

# Statistical methods for relative quantification of post-translational modifications in global proteomics experiments

Devon Kohler    Tsung-Heng Tsai    Ting Huang    Erik Verschueren    Trent Hinkle

Meena Choi    Olga Vitek

[Add Author list]

## Abstract

The scientific community widely utilizes mass spectrometry (MS)-based proteomics to quantify the abundance of proteins and their post-translational modifications (PTMs). Experiments targeting PTMs face several specific challenges. These include the low abundance of modified proteo-forms, few representative peptides that span modification sites, and convolution with abundance changes in the overall protein expression. Due to these challenges, a robust approach to estimate relative systematic changes in PTMs should combine information pertaining to PTM sites over several peptides, replicates in multiple conditions, and consider sources of confounding and variation present in the experiment. We propose a general statistical model and workflow that is both reproducible and comprehensive. The model measures PTM and protein abundance by summarizing intensities through Tukey’s median polish method. Then a model based on the family of linear mixed-effects models is fit. This model is automatically adjusted to the specific experimental design. Finally, the PTM abundances are adjusted to remove variance from changes in the overall protein. We implement this model in the free and open-source R package MSstatsPTM.

## Introduction

The signaling mechanisms that allow cells to mount a dynamic and fast response to a multitude of events are primarily facilitated by the modification of proteins at specific residues, acting as molecular on/off switches.[6] [5] Mass spectrometry-based label-free proteomics is broadly established as the tool-of-choice for unbiased and large-scale identification and quantification of proteins and their post-translational modifications (PTMs) using liquid chromatography coupled with mass spectrometry (LC-MS)[9] [13]. Studies targeting the post-translationally modified proteome focus either on the accurate localization of modification sites on proteins, relative or absolute quantification of a modification site’s occupancy repertoire, or relative changes in occupancy across experimental conditions.[11] Regardless of the question at hand, interrogating the modified proteome is challenging due to a number of reasons. First, the relatively lower abundance of modified proteo-forms dictates that a global interrogation can only be achieved through large-scale enrichment protocols with modification-specific antibodies or beads. Variability in the enrichment efficiency inevitably affects the reproducibility of the number of spectral features (e.g., peptide precursor ions or their fragments) and their intensities, which imposes challenges in both quantification and statistical modeling. Second, contrary to the often large number of identified peptides that can be used as features to model protein abundance changes, there are relatively few representative peptides that span a modification sites, which often results in sparse, and sometimes, inherently convoluted models (i.e., single versus multiple mod-

ified sites on a single peptide). Third, unless early signaling events are interrogated, the interpretation of the relative changes in modification occupancy are inherently convoluted with changes in the overall protein expression, making the interpretation of the results not straightforward. [12] Therefore, a robust approach to estimate systematic relative changes in post-translational modifications, at scale, should not only combine the quantitative information pertaining to a PTM site over peptides and replicates in multiple conditions, but take into account various sources of variations and confounding factors present in the experiments.

Despite the important implications of PTMs in biological functions, there is a lack of general framework to summarize the available quantitative information from LC-MS data, to perform statistical inference, and to draw conclusions to characterize the quantitative properties of PTM in a statistically rigorous manner. Many investigations performed differential expression analysis of PTMs using two-sample t-test or its extensions. [Add Ref about current PTM approaches] The approach takes as input intensities of individual features from modified peptides, or intensity ratios of modified and unmodified peptide features, and compares the mean abundance of a PTM site from one condition to another. Modifications of the t-test such as moderated t-test with limma were also proposed.[16] While simple, the approach does not fully account for the sources of variations, and it is not directly applicable to experiments with complex designs, e.g., comparisons of multiple conditions, acquisition in multiple batches, etc. Isobar-PTM was developed for experiments with MS/MS quantitative strategies that employ isobaric labels such as tandem mass tags (TMT) and isobaric tag for relative and absolute quantification (iTRAQ).[3] Isobar-PTM expresses MS measurements with a linear model and performs adjustment with respect to protein abundance using the difference between log-ratio of modified peptides in two channels and log-ratio of protein level. The modeling framework, however, is not applicable for either label-free workflows or experiments with complex designs.

