

Bayesian statistical modeling reveals missing value mechanisms in label-free Mass Spectrometry-based proteomics experiments

Devon Kohler and Olga Vitek

Khoury College of Computer Sciences
Northeastern University, Boston, MA

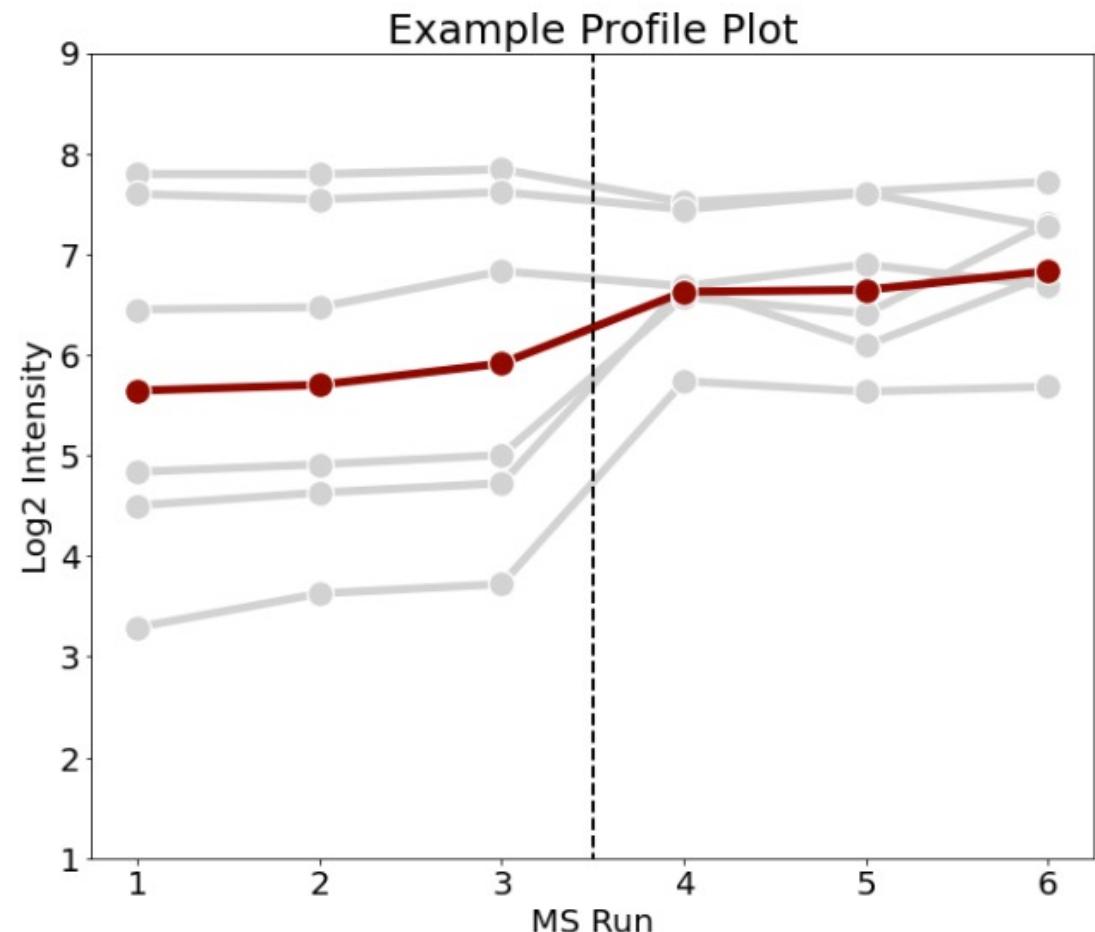
The authors declare no conflicts of interest

Presentation Outline

- Problem statement
- Background
 - MSstats model
 - Missing value imputation
 - Summarization
- Proposed method
- Evaluation
 - Missing mechanisms
 - Imputation and summarization
 - Differential analysis

Accurate relative quantification of proteins in MS-based proteomics is challenging and requires specialized techniques

- Challenges
 - Indirect measurements (peptide fragments)
 - Missing peptides across samples
- Analytical/statistical techniques
 - Normalization
 - Imputation of missing peptide intensities
 - Summarization
- Feature definition
 - precursor ions in DDA
 - transitions in SRM
 - fragments in DIA



Presentation Outline

- Problem statement
- Background
 - MSstats model
 - Missing value imputation
 - Summarization
- Proposed method
- Evaluation
 - Missing mechanisms
 - Imputation and summarization
 - Differential analysis

Experimental design and MSstats model for one protein

Whole plot											
Subplot	Condition ₁					...	Condition _I				
	BioReplicate ₁	BioReplicate ₂	...	BioReplicate _J	...		BioReplicate _{(I-1)J+1}	BioReplicate _{(I-1)J+2}	...	BioReplicate _{IJ}	
Feature ₁	y	y	...	y	...	y	y	y	...	y	
Feature ₂	NA	y	...	y	...	y	NA	y	
...	
Feature _L	y	NA	...	y	...	y	y	y	...	y	

$$y_{ijl} = \mu + Condition_i + BioReplicate(Condition)_{j(i)} + Feature_l + BioReplicate * Feature_{ijl}$$

where $\sum_{i=1}^I Condition_i = 0, \sum_{l=1}^L Feature_l = 0$

$$BioReplicate(Condition)_{j(i)} \sim^{iid} N(0, \sigma_{BioReplicate}^2)$$

$$Run * Feature_{ijl} = \epsilon_{ijl} \sim^{iid} N(0, \sigma_\epsilon^2)$$

Kohler et al. Journal of Proteome Research. 2023.

Two-step model in MSstats

Step 1 : Run-level subplot summarization

Impute censored intensities with Accelerated Failure Time 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 _{UK-2}	Run _{UK-1}	Run _{UK}
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	\hat{y}_{imp}	y	\hat{y}_{imp}	\hat{y}_{imp}	y	...	\hat{y}_{imp}	y	y	...	NA	y	y	y	y	y	...	y	\hat{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 _{UK-2}	Run _{UK-1}	Run _{UK}
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}



Step 2 : Whole plot modeling and inference

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

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

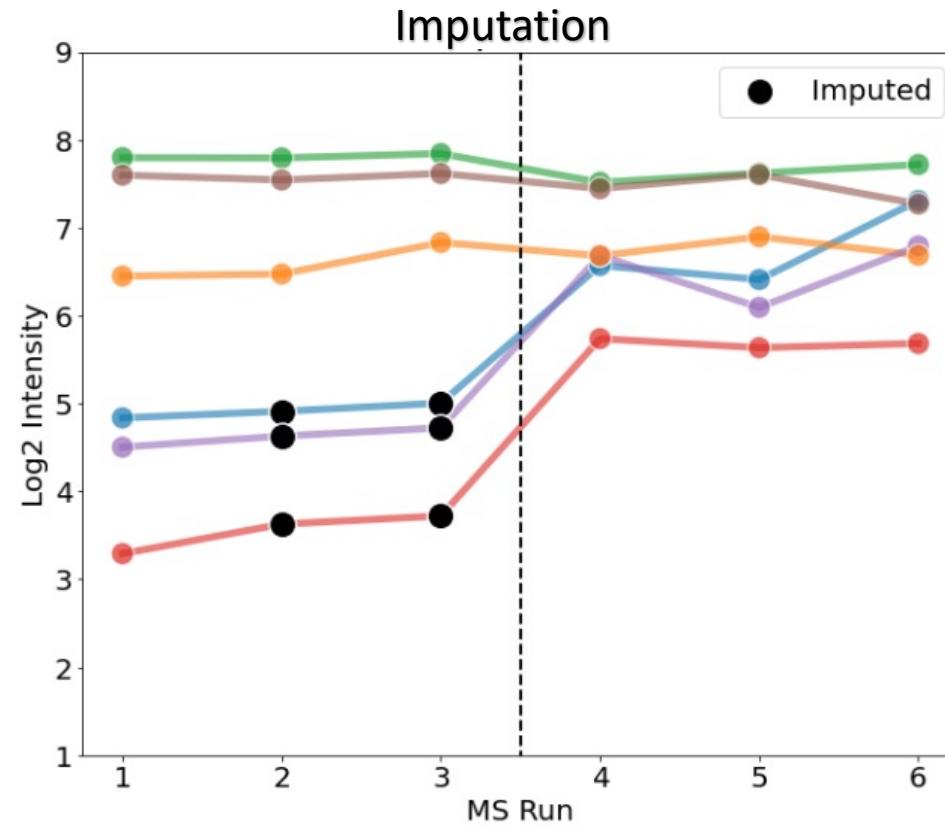
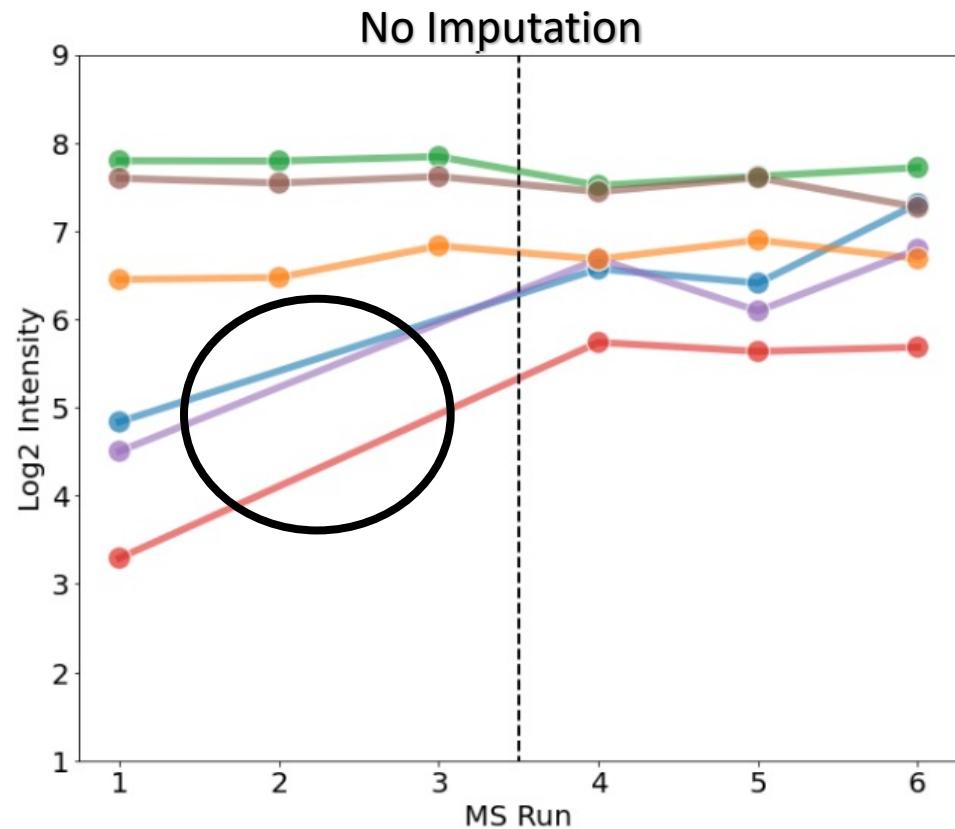
	Condition ₁						...	Condition _I						...	Condition _J						...	Condition _I					
	Subject ₁		Subject ₂		...	Subject _J		...	Subject _{(I-1)J+1}		Subject _{(I-1)J+2}		...	Subject _I		Subject ₁		Subject ₂		...	Subject _J		Subject _I				
	Run ₁	Run ₂	Run ₃	Run ₄	Run ₅	Run ₆	...	Run _{JK-2}	Run _{JK-1}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

Differential Analysis

Presentation Outline

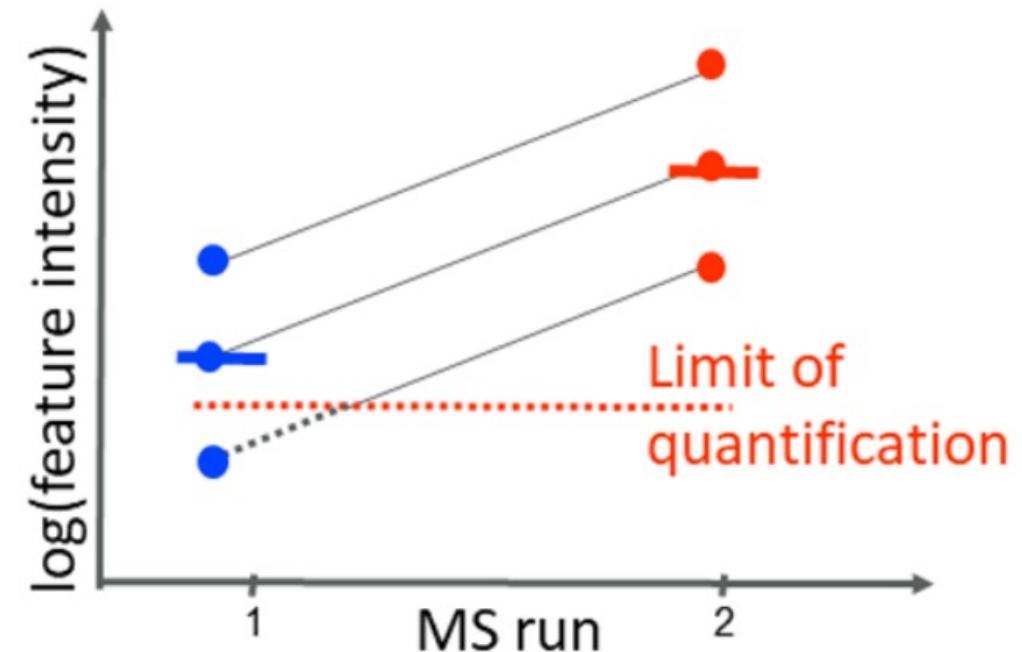
- Problem statement
- Background
 - MSstats model
 - Missing value imputation
 - Summarization
- Proposed method
- Evaluation
 - Missing mechanisms
 - Imputation and summarization
 - Differential analysis

