

# Supplementary Methods and Response to eLife Reviewers

Genetic Disruption of WASHC4 Drives Endo-lysosomal Dysfunction and Cognitive-Movement Impairments in Mice and Humans

Jamie Courtland<sup>1\*</sup>, Tyler W. A. Bradshaw<sup>1\*</sup>, Greg Waite<sup>2</sup>, Erik J. Soderblom<sup>2,3</sup>, Tricia Ho<sup>2</sup>, Anna Rajab<sup>4</sup>, Ricardo Vancini<sup>5</sup>, Il Hwan Kim<sup>2†</sup>, Ting Huang<sup>6</sup>, Olga Vitek<sup>6</sup>, Scott H. Soderling<sup>3</sup>

**Author correspondence:**

[jlc123@duke.edu](mailto:jlc123@duke.edu) (JC); [tyler.w.bradshaw@duke.edu](mailto:tyler.w.bradshaw@duke.edu) (TWAB); [greg.waite@duke.edu](mailto:greg.waite@duke.edu) (GW); [erik.soderblom@duke.edu](mailto:erik.soderblom@duke.edu) (EJB); [tricia.ho@duke.edu](mailto:tricia.ho@duke.edu) (TH); [drannarajab@gmail.com](mailto:drannarajab@gmail.com) (DR); [ricardo.vancini@duke.edu](mailto:ricardo.vancini@duke.edu) (RV); [ikim9@uthsc.edu](mailto:ikim9@uthsc.edu) (IK); [huang.tin@northeastern.edu](mailto:huang.tin@northeastern.edu) (TH); [o.vitek@northeastern.edu](mailto:o.vitek@northeastern.edu) (OV); [scott.soderling@duke.edu](mailto:scott.soderling@duke.edu) (SHS)

\*These authors contributed equally to this work.

**Present address:**

<sup>†</sup>Department of Anatomy and Neurobiology, University of Tennessee Health Science Center, Memphis, TN 38163, USA

<sup>1</sup>Department of Neurobiology, Duke University School of Medicine, Durham, NC 27710, USA; <sup>2</sup>Proteomics and Metabolomics Shared Resource, Duke University School of Medicine, Durham, NC 27710, USA; <sup>3</sup>Department of Cell Biology, Duke University School of Medicine, Durham, NC 27710, USA; <sup>4</sup>Burjeel Hospital, VPS Healthcare, Muscat, Oman; <sup>5</sup>Department of Pathology, Duke University School of Medicine, Durham, NC 27710, USA; <sup>6</sup>Khoury College of Computer Sciences, Northeastern University, Boston, MA 02115, USA

## Abstract

In the review of this manuscript significant concerns were raised by the reviewers about the validity of our approach to perform protein- and module-level inference from our TMT proteomics dataset. Here we address these concerns and provide additional description of our statistical approach.

## Reanalysis of SWIP<sup>P1019R</sup> TMT Proteomics

At the center of the cogent critique of our manuscript were questions about statistical validity of our previously described approach. Succinctly, the issue at question is whether or not the R package edgeR is an appropriate tool for analysis of protein mass spectrometry data.

PreviousAt the Plubell et al., 2017 Several tools for analysis of Proteomics data exist. At least one prior publicTwo previous publications describe the use of edgeR for analysis of TMT

Previously we used a customized workflow<sup>1</sup> to preprocess and normalize the data prior to performing statistical testing using edgeR's flexible GLM framework.

At the core of our normalization approach is the use of common quality control (QC) samples, analyzed in technical duplicate in each TMT mixture. The sample pool QC sample was formed by combining aliquots of all biological replicates. This sample therefore is representative biological and technical variation. Our normalization approach, termed IRS normalization, is essential to correct for