We propose a general statistical approach, which explicitly characterizes the variations and confounding factors present in bottom-up PTM experiments. The proposed approach is aimed at the detection of quantitative changes in PTMs between conditions utilizing procedures developed for summarization of LC-MS data, quantitative characterization of site-specific PTMs, and adjustment with respect to protein abundance. Quantitative analyses of PTMs often involve comparisons between multiple inter-related conditions of the same biological system. The general statistical framework underlying the proposed approach allows for analyzing experiments with complex designs, including those with multiple conditions, multiple batches, and measured at different points in time.

The proposed approach was evaluated using datasets from computer simulations, benchmark controlled mixtures, and biological investigations. The approach was then compared against the commonly applied t-test and limma package. The results demonstrate that by appropriately leveraging the information from the entire dataset the proposed approach improves the reproducibility and accuracy of the estimates of PTM fold changes, results in a better calibrated type I error rate, and improves the statistical power of detecting changes in PTMs. The proposed approach is implemented as an open source R package MSstatsPTM, which employs similar input format as in MSstats and MSstatsTMT [4] [8].

## Proposed Approach

Figure 6.3 schematically illustrates a simplified version of the data structure resulting from a typical bottom-up experiment for quantitative analysis of PTMs, in which there are multiple layers of variation present. A PTM site is quantified with multiple spectral features, which vary in sequence (e.g., fully or partially cleaved peptides), ionization efficiency, charge states, etc. The number quantified features vary across replicate LC-MS/MS runs of the same sample, and across conditions. To perform adjustment with respect to protein abundance, features of unmodified peptides are used for the inference of underlying protein abundance. Typically, because of the enrichment step for PTMs, very few of those features are present in original LC-MS runs. For more accurate estimation of protein abundance, separate global proteomics data of unenriched samples are often acquired. If unmodified features are unavailable for any given modified feature, unmodified intensity adjustment cannot be performed. As different levels of variability are present in the data, the log-intensities of the features for modified and unmodified peptides are modeled separately using two linear mixed models.

## Statistical modeling and parameter estimation

The proposed approach takes as input a list of log-transformed intensities of spectral features, identified and quantified across LC-MS runs. The features, which are precursor ions of modified or unmodified peptides, are used to characterize the identified PTM sites and proteins. For each PTM site, the feature log-intensities of the modified peptides spanning the site are expressed using a linear mixed model in consideration of the effects of condition, run, feature and interaction between run and feature. In addition mixture may be included for features acquired via tandem mass tag (TMT) methods. The model parameters are estimated using the split-plot approach as in MSstats, where the feature log-intensities are first summarized into a single value per site per run in the subplot model, and the site-level summaries are then used for the inference of the PTM site abundance.[4] In the site-level summarization, Tukey’s median polish (TMP), a simple and robust procedure is applied to iteratively fit a two-way additive model with the effects of run and feature, which in turn summarizes the log-intensities for each site.[15] After summarization, the inference of the PTM site abundance in each condition is carried through fitting a model based on the family of linear mixed-effects models, taking into account the specific experimental design 6.3.[2] [7] Statistical modeling and quantification for global proteomics data are performed by the same procedure as for PTM data.

## Detection of changes in PTMs

Detection of differentially modified PTM sites is performed through testing the null hypothesis of ‘no change’ against the alternative. The null hypothesis states that there is no difference in log-abundance of the PTM site between conditions, adjusted with respect to protein abundance 6.3. Specifically, the adjusted difference is given by the difference in log-abundance of the PTM site, subtracted by the difference in log-abundance of the underlying protein, which is equivalent to the log of the ratio of PTM abundance difference to protein abundance difference. The estimate of the adjusted difference and the standard error (SE) of the estimate are obtained by combining the difference estimates and the associated SEs from both counterparts.

