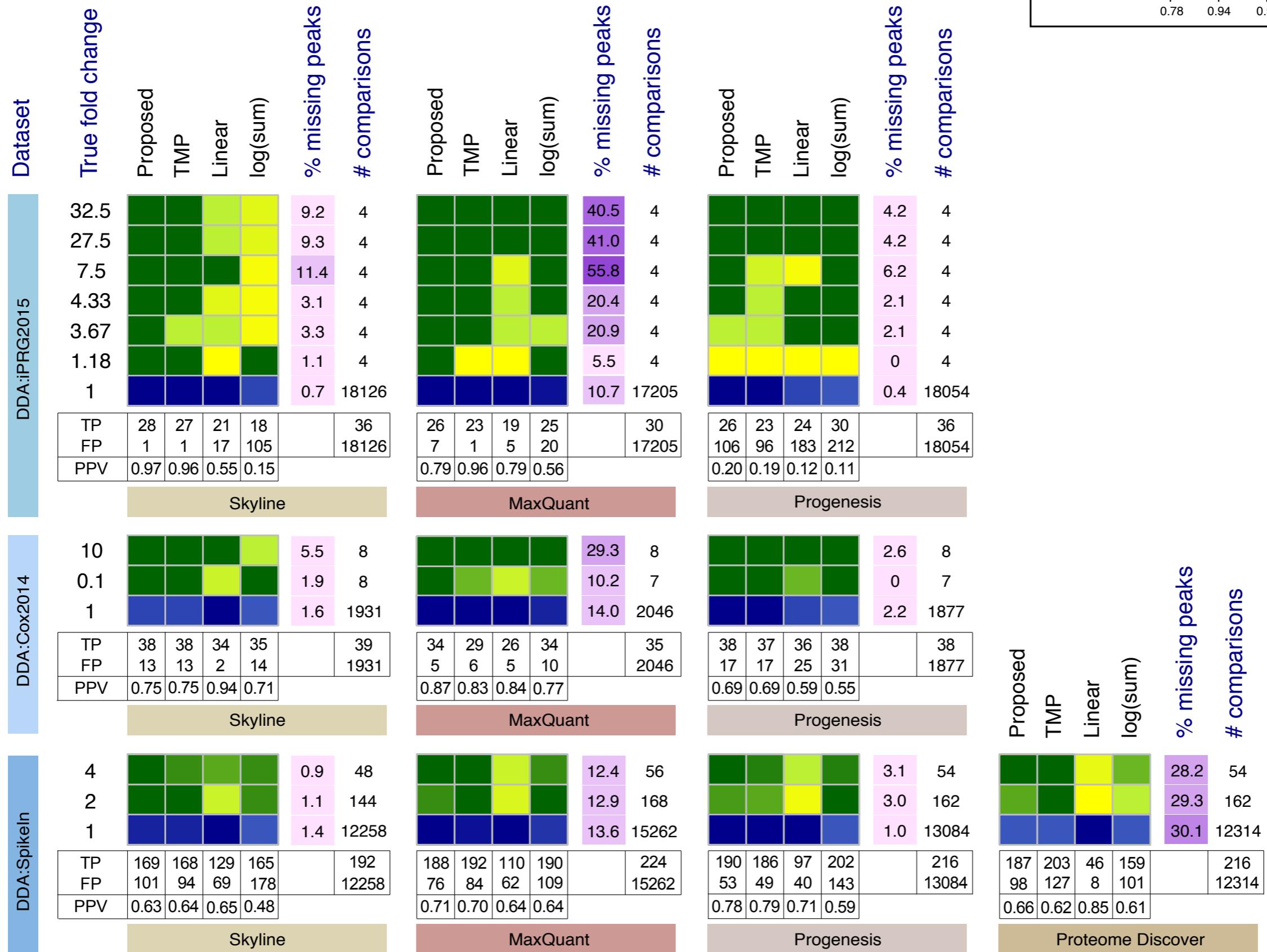
- iPRG2015: detection of diff. abundance
 - Community effort with label-free shotgun proteomics
- MSstats
 - Normalization, statistical modeling, inference
 - Evaluation
- Extensions to MSstats
 - Assay characterization, longitudinal monitoring