Missing values can have a major impact on protein-level inference



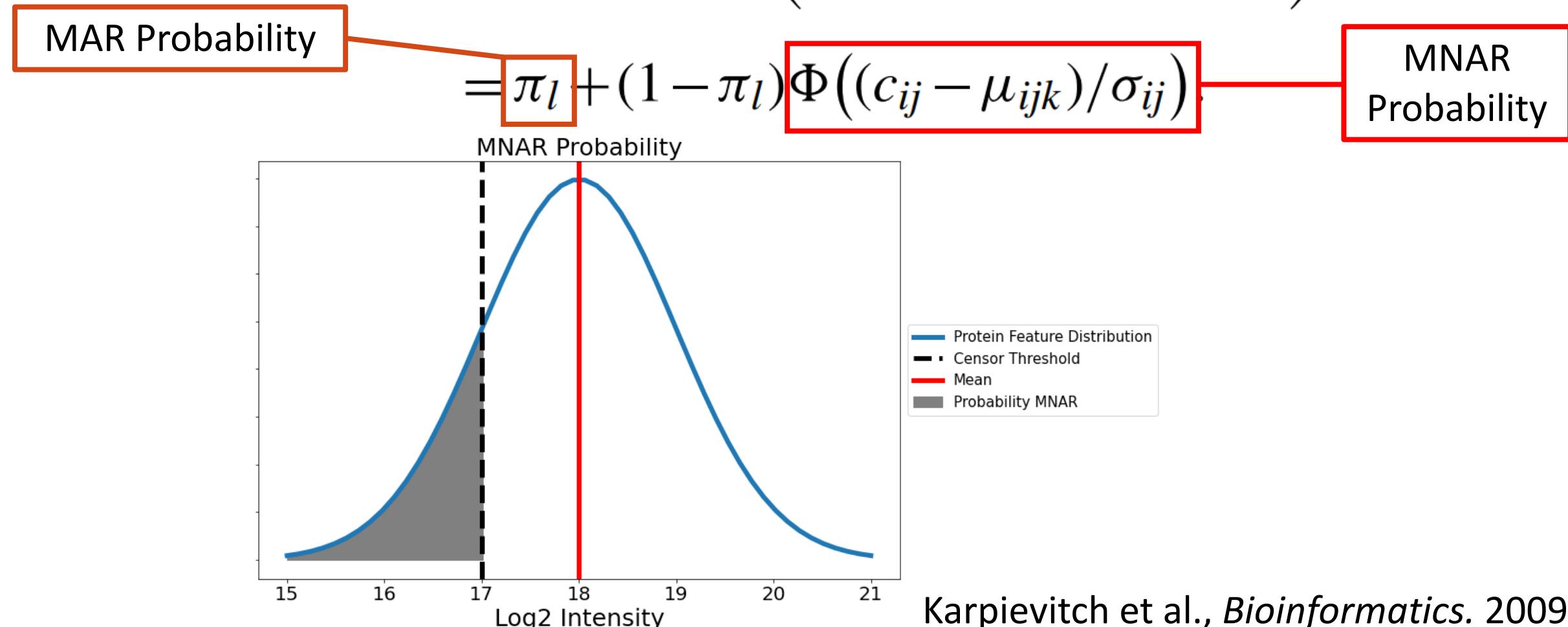
Missing value imputation can create bias, data homogeneity, and are overly confident

- Can create bias if the imputation method is incorrect
 - Missing completely at random (MCAR)
 - Missing not at random (MNAR)
- Causes data to be more homogenous
- Classic imputation returns a point estimate, removing the uncertainty in the estimation process

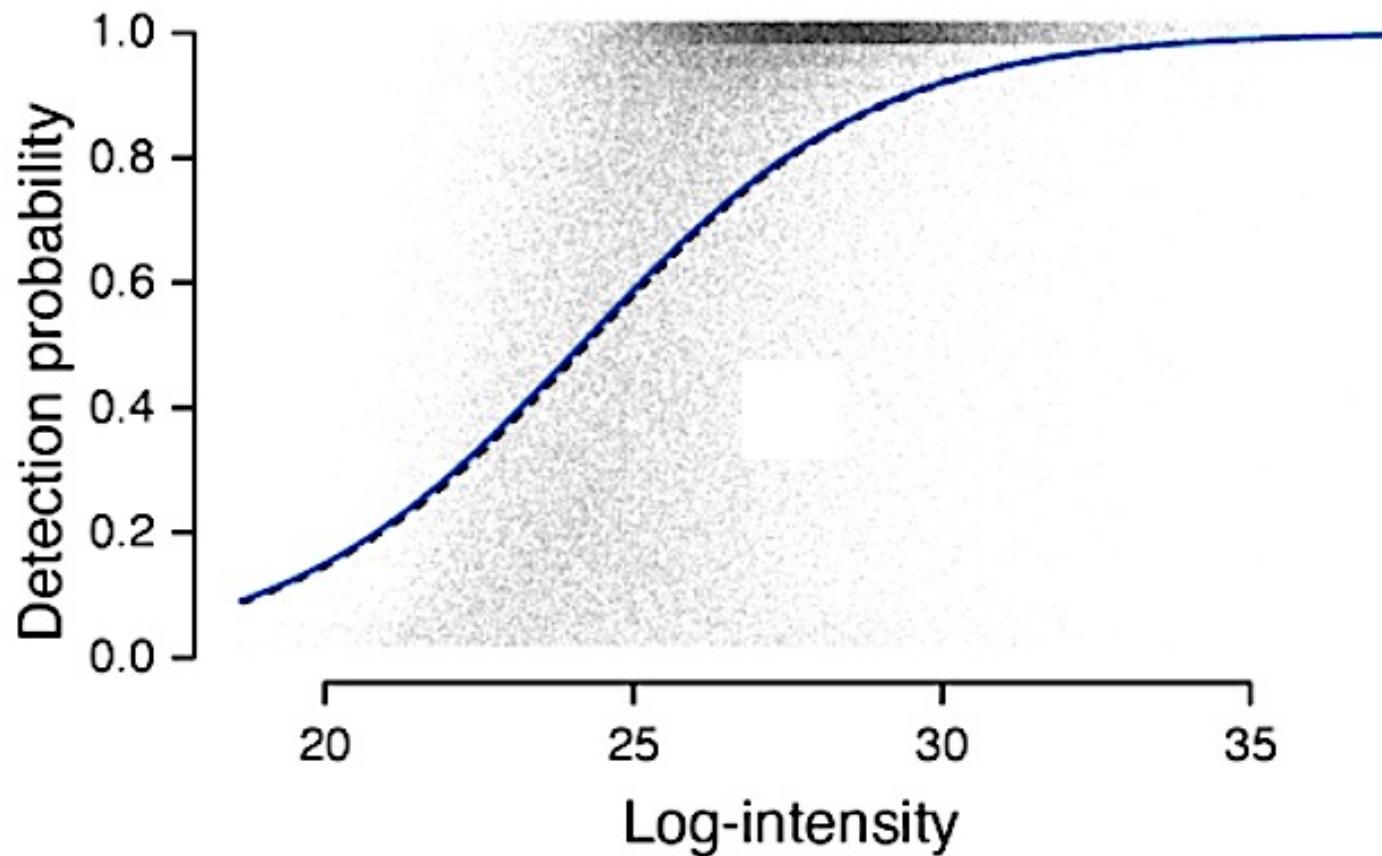


Build missing mechanisms into model

$$P(W_{ijkl} = 1) = 1 - (1 - \pi_l) \left(1 - \Phi\left((c_{ij} - \mu_{ijk})/\sigma_{ij}\right) \right)$$



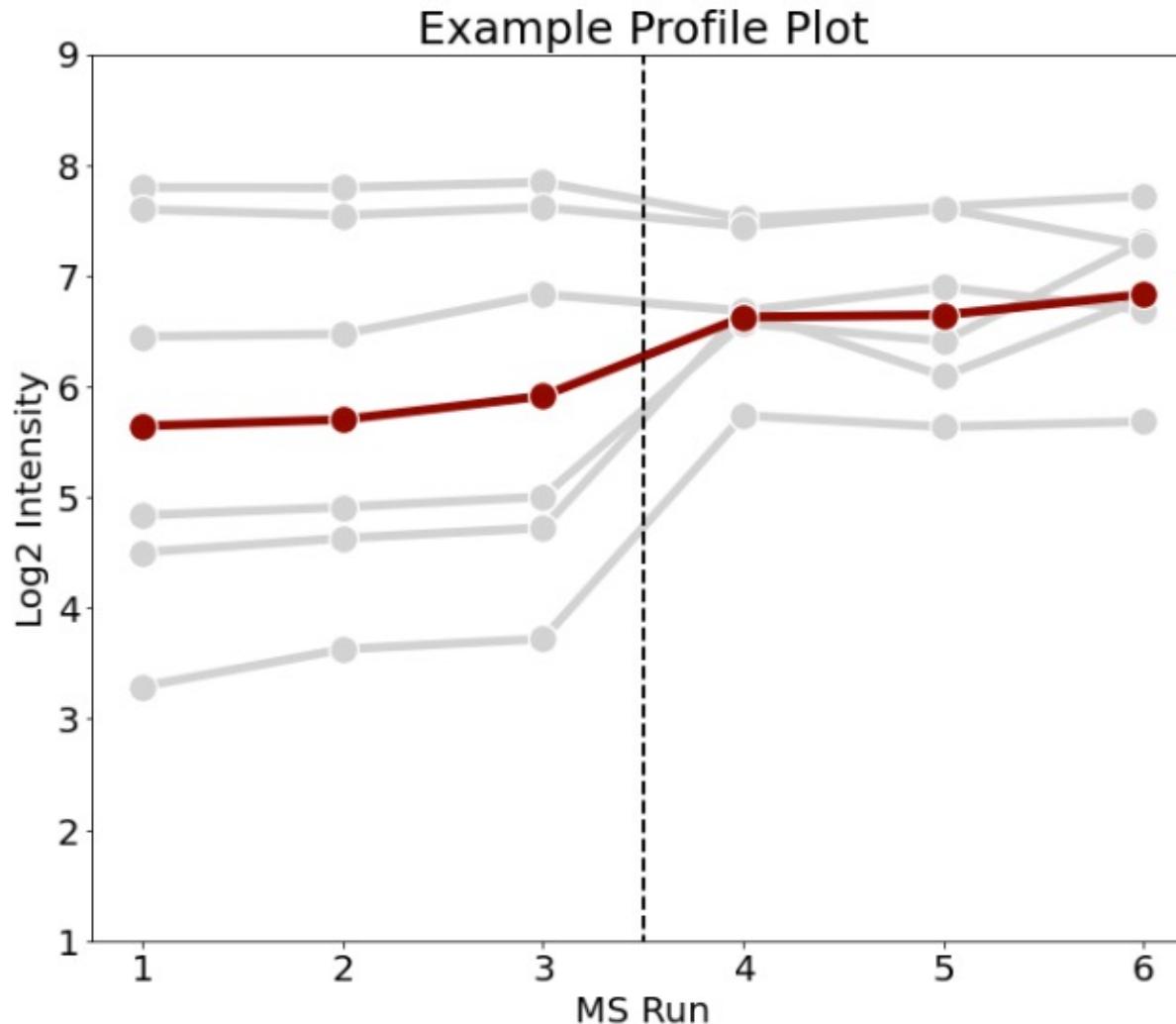
Detection probability curve better represents MNAR data generating process



Presentation Outline

- Problem statement
- Background
 - MSstats model
 - Missing value imputation
 - Summarization
- Proposed method
- Evaluation
 - Missing mechanisms
 - Imputation and summarization
 - Differential analysis

Protein summarization results in point estimates that eliminate uncertainty in the estimation