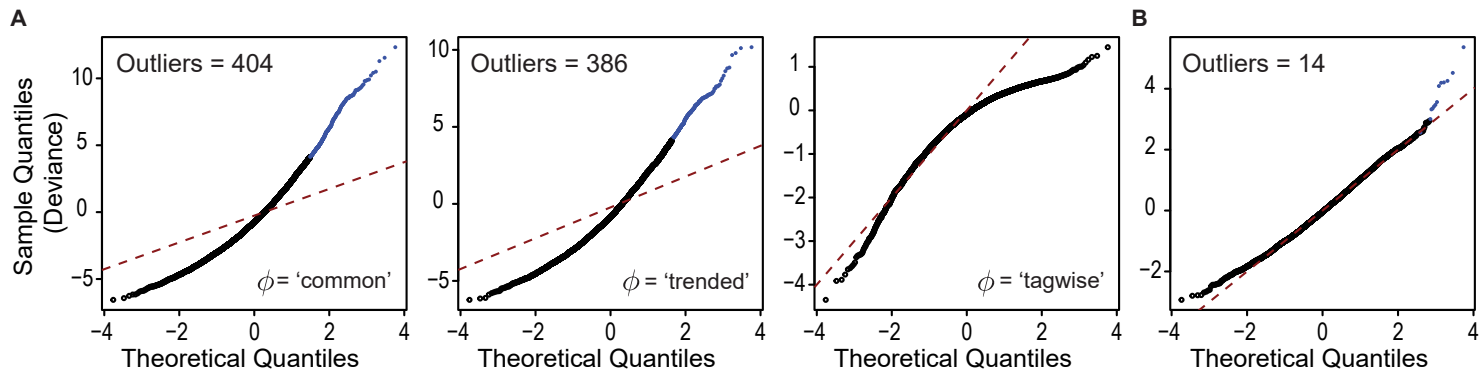
The most important step in our normalization approach is IRS normalization. MS2 random sampling results in identification and quantification of proteins by different peptides in each MS experiment. To account for this source of variability, protein measurements are adjusted by a scaling factor such that the geometric mean of all internal reference standards are equal (Plubell et al., 2017). This is essential to account for the stochasticity of peptide quantification in MS experiments. Phillip Wilmarth's [GitHub] offers an excellent exploration of IRS normalization.

However, we failed to thoughtfully consider the overall adequacy of the NB framework for mass spectrometry data. Here we reconsider its appropriateness for our TMT proteomics dataset.

Our previous approach can be summarized as the "Sum + IRS Normalization" method described by Huang et al., 2020. We summarize proteins as the sum of its features and use IRS normalization. A potential weakness of the sum method is that outlier peptides can strongly influence the summary value. Median-based methods (median and median-polish) avoid this problem. We examined outliers in a method described by [xxx]

Statistical inference in edgeR is built on a negative binomial (NB), generalized linear model (GLM) framework. Therefore, the data are assumed to be adequately described by a NB distribution parameterized by a dispersion parameter,  $\phi$ .<sup>2</sup>

<sup>2</sup> The dispersion parameter can take several forms. edgeR supports three dispersion models: 'common', 'trended', and 'tagwise'. However, when using edgeR's robust quasi-likelihood test methods, only global (i.e. 'common' or 'trended') dispersion metrics are appropriate (see



**Figure 1. Goodness-of-fit of edgeR (A), and MSstats (B) statistical approaches.** The overall adequacy of the linear models fit to the data were assessed by plotting the residual deviance for all proteins as a quantile-quantile plot (McCarthy *et al.*, (2012)). **(A)** The normalized protein data were fit with a NB GLM of the form:  $\text{Abundance} \sim \text{Mixture} + \text{Condition}$ . Where *Mixture* is a blocking factor that accounts for sources of variability between experiments. Protein-wise deviance statistics were transformed to normality and plotted against theoretical normal quantiles using edgeR: `:gof`. **(B)** The normalized protein data were fit with a linear mixed-effects model (LMM) of the form:  $\text{Abundance} \sim 0 + \text{Condition} + (1|\text{Mixture})$ . Where *Mixture* indicates the random effect of *Mixture*. The residual deviance and degrees of freedom were extracted from the fitted models, z-score normalized, and plotted as in (A). Proteins with a significantly poor fit are indicated as outliers in blue (Holm-adjusted P-value < 0.05).

We evaluated the overall adequacy of the edgeR model by plotting the residual deviance of all proteins against their theoretical, normal quantiles in a quantile-quantile plot. **Figure 1** illustrates the overall lack of fit for the three dispersion models fit by edgeR. As an alternative to edgeR we considered MSstatsTMT, an extension of MSstats for analysis of TMT proteomics experiments.