TESTING: DDA

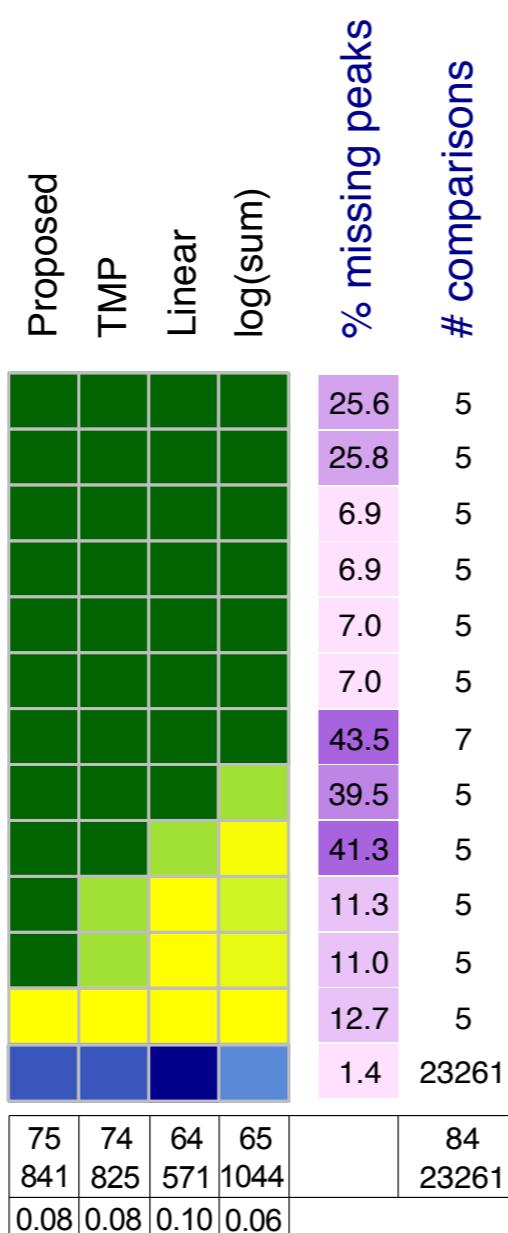
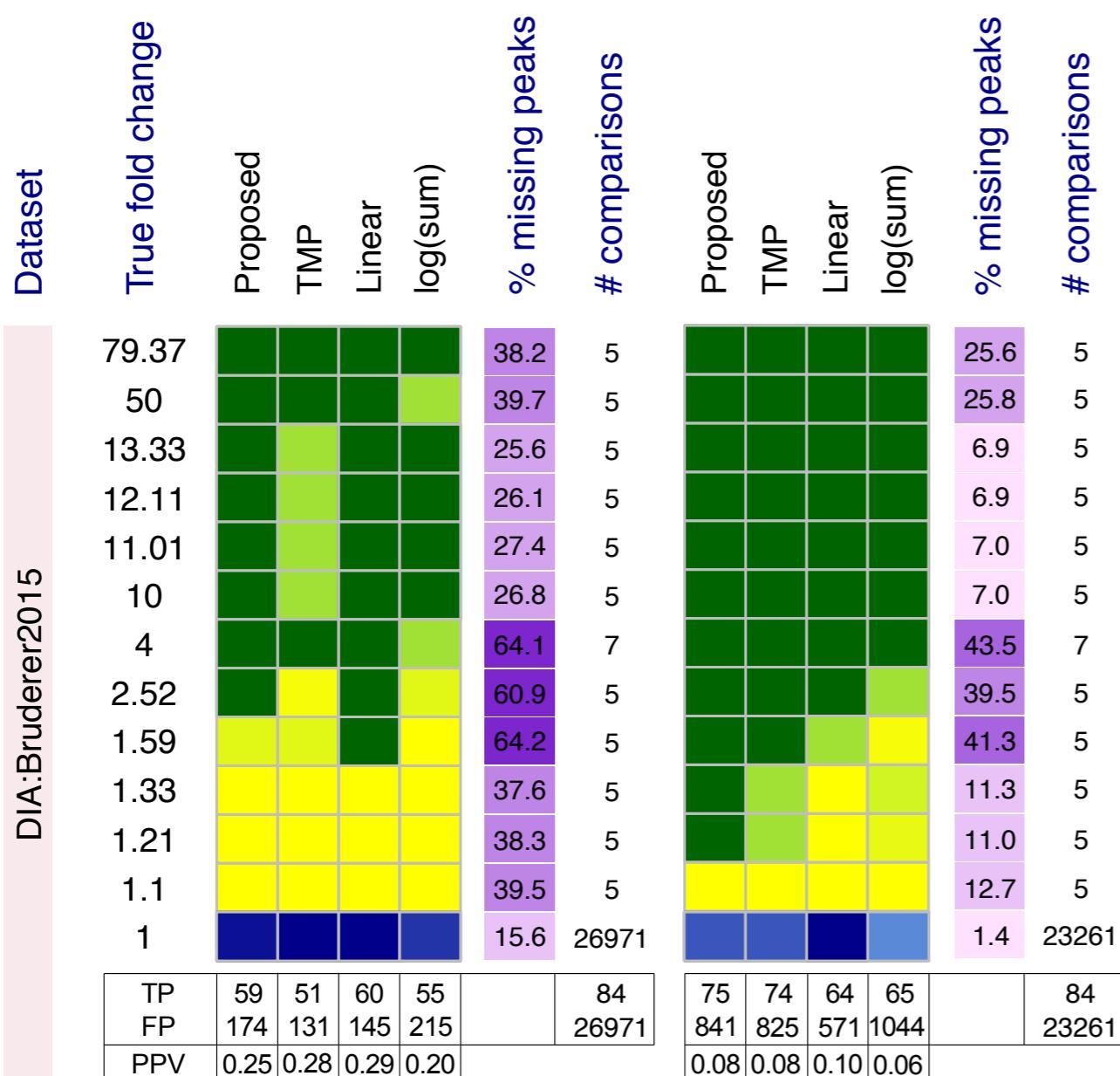