Presentation Outline

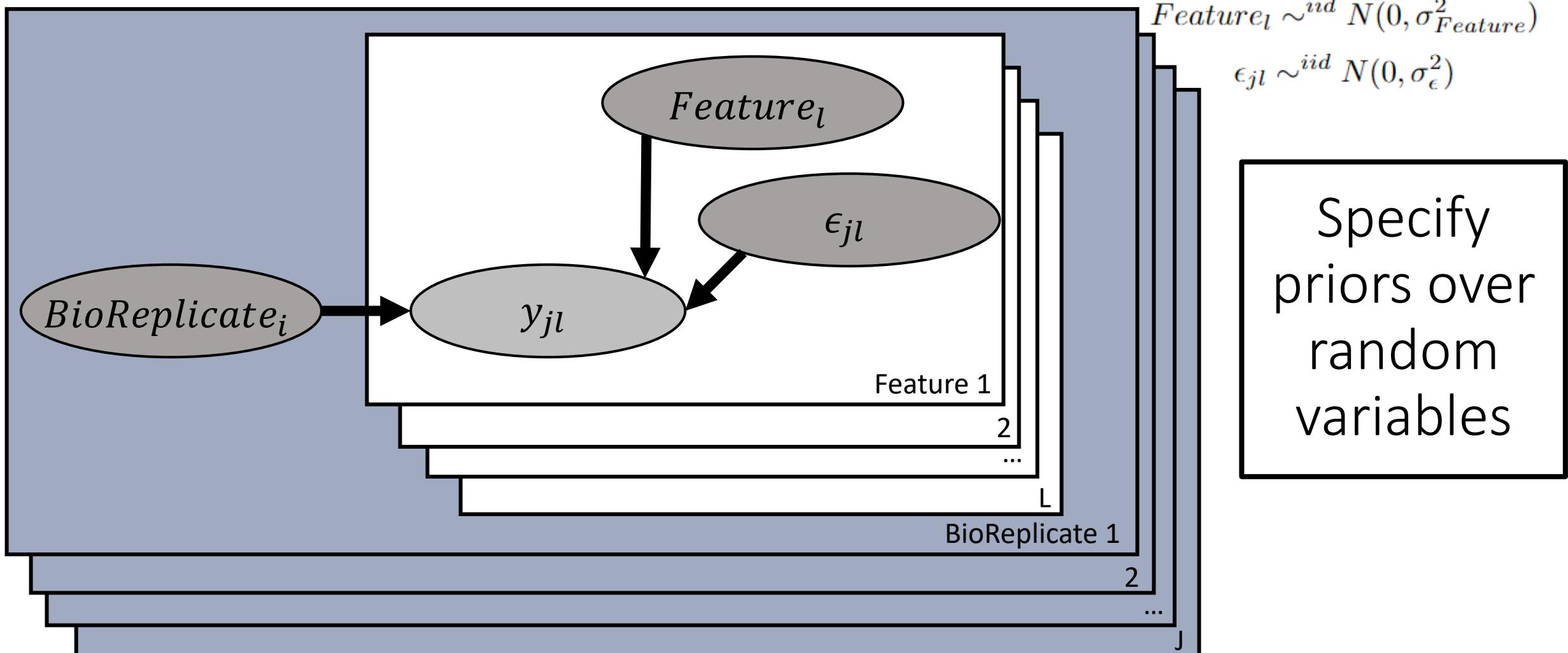
- Problem statement
- Background
 - MSstats model
 - Missing value imputation
 - Summarization
- Proposed method
- Evaluation
 - Missing mechanisms
 - Imputation and summarization
 - Differential analysis

We propose a Bayesian statistical modeling method that quantifies the uncertainty in the imputation and summarization workflow

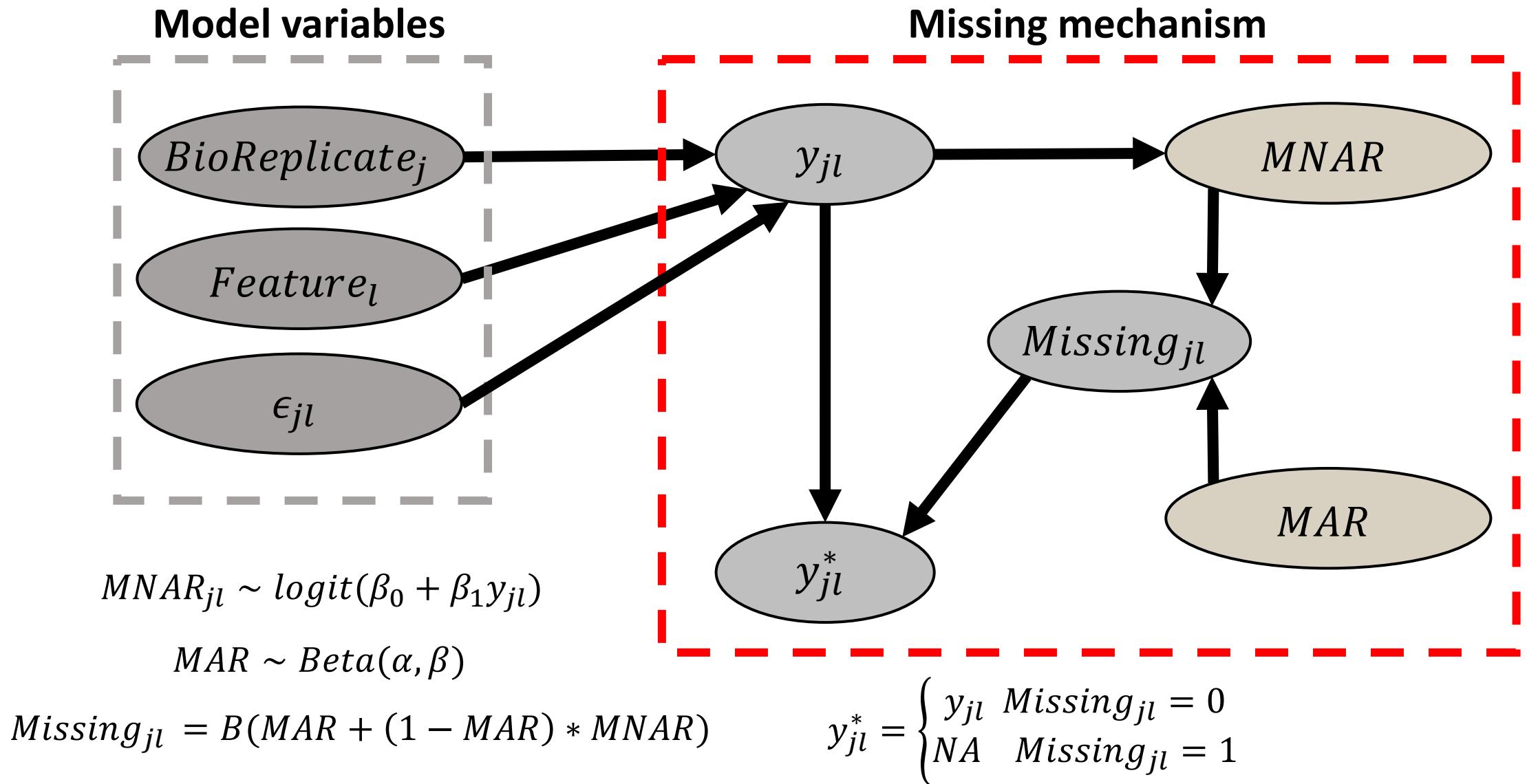
- Counters bias stemming from missing value imputation
- Flexible to a variety of experimental designs and acquisition methods
- Quantifies the uncertainty via posterior estimation

Reformulate subplot model into graphical format

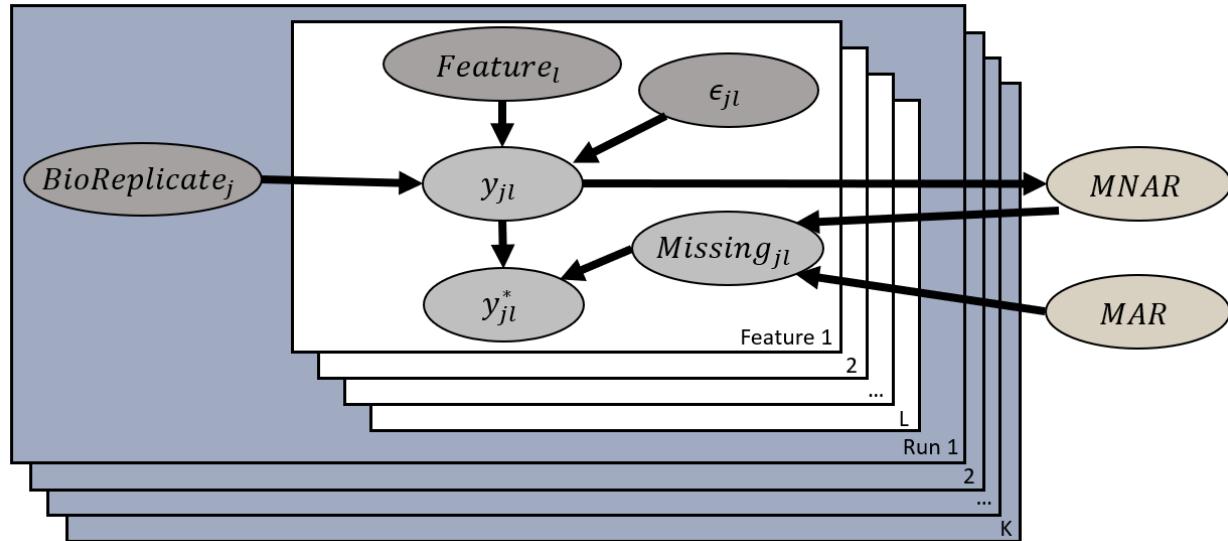
$$y_{jl} = \mu + BioReplicate_j + Feature_l + \epsilon_{jl} \quad \text{where} \quad BioReplicate_j \sim^{iid} N(0, \sigma_{BioReplicate}^2)$$



Graphical encoding of missing mechanism



Graphical and mathematical model



$$y_{jl} \sim N(\mu + BioReplicate_j + Feature_l, \epsilon_{jl})$$

where $BioReplicate_j \sim^{iid} N(0, \sigma_{BioReplicate}^2)$

$$Feature_l \sim^{iid} N(0, \sigma_{Feature}^2), \epsilon_{jl} \sim^{iid} N(0, \sigma_\epsilon^2)$$

Missing mechanism

$$y_{jl} \sim N(\mu + BioReplicate_j + Feature_l, \epsilon_{jl})$$

$$MNAR_{jl} \sim \text{logit}(\beta_0 + \beta_1 y_{jl})$$

$$MAR \sim \text{Beta}(\alpha, \beta)$$

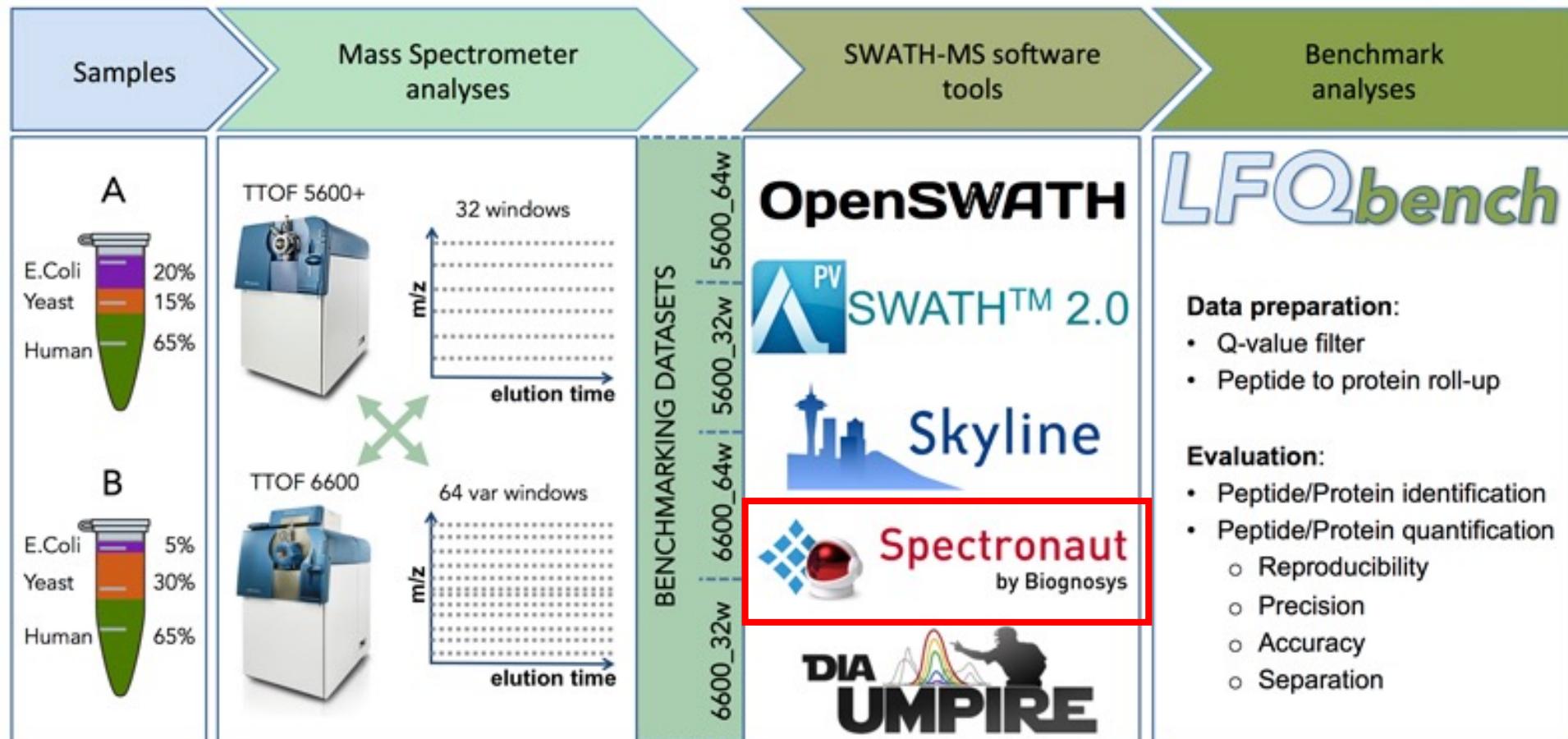
$$Missing_{jl} = B(MAR + (1 - MAR) * MNAR)$$

$$y_{jl}^* = \begin{cases} y_{jl} & Missing_{jl} = 0 \\ NA & Missing_{jl} = 1 \end{cases}$$