MSstatsTMT utilizes a linear mixed-model framework. The strength of linear mixed models (LMMs) is in their ability to account for complex sources of variation in an experimental design.

In a mixed model one or more covariates are a categorical variable representing experimental or observational "units" in the data set. [...] If the set of possible levels of the covariate is fixed and reproducible we model the covariate using fixed-effects parameters. If the levels that we observed represent a random sample from the set of all possible levels we incorporate random effects in the model.

A TMT proteomics experiment consists of  $m = 1 \dots M$  concatenations of isobaric-TMT labeled samples or *Mixtures*. Each TMT channel is dedicated to the analysis of  $c = 1 \dots C$  individual biological or treatment *Conditions* prepared from one  $b = 1 \dots B$  biological replicates or *Subjects*. A single mixture may be profiled in  $t = 1 \dots T$  technical replicate mass spectrometry runs.

edgeR: `:glmQLFit`'s documentation).

We prepared 7 subcellular fractions (BioFraction) from 2 Conditions: Control and SWIP<sup>P1019R</sup> Mutant mice. There were 6 Subjects, three bioreplicate Control and SWIP<sup>P1019R</sup> Mutant mice.

	C1	C2	C3	C4	C5	C6	C7	C8	C9	C10	C11	C12	C13	C14	C15	C16
Mix1	WT-5K	WT-9K	WT-12K	WT-15K	WT-30K	WT-79K	WT-129K	QC1	MUT-5K	MUT-9K	MUT-12K	MUT-15K	MUT-30K	MUT-79K	MUT-129K	QC2
Mix2	WT-5K	WT-9K	WT-12K	WT-15K	WT-30K	WT-79K	WT-129K	QC1	MUT-5K	MUT-9K	MUT-12K	MUT-15K	MUT-30K	MUT-79K	MUT-129K	QC2
Mix3	WT-5K	WT-9K	WT-12K	WT-15K	WT-30K	WT-79K	WT-129K	QC1	MUT-5K	MUT-9K	MUT-12K	MUT-15K	MUT-30K	MUT-79K	MUT-129K	QC2

**Figure 2. Experimental Design.** We utilized 16-plex TMT tags to label samples prepared from 6 mice.

	F4	F5	F6	F7	F8	F9	F10	F4	F5	F6	F7	F8	F9	F10	
L1	-1	0	0	0	0	0	0	+1	0	0	0	0	0	0	Mutant.F4-Control.F4
L2	0	-1	0	0	0	0	0	0	+1	0	0	0	0	0	Mutant.F4-Control.F4
L3	0	0	-1	0	0	0	0	0	0	+1	0	0	0	0	Mutant.F4-Control.F4
L4	0	0	0	-1	0	0	0	0	0	0	+1	0	0	0	Mutant.F4-Control.F4
L5	0	0	0	0	-1	0	0	0	0	0	0	+1	0	0	Mutant.F4-Control.F4
L6	0	0	0	0	0	-1	0	0	0	0	0	0	+1	0	Mutant.F4-Control.F4
L7	0	0	0	0	0	0	-1	0	0	0	0	0	0	+1	Mutant.F4-Control.F4
L8	-1/7	-1/7	-1/7	-1/7	-1/7	-1/7	-1/7	+1/7	+1/7	+1/7	+1/7	+1/7	+1/7	+1/7	Mutant.F4-Control.F4

**Figure 3. Contrast matrix indicating 'Intra-BioFraction' and 'Mutant-Control' comparisons.**

In an experiment such as ours with multiple mixtures and biological replicates, but no technical replication of mixture ( $T = 1$ ) MSstatsTMT fits a linear mixed model of the following form to each protein:

Where Mixture is a mixed-effect and quantifies variation between TMT mixtures. Condition is a fixed effect (mean = 0) and in our experiment represents the interaction of terms Genotype and BioFraction.  $\epsilon$  is a random effect representing both biological and technical variation, quantifying any remaining error.

$$Y_{mcbt} = \mu + Mixture_m + TechRep/Mixture_{m/t} + Condition_c + Subject_{mcb} + \epsilon_{mcbt} \quad (1)$$

$$Y_{mcbt} = \mu + \text{Mixture}_m + \text{Condition}_c + \text{Subject}_{mcb} + \epsilon_{mcbt} \quad (2)$$

Where Mixture represents the random-effect of mixture and Condition is a fixed-effect and in our experiment is interaction of Genotype and BioFraction—the 14 combinations of 7 BioFractions fractions from Control and Mutant mice.  $\epsilon_{mcbt}$  is the residual error ( $\sigma^2$ ).