TESTING: DDA

a

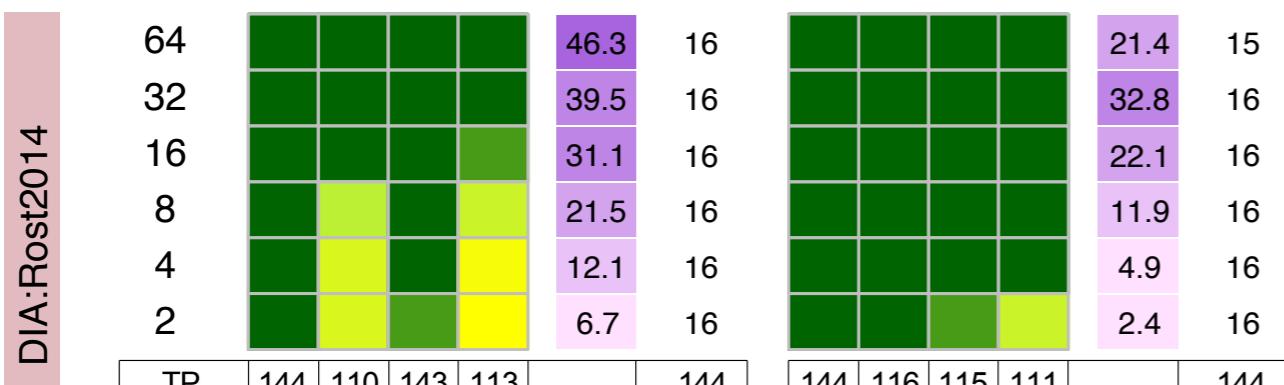


TESTING: DIA



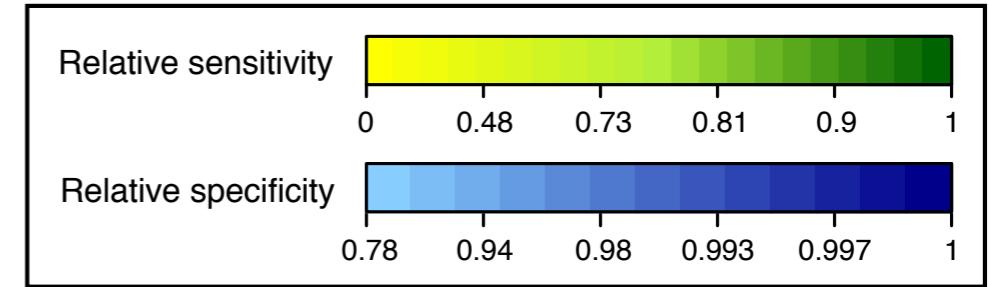
Skyline

Spectronaut

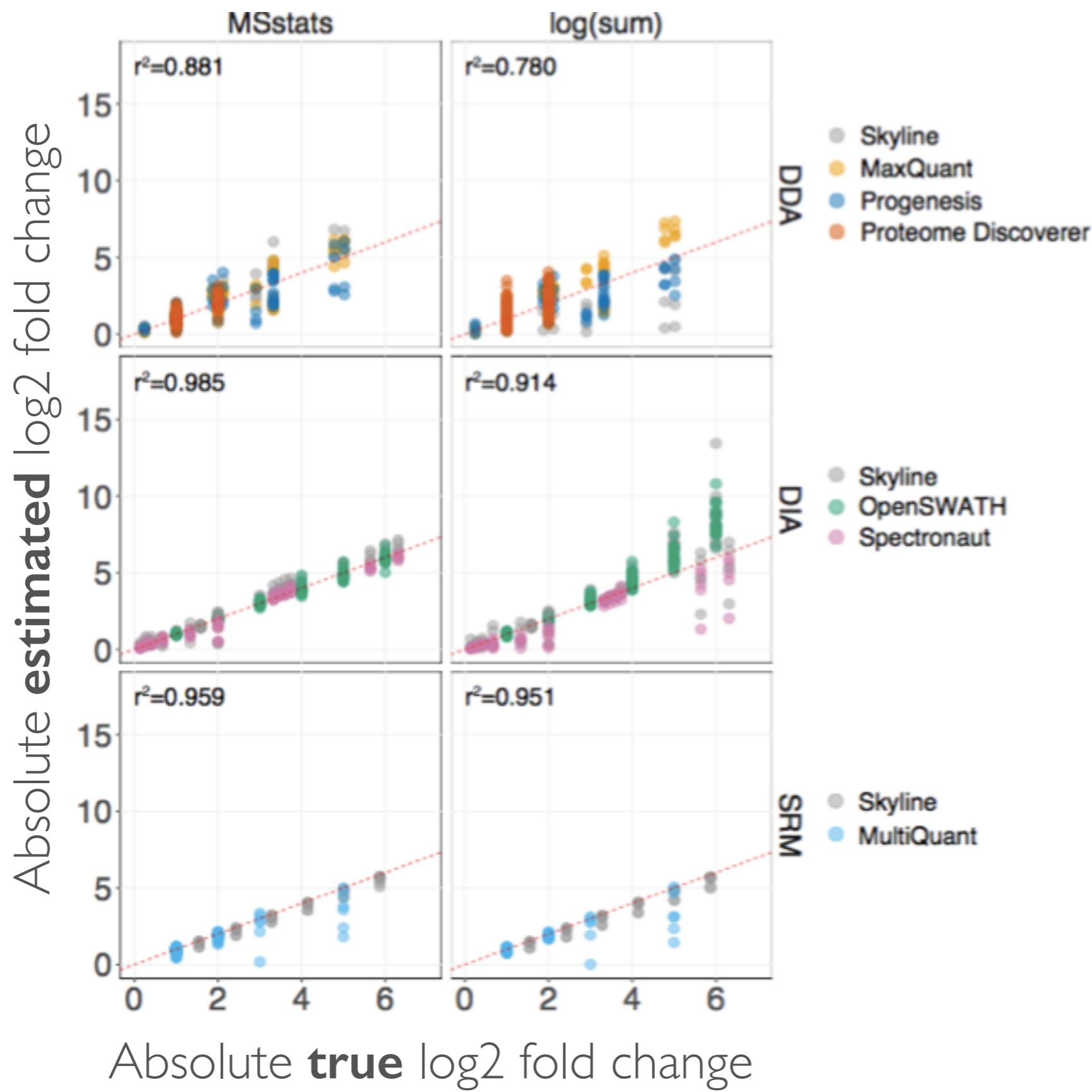


Skyline

OpenSWATH



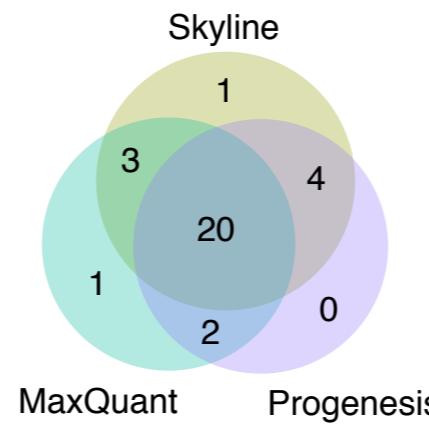
ESTIMATION OF LOG-FOLD CHANGE



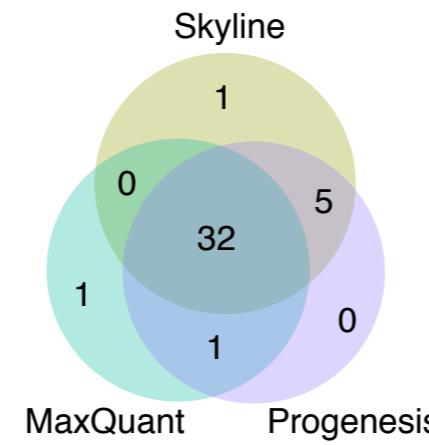
BETTER STATISTICAL METHODS ENHANCE REPRODUCIBLE RESEARCH