Presentation Outline

- Problem statement
- Background
 - MSstats model
 - Missing value imputation
 - Summarization
- Proposed method
- Evaluation
 - Missing mechanisms
 - Imputation and summarization
 - Differential analysis

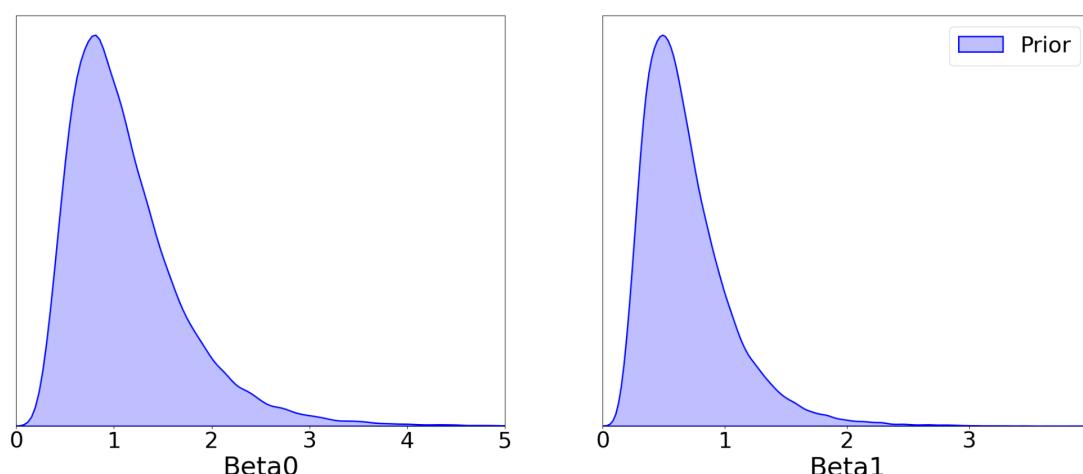
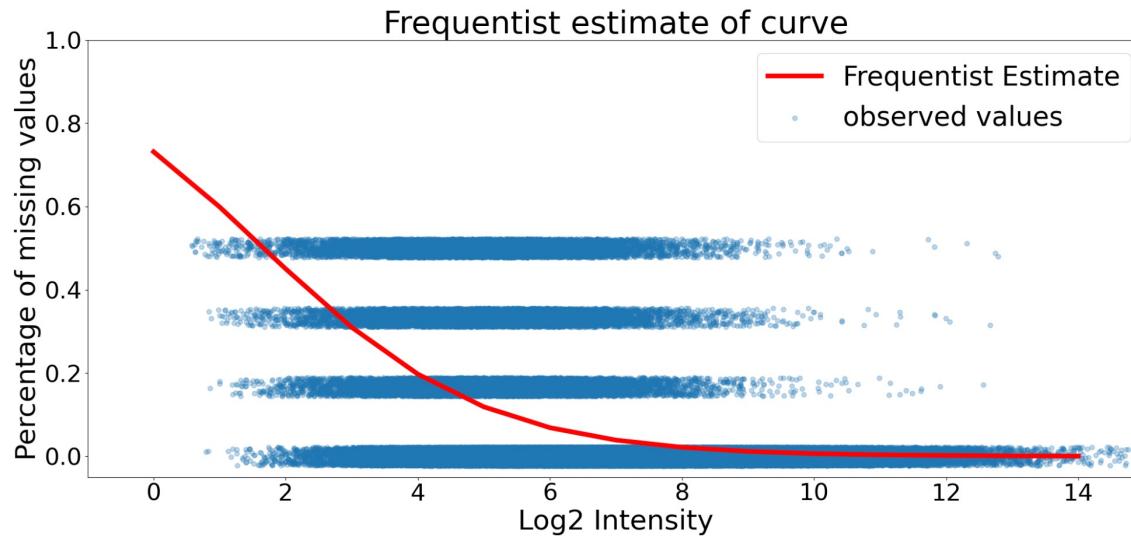
Benchmark dataset