In our experimental design, we made measurements from seven BioFractions from each subject. Thus, we should include the term Subject, representing the 6 individual mice or subjects analyzed in our experiment.

However, in our design Mixture is confounded with the term Subject – in each mixture we analyzed all BioFractions from a single Control and Mutant mouse. Thus we can choose to account for the effect of Mixture or Subject, but not both. Assuming Mixture contributes greater to the variance, we drop the term Subject, and the reduced model is equivalent to ??.

Model based testing of differential abundance between pairs of conditions is assessed through contrast of conditioned means estimated by fitting the parameters of the model by REML to obtain  $\hat{\beta}$ ,  $\sigma^2$  and  $\hat{V}$ .

The degrees of freedom are determined by the Satterthwaite approximation[REF], and the T-statistic for the contrast is taken to be (lmerTest ref):

$$t = \frac{l^T * \hat{\beta}}{\text{sqrt}(l * \sigma^2 * \hat{V} * l^T)} \quad (3)$$

$\sigma^2$  is the error from **Equation ??**.  $l^T$  is a vector specifying a contrast between positive and negative coefficients in the model.

Together, the denominator  $\sqrt{l * \sigma^2 * \hat{V} * l^T}$  is the standard error of the contrast.

```
suppressPackageStartupMessages({
  library(dplyr)
  library(data.table)
})

## load SwipProteomics data
data(swip)
data(gene_map)
data(msstats_prot)
data(alt_contrast)
```

```

data(msstats_contrasts)

## formula to be fit:
fx0 <- formula("Abundance ~ 0 + Condition + (1|Mixture)")

# fit the model
idx <- msstats_prot$Protein == swip
fm <- lmerTest::lmer(fx0, msstats_prot[idx,])

# calculate model statistics
model_summary <- summary(fm,ddf="Satterthwaite")

df <- model_summary$coefficients
df %>% as.data.table(keep.rownames="Coefficient") %>% knitr::kable()

```

Coefficient	Estimate	Std. Error	df	t value	Pr(> t )
ConditionControl.F10	7.619237	0.1213002	17.24812	62.81305	0
ConditionControl.F4	6.711692	0.1213002	17.24812	55.33125	0
ConditionControl.F5	6.946177	0.1213002	17.24812	57.26434	0
ConditionControl.F6	7.240695	0.1213002	17.24812	59.69235	0
ConditionControl.F7	7.321630	0.1213002	17.24812	60.35958	0
ConditionControl.F8	7.129848	0.1213002	17.24812	58.77853	0
ConditionControl.F9	6.954883	0.1213002	17.24812	57.33611	0
ConditionMutant.F10	5.785004	0.1213002	17.24812	47.69163	0
ConditionMutant.F4	5.405091	0.1213002	17.24812	44.55962	0
ConditionMutant.F5	5.568104	0.1213002	17.24812	45.90349	0
ConditionMutant.F6	5.641531	0.1213002	17.24812	46.50883	0
ConditionMutant.F7	5.633054	0.1213002	17.24812	46.43895	0
ConditionMutant.F8	5.493008	0.1213002	17.24812	45.28440	0
ConditionMutant.F9	5.781281	0.1213002	17.24812	47.66093	0

```

# evaluate goodness-of-fit
r2_nakagawa <- r.squaredGLMM.merMod(fm)
r2_nakagawa %>% round(3) %>% knitr::kable()

```

R2m	R2c
0.935	0.949

```

contrast <- msstats_contrasts[1,]
lmerTestContrast(fm,contrast) %>% knitr::kable()

```

Contrast	log2FC	percentControl	Pvalue	Tstatistic	
Mutant.F4-Control.F4	-0.9075446	0.5330916	2.5e-06	-5.986412	0.1516