Better agreement
in #differentially
abundant proteins
between tools

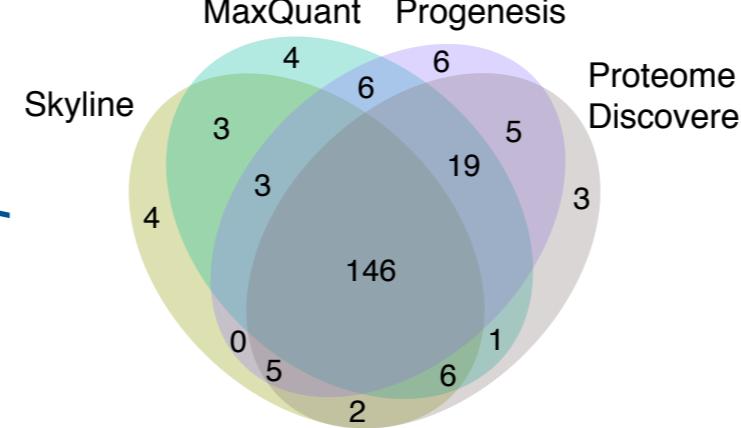
DDA: iPRG2015



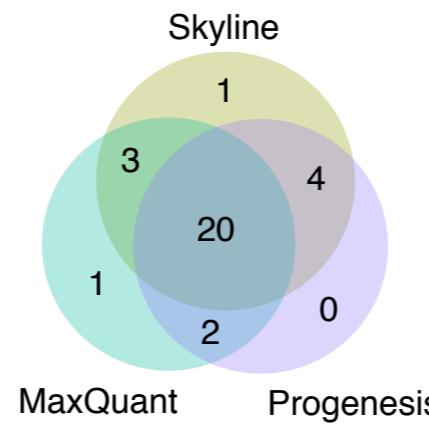
DDA: Cox 2014



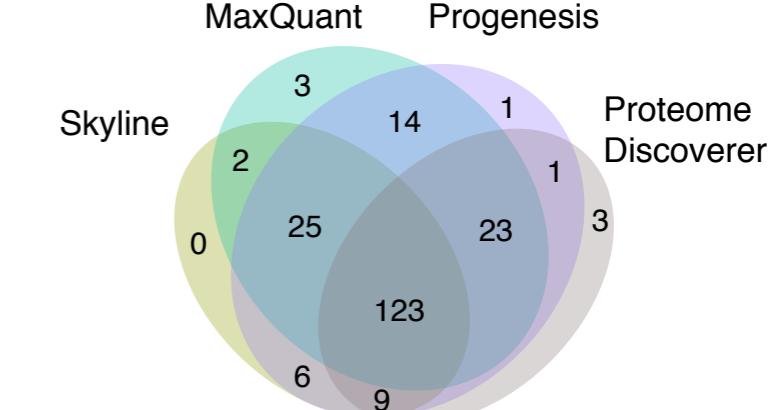
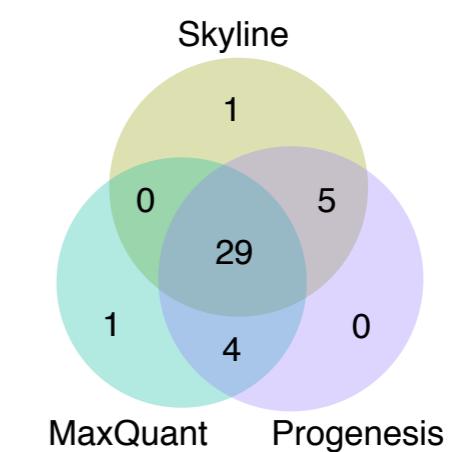
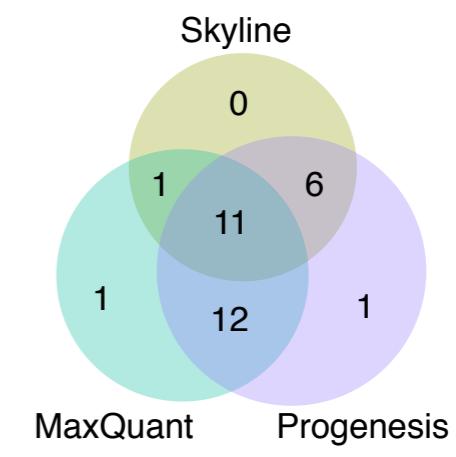
DDA: Spike-in