Navarro et al. Nat. Biotech. 2016

Estimate MNAR missing mechanism

Step 1. Learn prior from the data

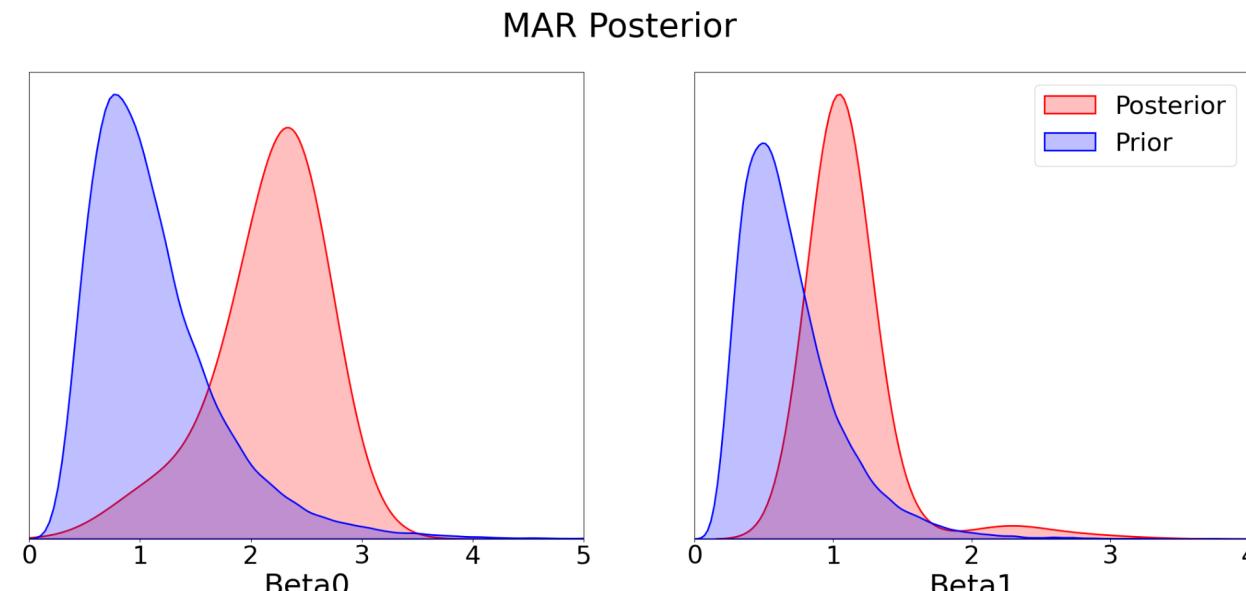


$$MNAR_{jl} \sim \text{logit}(\beta_0 + \beta_1 y_{jl})$$

$$\beta_0 \sim \text{LogNormal}(\mu_{\beta_0}, \sigma_{\beta_0}^2)$$

$$\beta_1 \sim \text{LogNormal}(\mu_{\beta_1}, \sigma_{\beta_1}^2)$$

Step 2. Estimate posterior

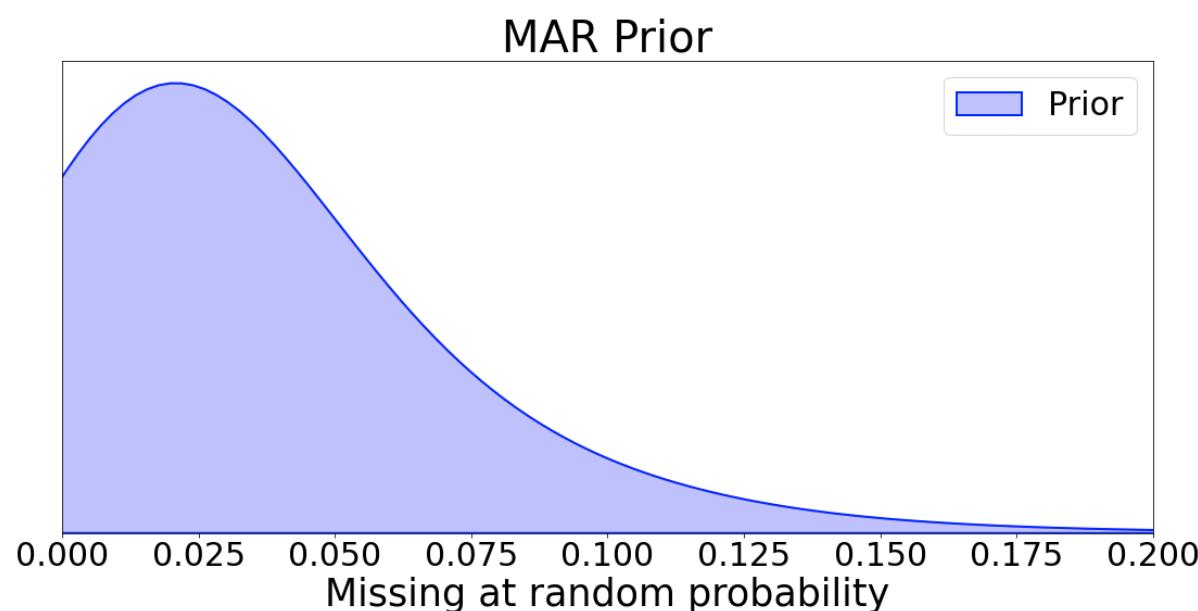


Estimate MAR mechanism

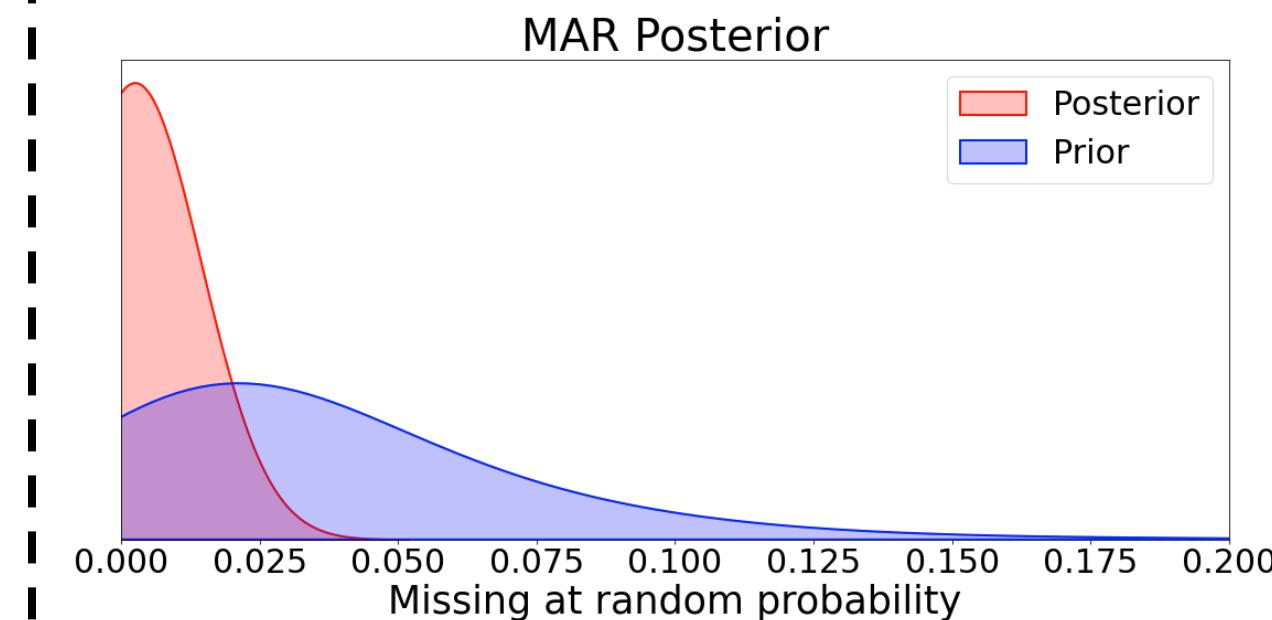
$$MAR \sim Beta(\alpha, \beta)$$

Step 1. Learn prior from the data

Quantile	Avg. log intensity	Avg. missing %
90 th	8.48	2.8%
95 th	9.15	2.2%



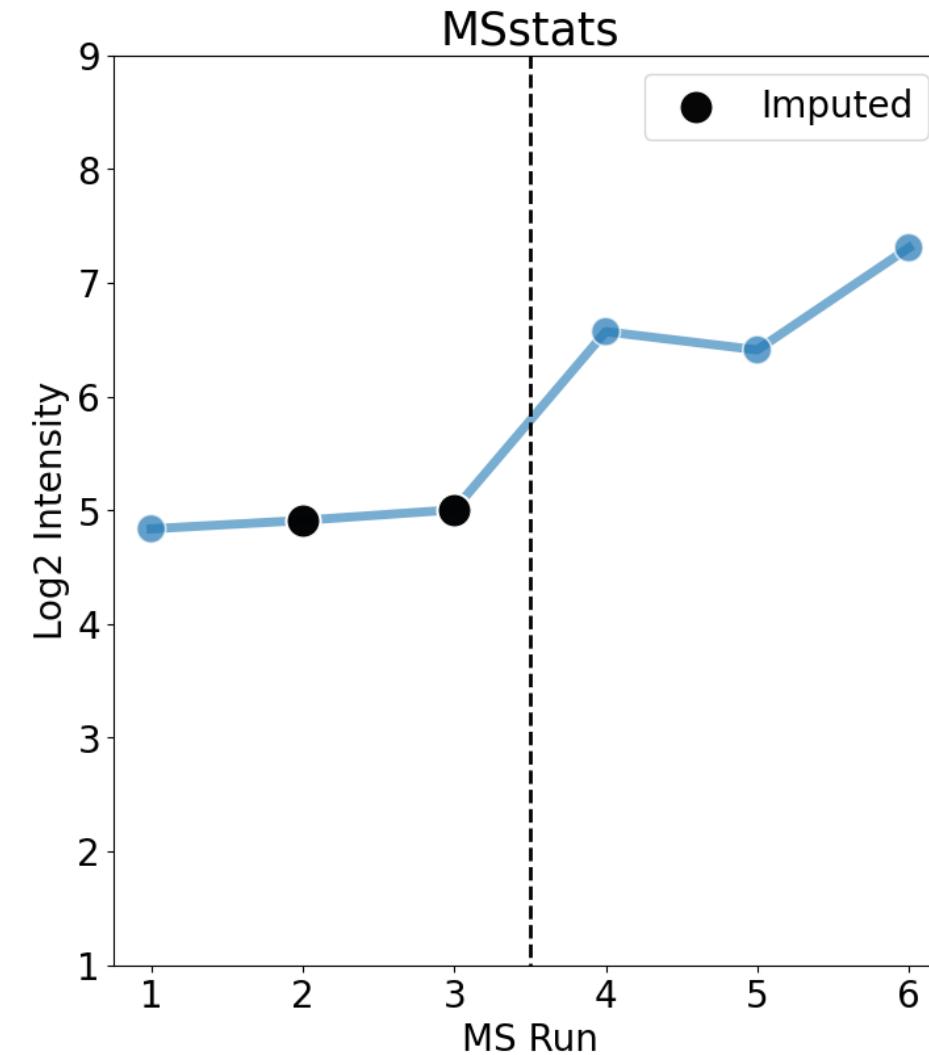
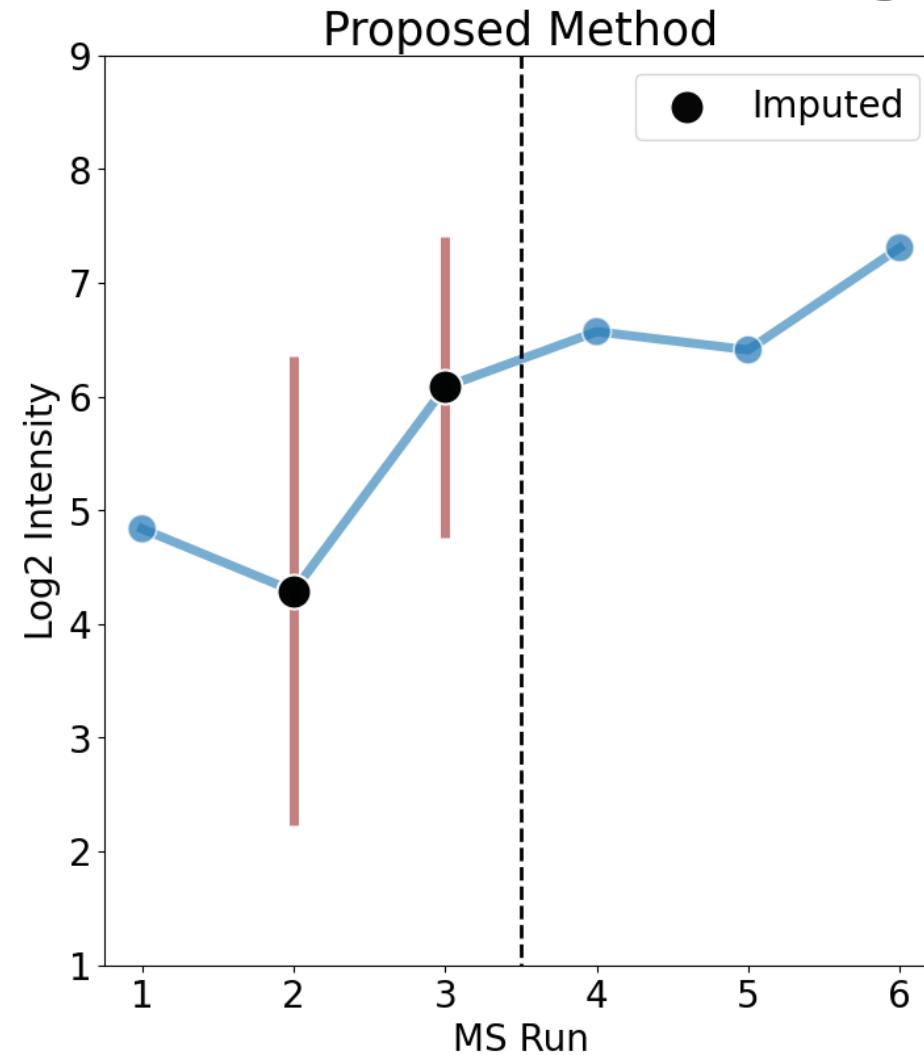
Step 2. Estimate posterior



What do these posteriors
mean for imputation?

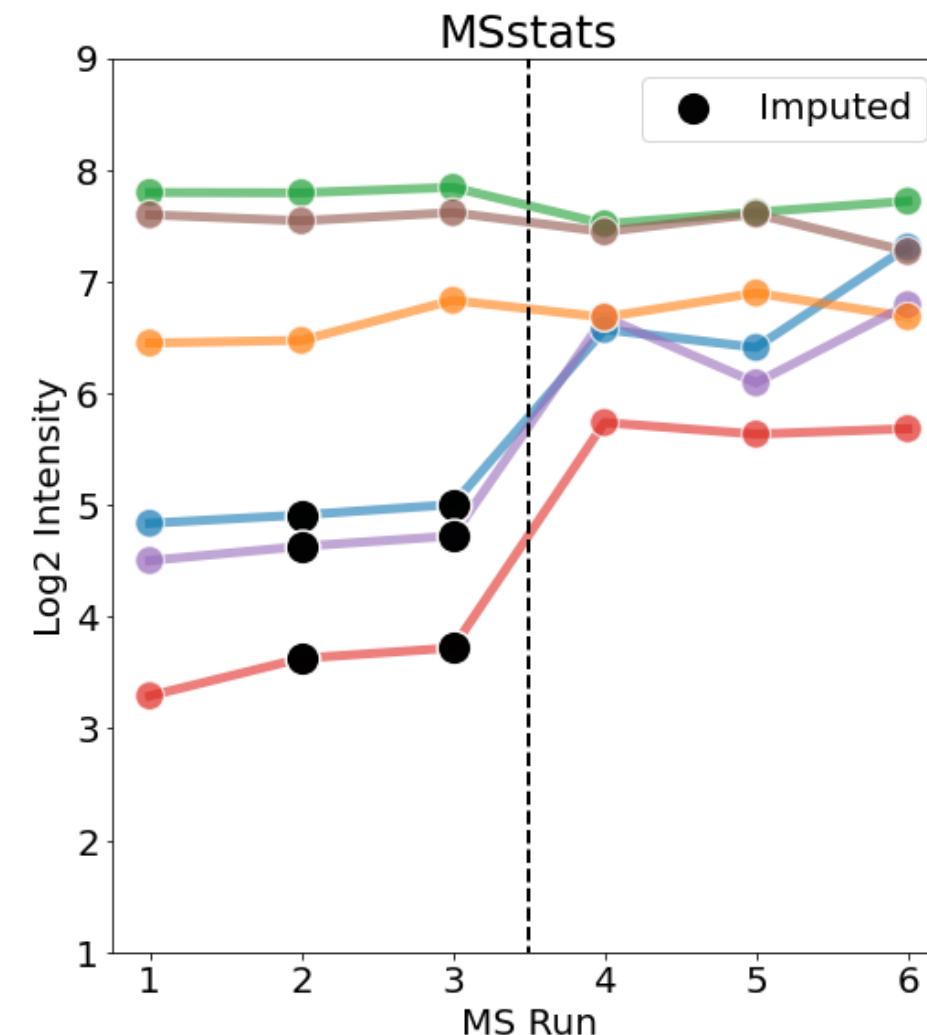
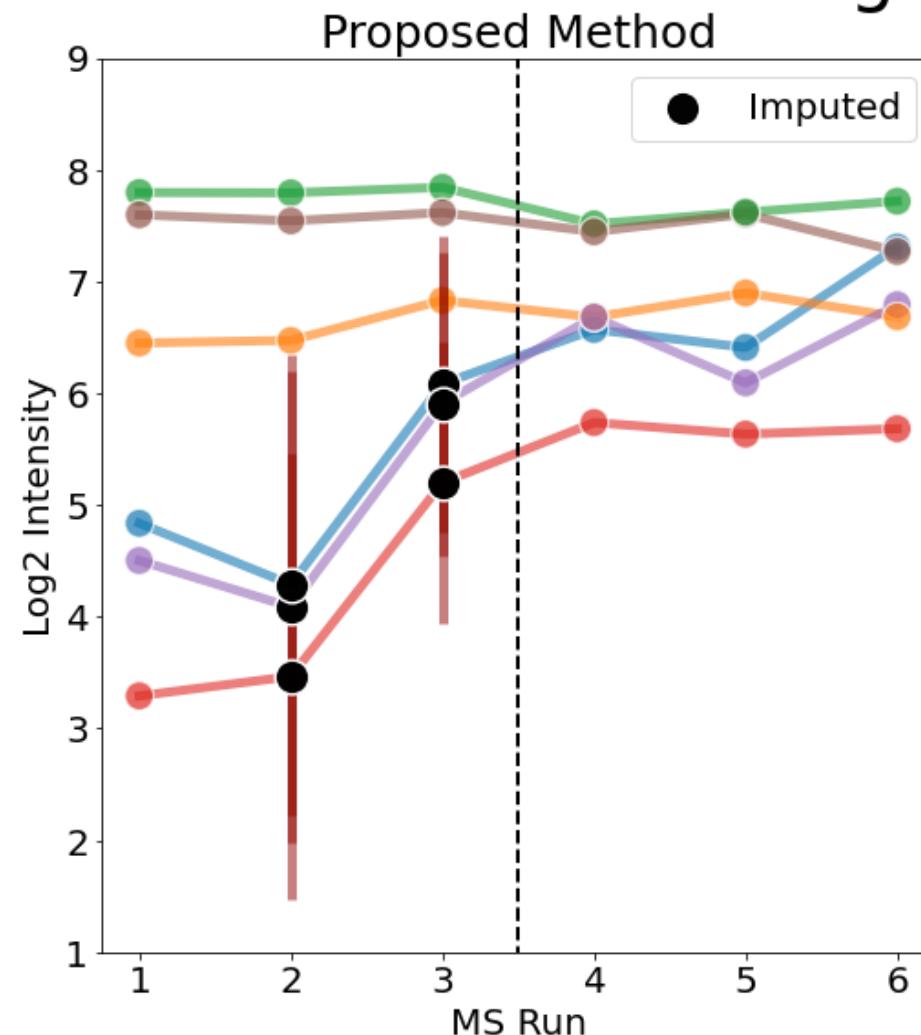
Imputation comparison with MSstats

No change in conditions



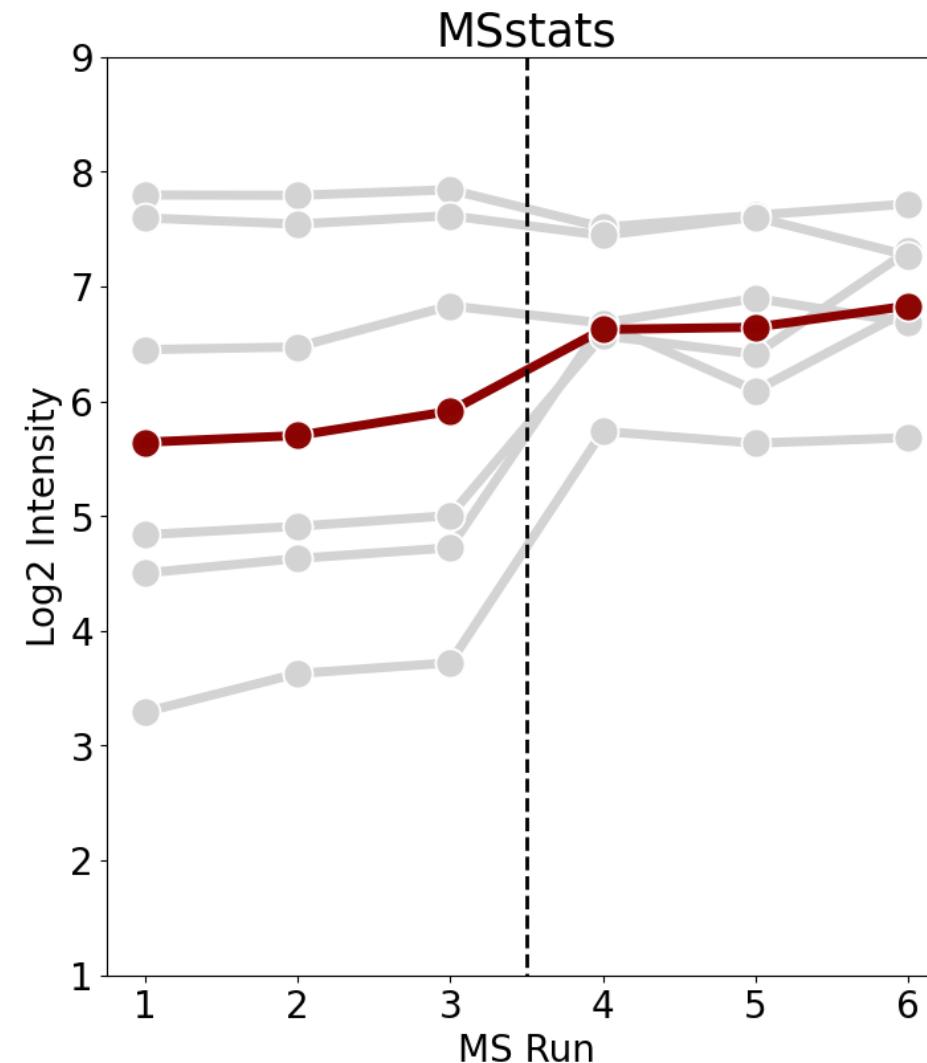
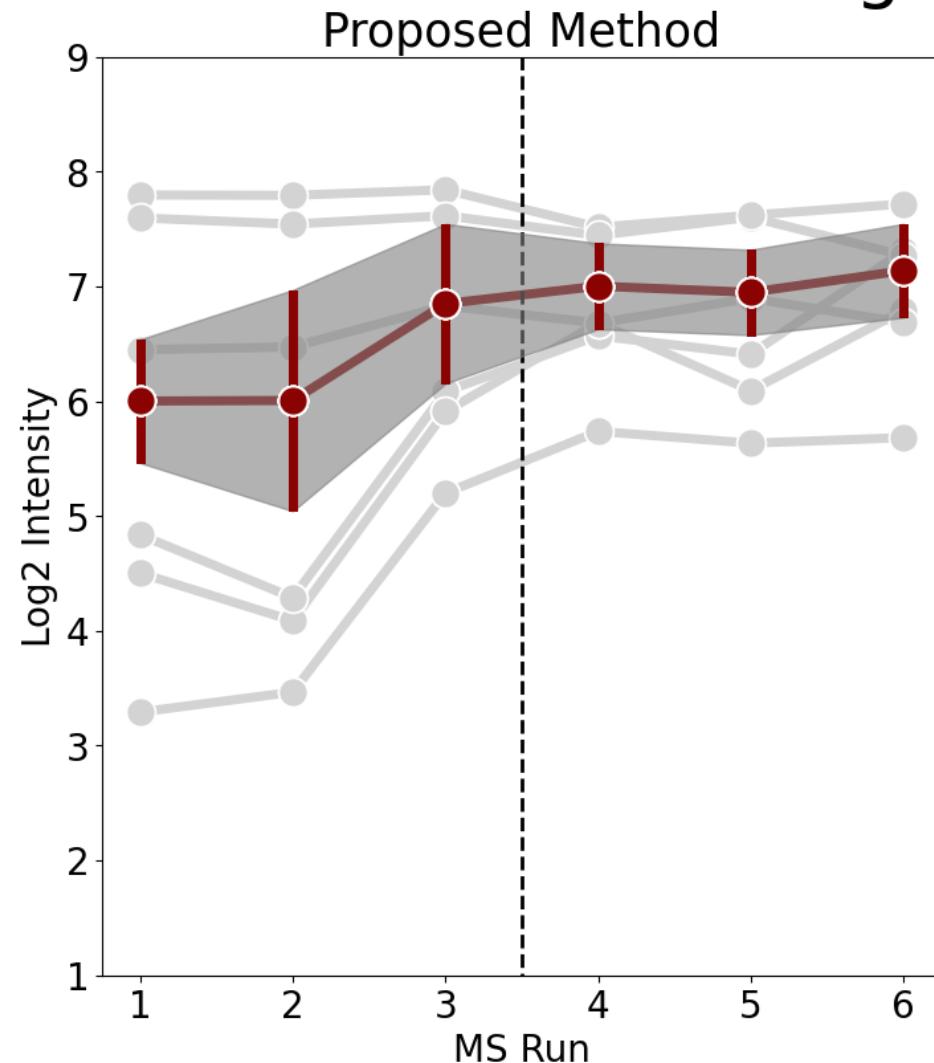
Imputation comparison with MSstats

No change in conditions



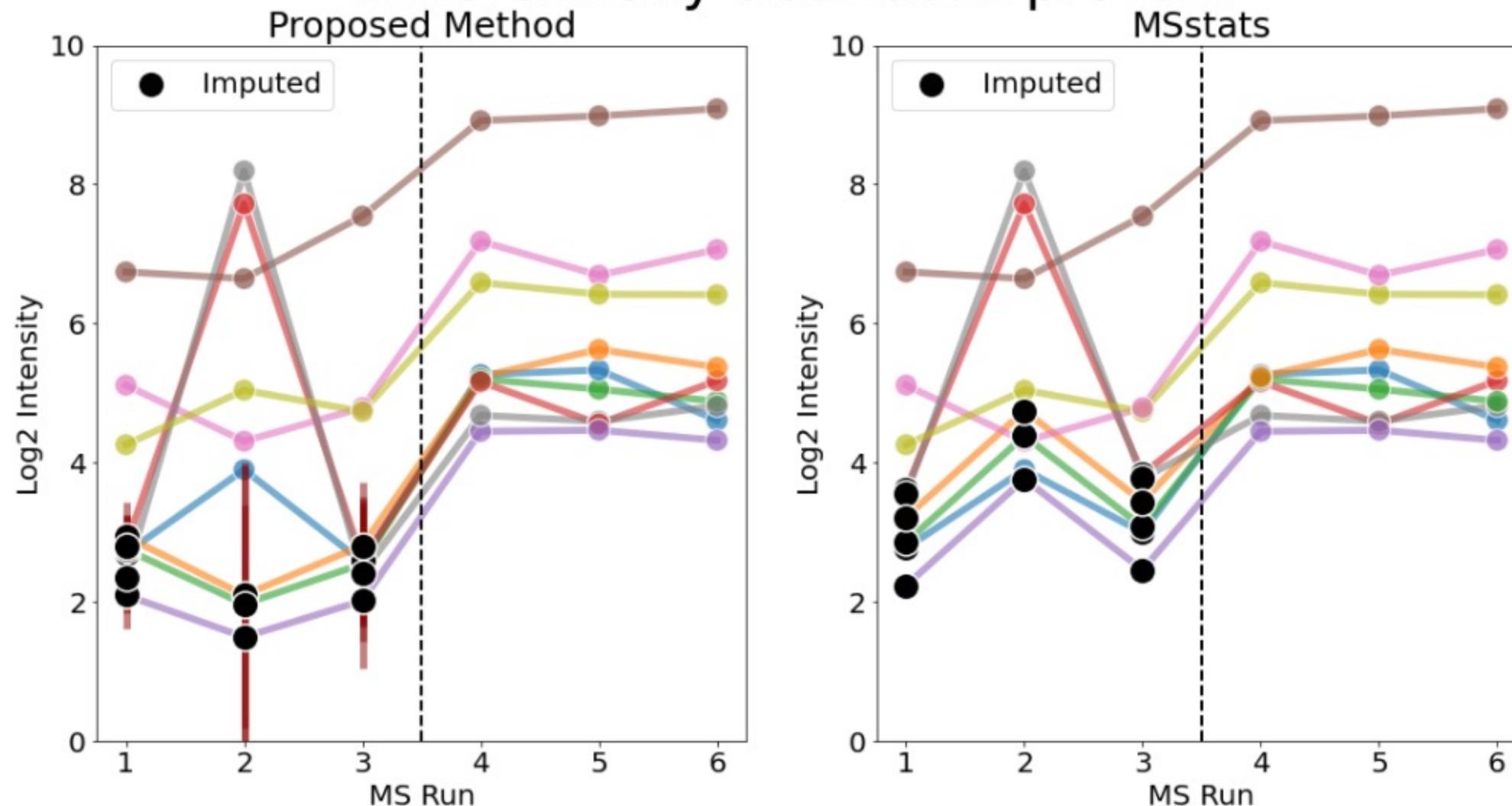
Summarization comparison with MSstats

No change in conditions



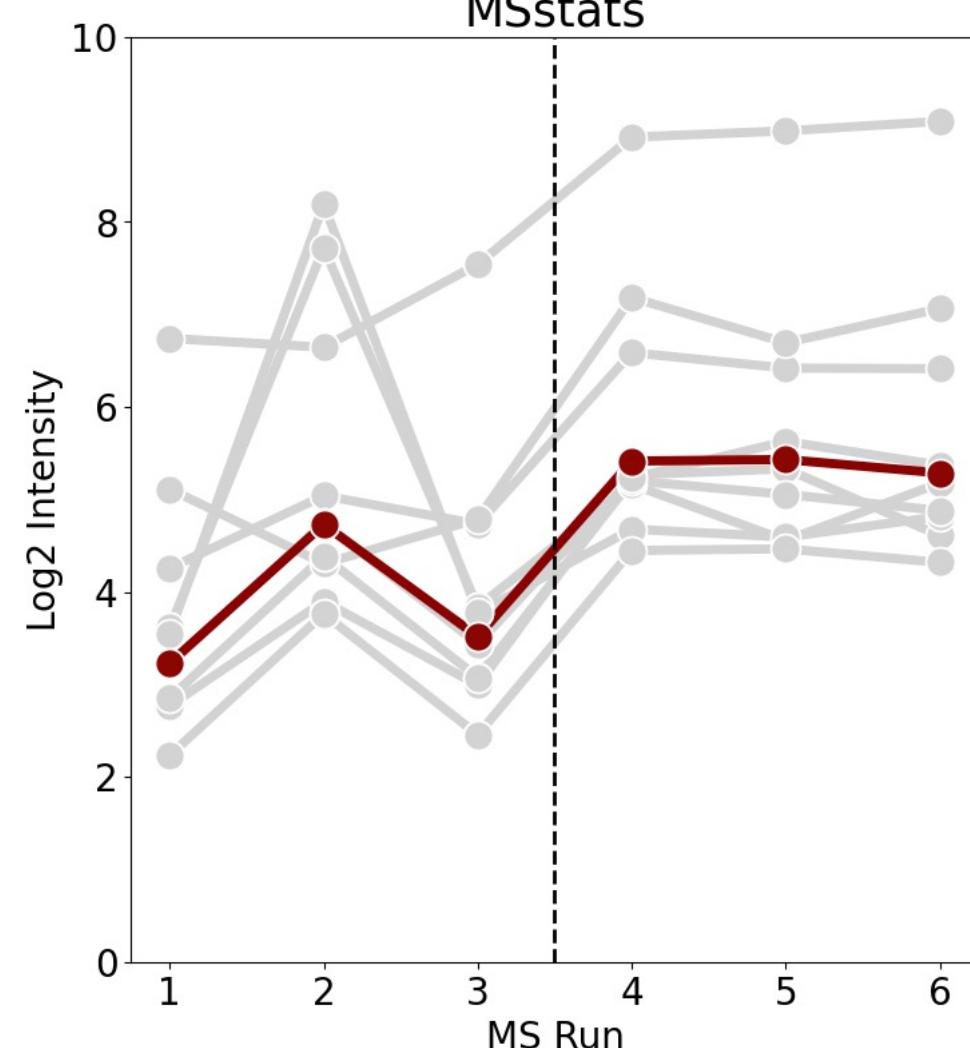
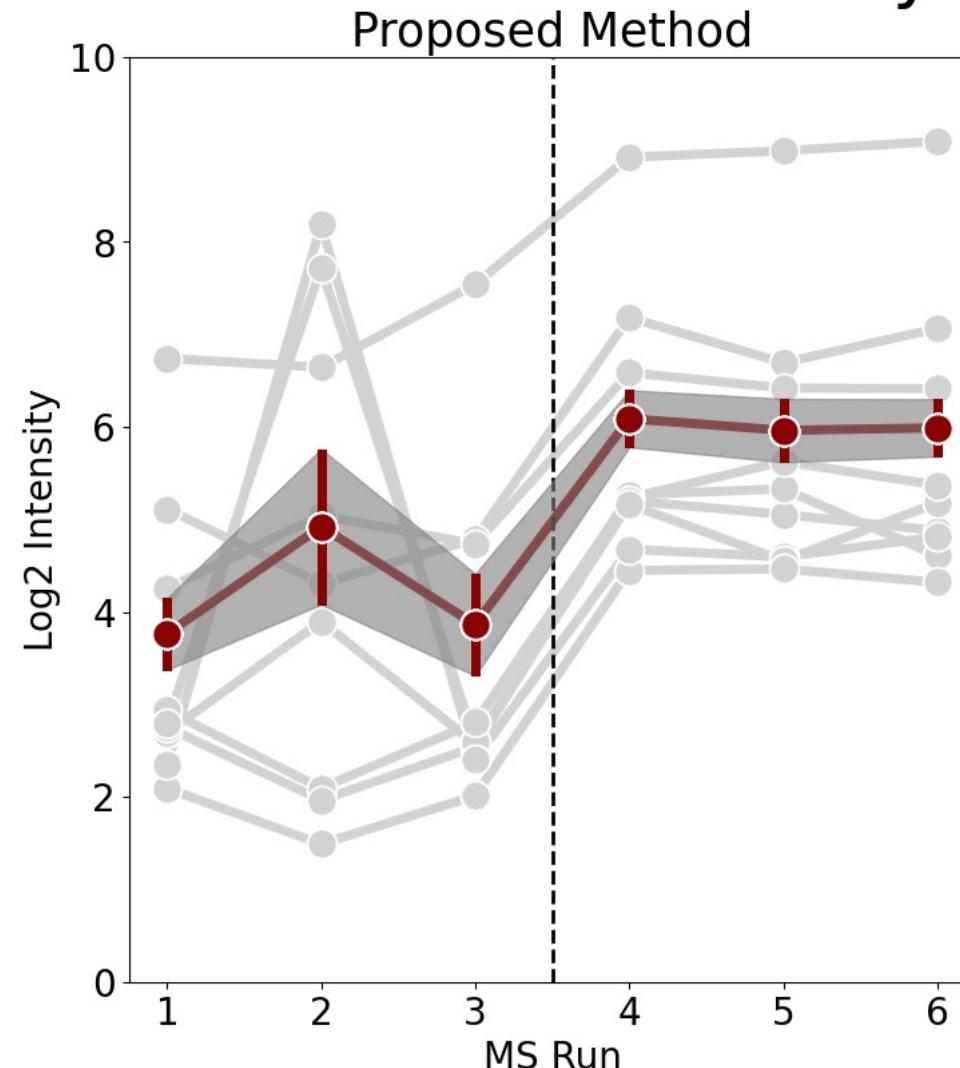
Imputation comparison with MSstats