Proposed

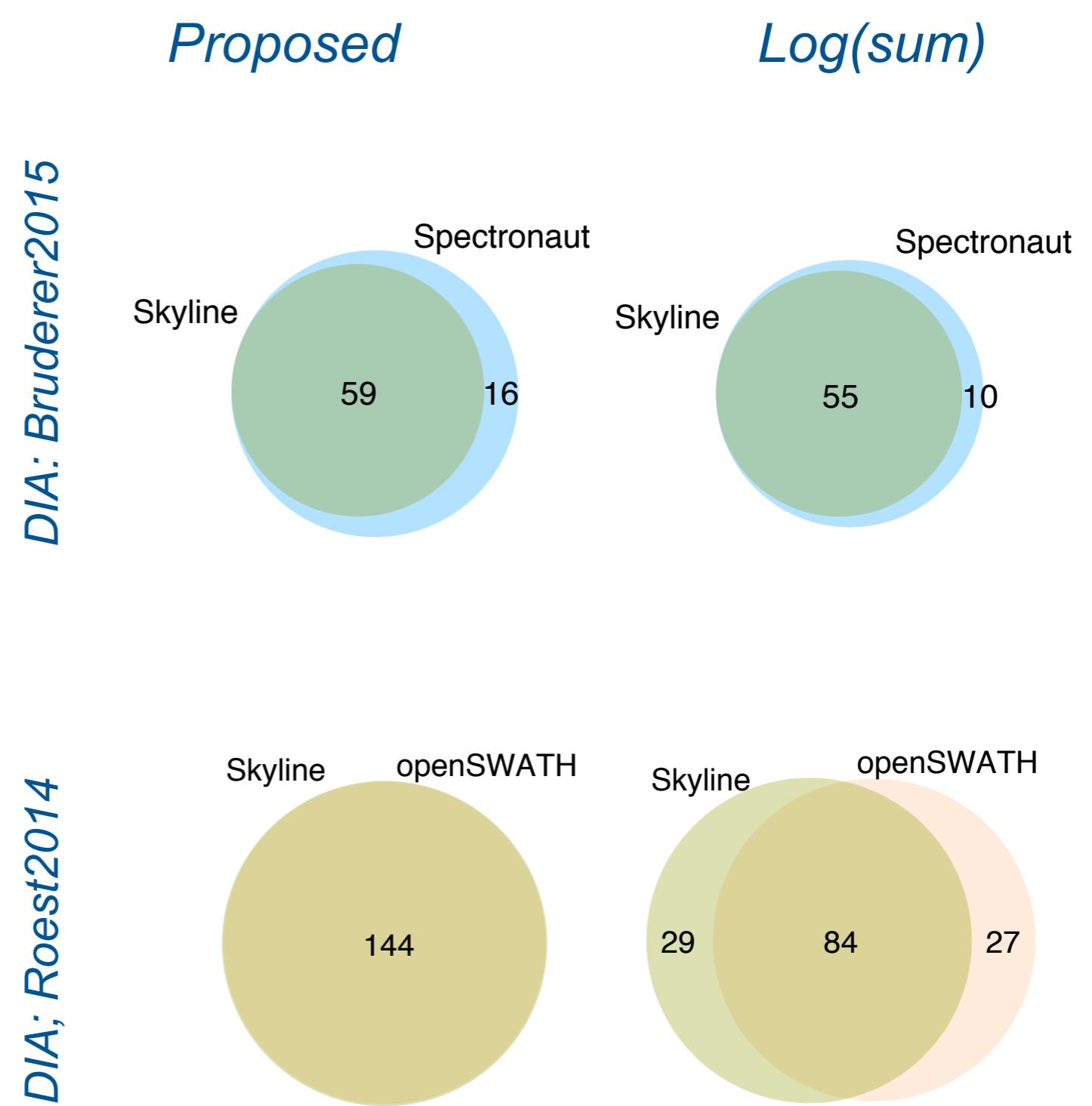


Log(sum)