Differentially abundant protein



Summarization comparison with MSstats

Differentially abundant protein



Overview

- We now have uncertainty measurements associated with both the imputed values, as well as the protein-level intensity estimation per run
- What can we do with these posteriors?
 - Use in other analysis applications
 - Incorporate into differential analysis

Differential analysis with uncertain inputs

- Several options to do this
 - Building condition effect directly into the Bayesian model (one step)
 - Weighted linear regression with unequal error variances (two step)

Weighted least squares regression with unequal variances

MLE linear regression equation
with unequal variances

$$z_{ij} = \mu + Condition_i + \epsilon_{ij}$$

where $\sum_{i=1}^I Condition_i = 0$

$$\epsilon_{ij} = \begin{bmatrix} \sigma_{11}^2 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \sigma_{ij}^2 \end{bmatrix}$$

Define weight w_i as the reciprocal
of the σ_i^2 in maximum likelihood

$$w_{ij} = \frac{1}{\sigma_{ij}^2}$$

Define weighted loss function

$$\sum_{n=1}^{ij} w_n (z_n - Condition_i)^2$$

Weighted least squares regression with unequal variances

MLE linear regression equation
with unequal variances

$$z_{ij} = \mu + \text{Condition}_i + \epsilon_{ij}$$

where $\sum_{i=1}^I \text{Condition}_i = 0$

$$\epsilon_{ij} = \begin{bmatrix} \sigma_{11}^2 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \sigma_{ij}^2 \end{bmatrix}$$

Define weight w_i as the reciprocal
of the σ_i^2 in maximum likelihood

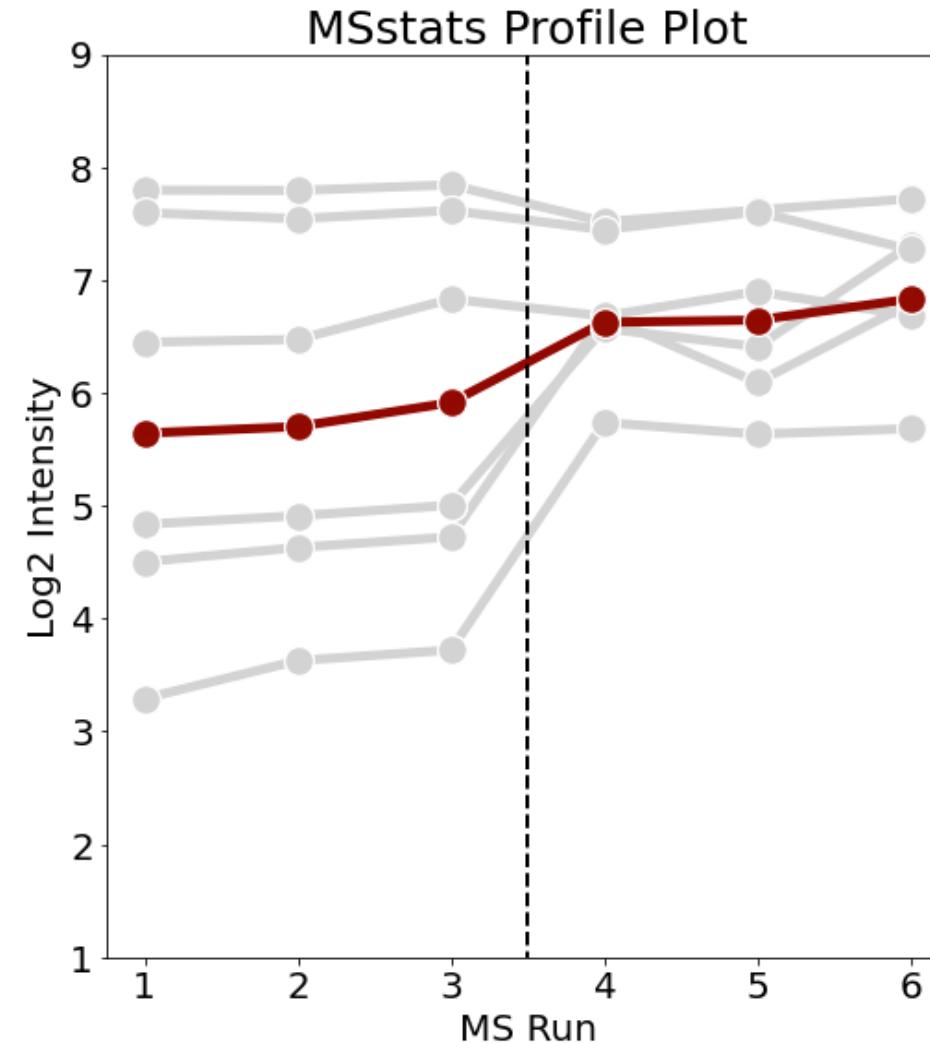
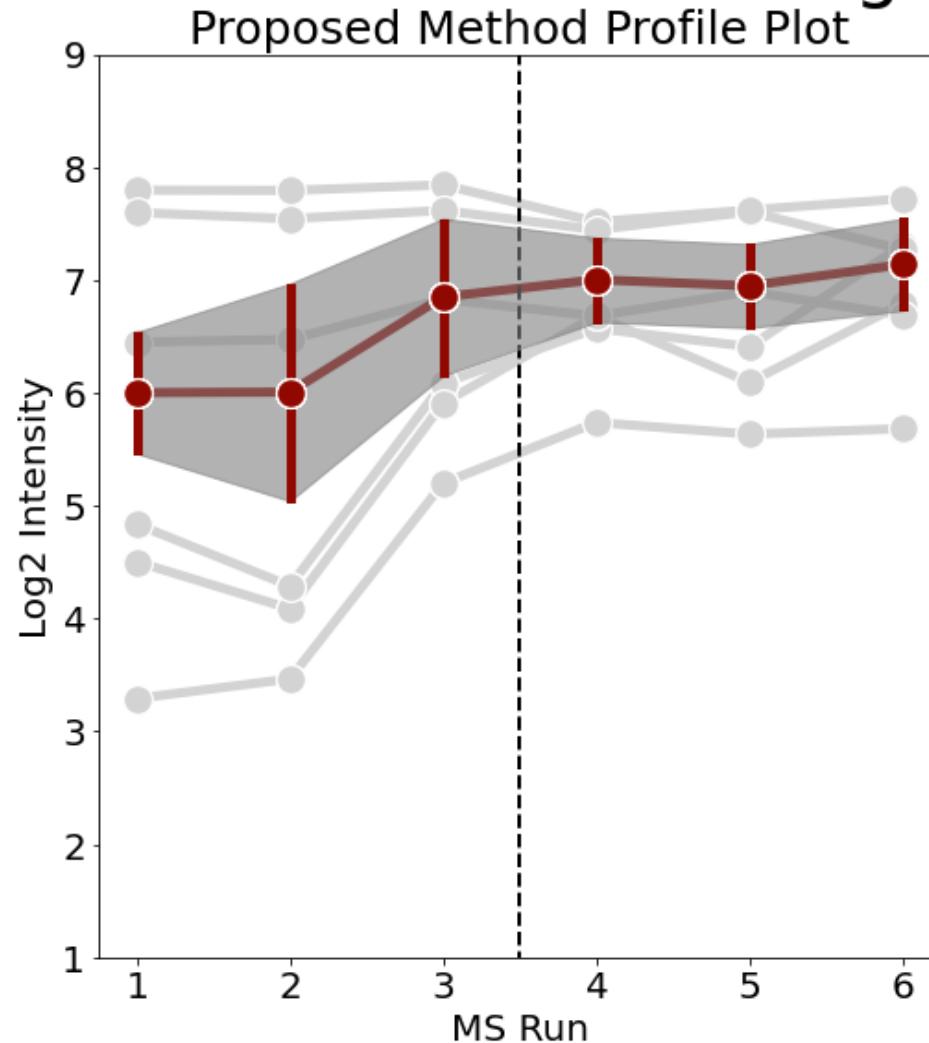
$$w_{ij} = \frac{1}{\sigma_{ij}^2}$$

Define weighted loss function

$$\sum_{n=1}^{ij} w_n (z_n - \text{Condition}_i)^2$$

Summarization comparison

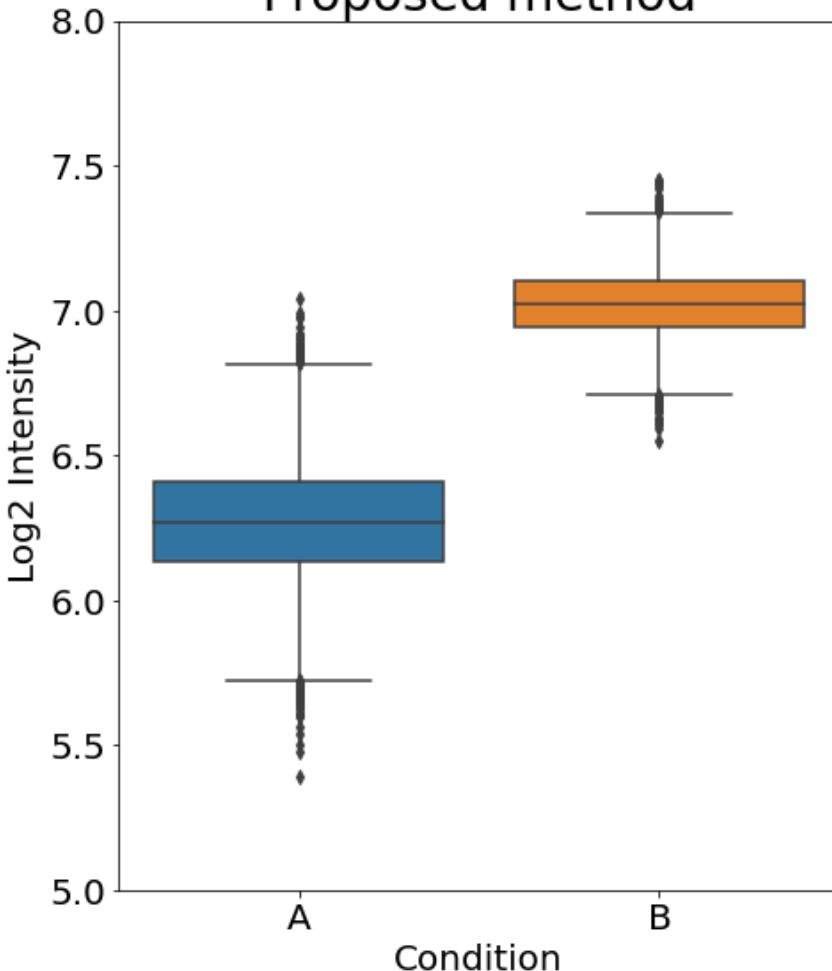
No change in conditions



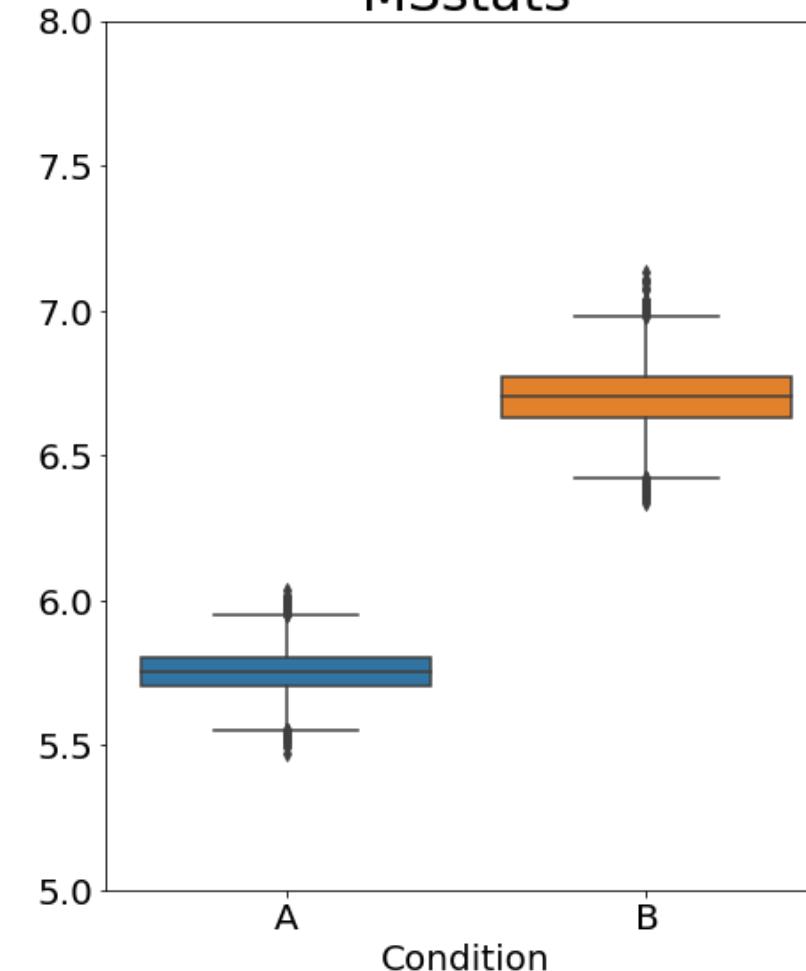
Statistical testing results

No change between conditions

Proposed method



MSstats



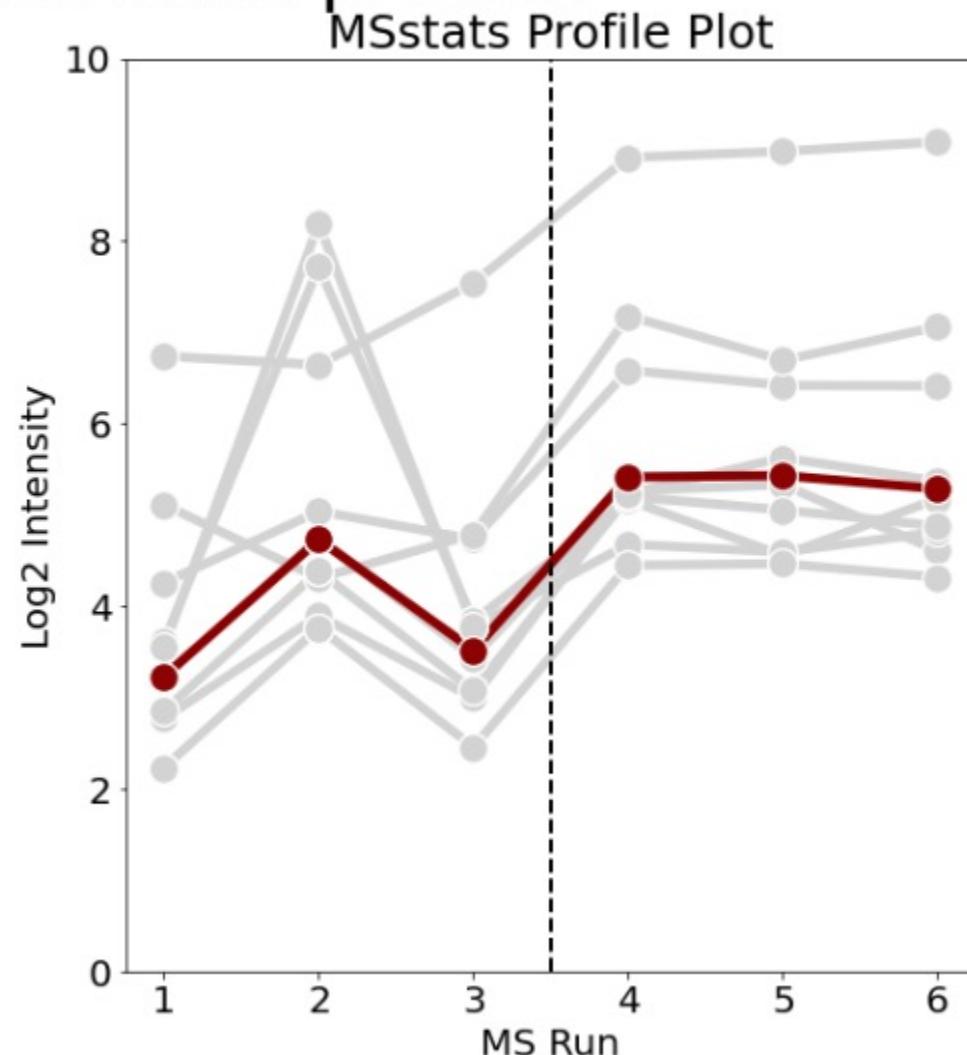
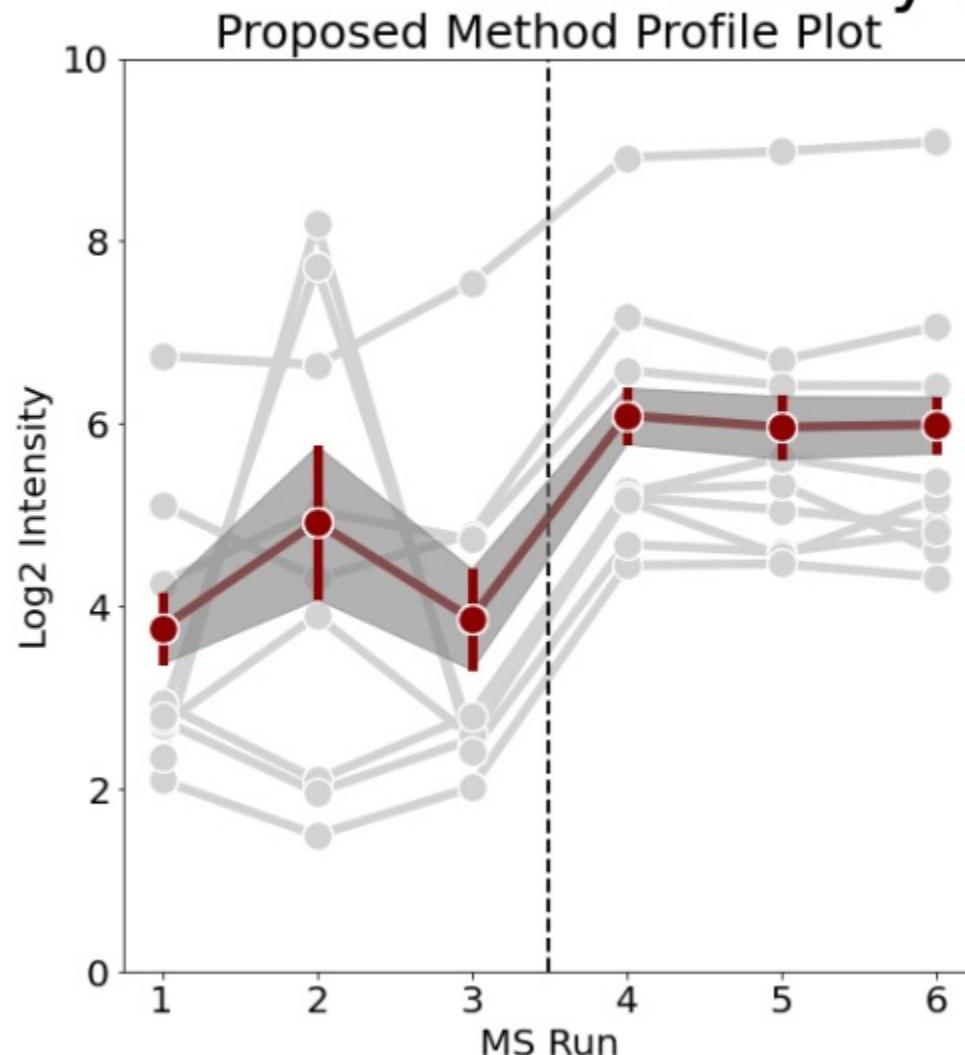
Adjusted p-value

.062

.0023

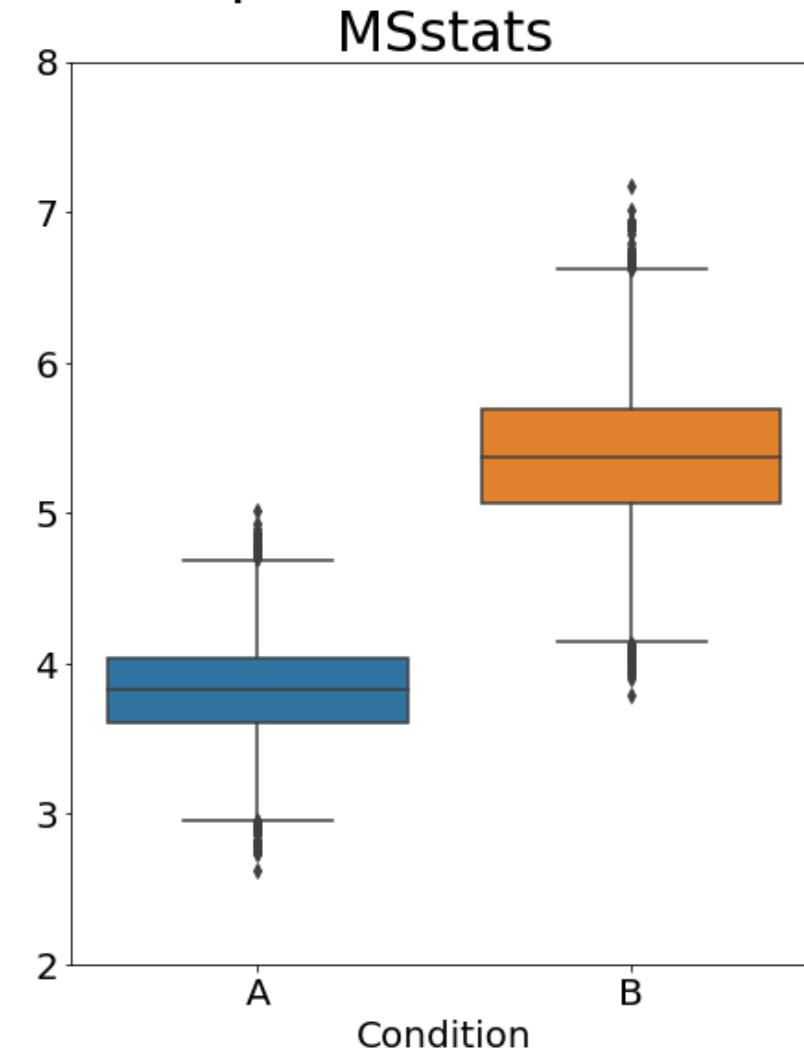
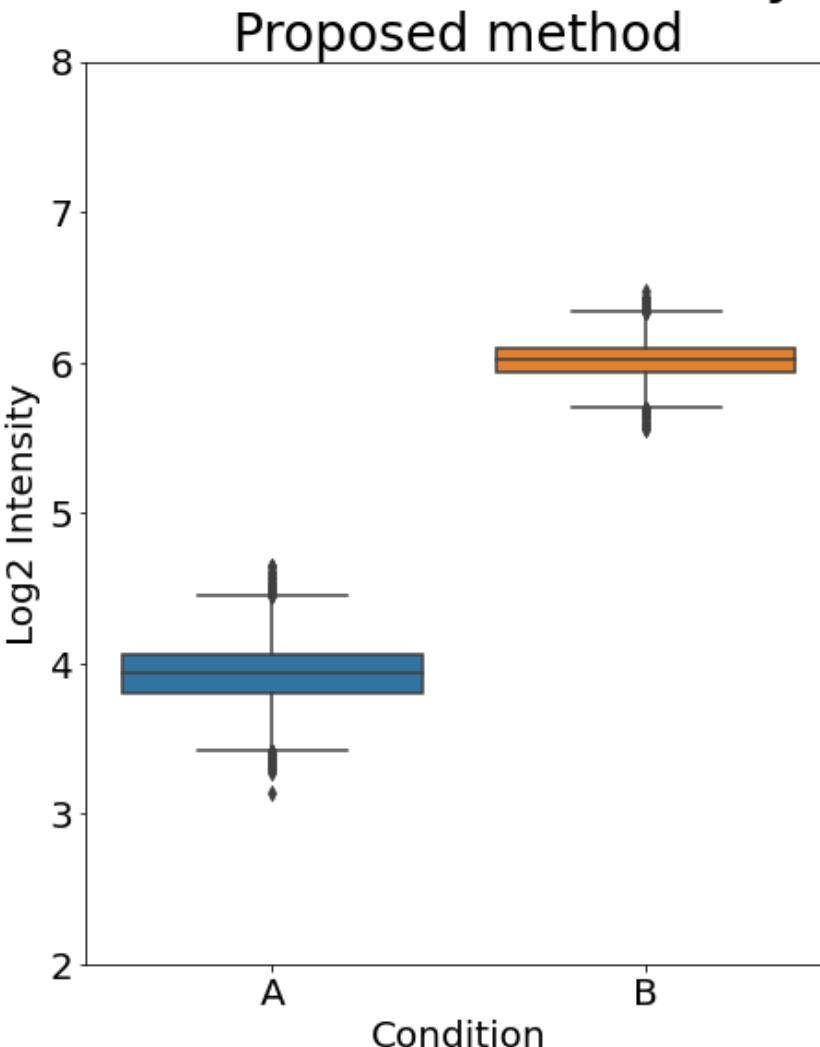
Summarization comparison

Differentially abundant protein



Statistical testing results

Differentially abundant protein



Adjusted p-value

.0024

.053

Concluding points

- Incorporate underlying uncertainty across upstream data processing workflow which was previously lost
- Learns the underlying probabilities in the missing value mechanisms
- Corrects differential analysis in cases of high uncertainty
- Flexible and applicable to a variety of complex experimental designs

Acknowledgements



Northeastern University
OLGA VITEK LAB
Statistical Methods For Studies Of Biomolecular Systems

Lab Members

Kylie Bemis
Sara Mohammad Taheri
Ritwik Anand
Vartika Tewari
Sai Srikanth Lakkimsetty
Mateusz Stankiak



Collaborators

Meena Choi

Jeremy Zucker

Karen Sachs

Barnett Institute travel
award