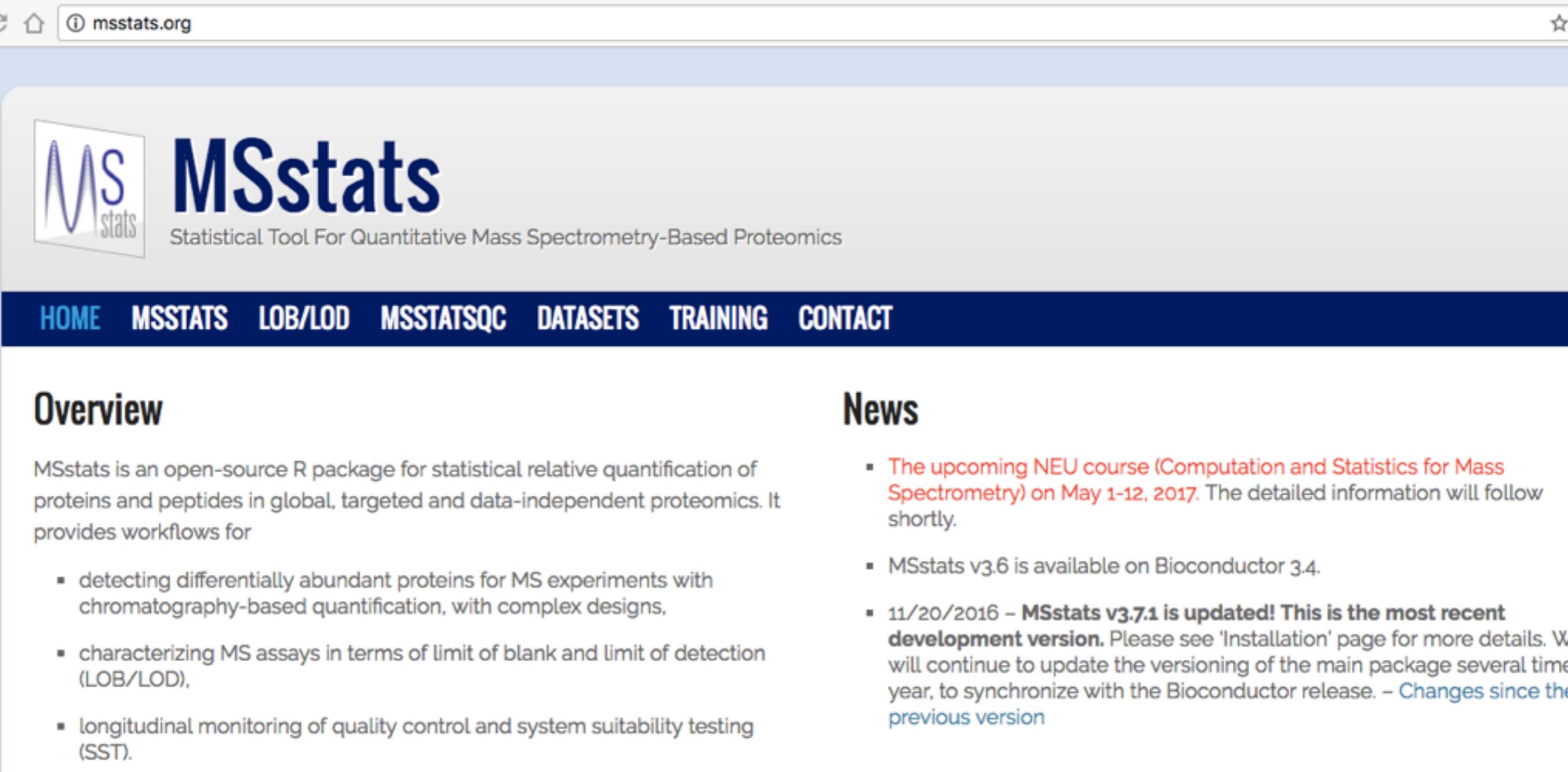


BETTER STATISTICAL METHODS ENHANCE REPRODUCIBLE RESEARCH

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 user account, followed by the URL 'msstats.org'. Below the header is the MSstats logo, which consists of a stylized 'M' and 'S' icon next to the word 'MSstats'. The tagline 'Statistical Tool For Quantitative Mass Spectrometry-Based Proteomics' is displayed. A navigation menu below the logo includes links for HOME, MSSTATS, LOB/LOD, MSSTATSQC, DATASETS, TRAINING, and CONTACT. The main content area has two sections: 'Overview' on the left and 'News' on the right. The 'Overview' section contains text about the package's purpose and features, along with a bulleted list of its capabilities. The 'News' section lists recent developments, including the upcoming NEU course, the availability of MSstats v3.6, and the update to v3.7.1.

MSstats

Statistical Tool For Quantitative Mass Spectrometry-Based Proteomics

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



MSstats

Statistical Tool For Quantitative Mass Spectrometry-Based Proteomics

www.msstats.org

M. Choi et al. *BMC Bioinformatics.*, 2014

Skyline external tool

[Download MSstats](#)

Downloaded: 10036



Tool Information

Organization: Vitek Lab, Purdue University

Authors: Meena Choi, Cyril Galitzine, Tsung-Heng Tsai, Olga Vitek

Languages: R(3.3.1), C#

More Information: <http://www.msstats.org/>

Bioconductor

Month	Nb of distinct IPs	Nb of downloads
Jan/2016	160	259
Feb/2016	160	295
Mar/2016	222	390
Apr/2016	255	366
May/2016	193	328
Jun/2016	220	378
Jul/2016	170	290
Aug/2016	153	241
Sep/2016	133	215
Oct/2016	172	368
Nov/2016	234	398
Dec/2016	192	363
2016	1808	3891

[MSstats 2016 stats.tab](#)

platforms all

downloads top 20%

posts 0

in Bioc 3.5 years

build ok

commits 0.33

test coverage unknown

OUTLINE

- iPRG2015: detection of diff. abundance
 - Community effort with label-free shotgun proteomics
- MSstats
 - Normalization, statistical modeling, inference
 - Evaluation
- Extensions to MSstats
 - Assay characterization, longitudinal monitoring

STATISTICAL METHODS FOR ASSAY CHARACTERIZATION

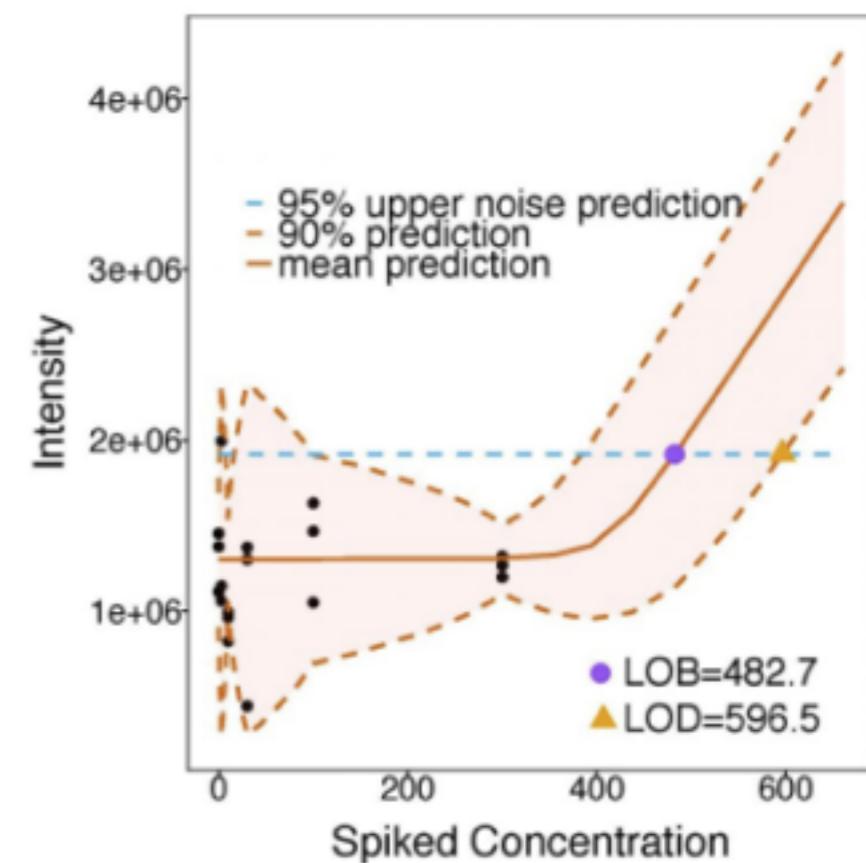


[HOME](#) [MSSTATS](#) [LOB/LOD](#) [MSSTATSQC](#) [DATASETS](#) [TRAINING](#) [CONTACT](#)

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

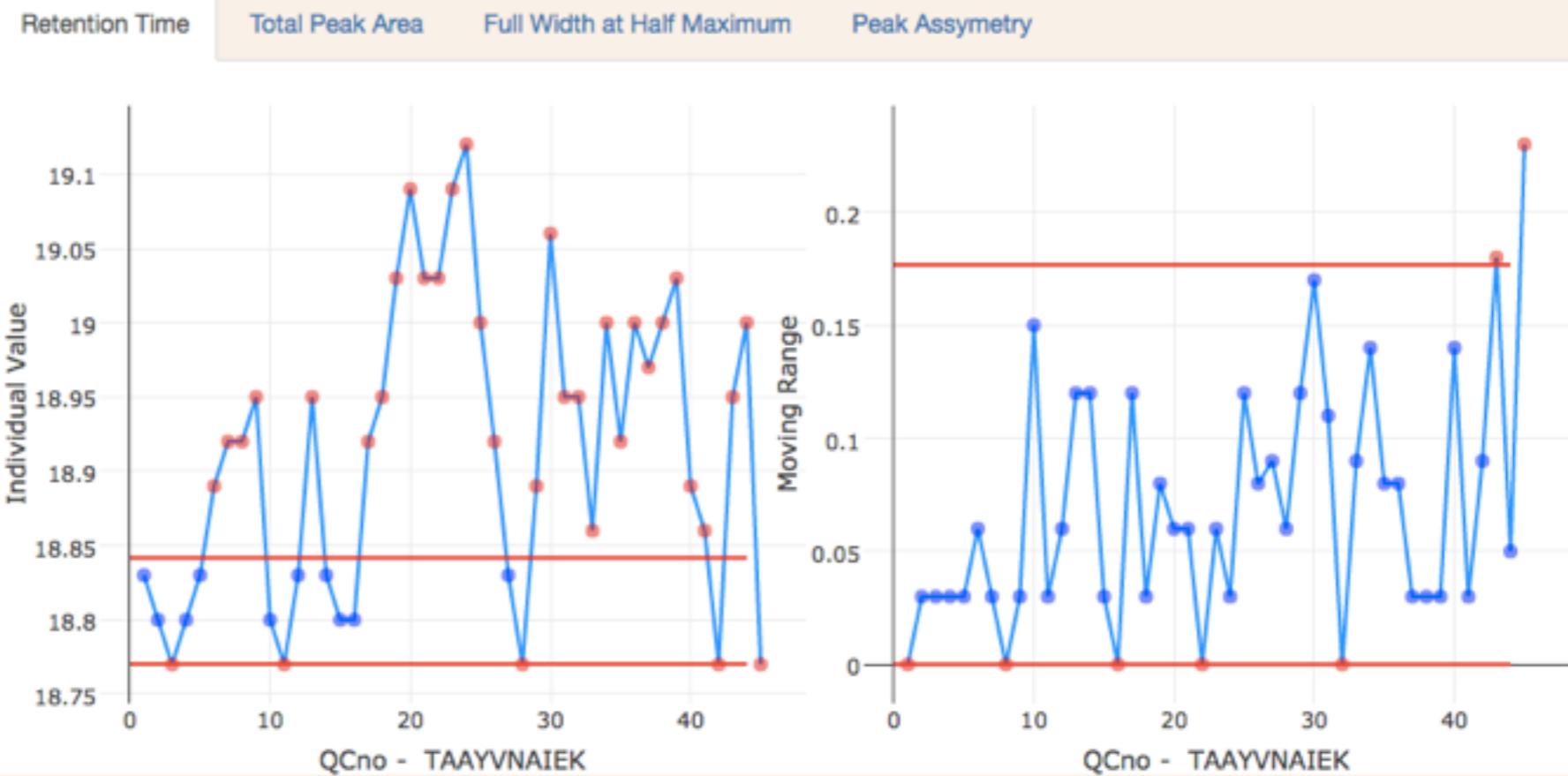
MSstatsQC

System suitability monitoring tools for quantitative mass spectrometry based proteomic experiments

Home Data Import Metric Summary Control Charts ▾ Change Point Analysis Help

choose your preferred metric to view plots

- Retention Time
- Total Peak Area
- Full Width at Half Maximum
- Peak Assymetry



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,