The test statistic for the hypothesis testing is the ratio of the estimate of the adjusted difference to its SE. To determine the statistical significance of the difference in terms of p-value, the test statistic is compared against the t distribution with degrees of freedom approximated by the Satterthwaite method.[14] Adjustment for multiple comparisons is performed using the Benjamini-Hochberg procedure to control the false discovery rate at a desired level, e.g., 0.05.[1] The details are provided in Supplementary Sec. S3.2 [4Devon: [link to supplementary using sfigref](#)].

## Analysis with multiple batches

The proposed statistical framework allows to analyze data from experiments of complex designs such as factorial design. We discuss below a specific design commonly considered in PTM experiments, in which data are acquired in multiple batches, often as a result of multiple repeats for the antibody enrichment step. Typically, the number of spectral features and their intensities are different across batches, and the run-to-run variation sometimes varies across batches as well. The proposed approach summarizes log-intensities in each batch separately, and considers two ways to perform the statistical inference and testing. When variability in site-level summaries and in PTM changes is expected across batches, the proposed approach performs inference of the PTM site abundance and the adjusted difference for each batch separately. The evidence about differential PTM sites between conditions is averaged over batches. Alternatively, assuming identical variability in the site-level summaries and identical differences between conditions across batches (i.e., no interaction effect between condition and batch), the statistical inference and testing can be performed based on one single model with an additional fixed effect of batch. Details about these modeling considerations are discussed in Supplementary [\[link to Sec. S3.4\]](#).

## Missing Value imputation

MSstatsPTM gives the ability to impute missing feature intensities if desired. When values are imputed, it is assumed they are missing for reasons of low abundance. Missing features are imputed in each MS run using the Accelerated Time Failure model.[15] In order to impute a feature’s missing values, the feature must be present in at least one MS Run. If the feature is not present across all MS Runs, the value will be left missing. Missing value imputation is done before summarization with Tukey’s Median Polish in order to correct for the impact of outliers. Missing value imputation is done separately for both the PTM and global protein datasets.

## Extension to TMT experiments

The statistical modeling approaches discussed above can also be extended to Tandem Mass Tag (TMT) labeling methods. TMT experiments introduce an additional source of variation in the form of different mixtures. To account for the mixture variable, first feature intensities are summarized as described previously, then a new linear mixed effects model with an added term for the mixture is included. This workflow and model follow the methods used for MSstatsTMT. [8] These methods are repeated to quantify and model global protein abundance. Once both PTM and global protein are modeled, the PTM model is adjusted for changes in global protein abundance using the same methods are before.

[add figref for TMT model] [add figref for TMT experiment table]

## Results

### Computer simulations

The proposed statistical approach methods were evaluated using computer simulation. Specifically, their properties under adjustment with respect to protein abundance and batch effects were evaluated. Two simulations were ran, one including a high number of features and no missing values, and one with few modified features and including missing values.

**Computer simulation: protein-level adjustment.** Differential intensity levels of modified peptides may be due to changes in modification, change in protein abundance, or both. The proposed approach adjusts the abundance with respect to unmodified peptides by combining the inference of modified and unmodified peptide abundances. Alternatively, two-sample t-test or limma that takes as input the ratio between modified and unmodified peptide intensities (difference on log scale) is commonly applied for the same purpose. In real experiments, multiple inter-related conditions are often compared together. Whereas t-test uses measurements from the two conditions being compared, the proposed approach and limma leverages measurements in all conditions for the inference of underlying abundance. We evaluated the impact with such adjustment by computer simulation, considering the following factors: with/without protein-level changes, effect size, number of replicates and number of conditions. Details of the simulation can be found in Supplementary supp-sec:sim.

In this simulation, when there was no change in protein abundance, all the considered approaches well calibrated the Type I error rate [figref Fig. 4a]. However, when the modification changes were entirely due to the changes in protein abundance across conditions, analysis without accounting for the protein-level changes resulted in off-target, high false positive rates [figref (Fig. 4b)]. As shown in [figref Fig. 4c-e], in detection of systematic changes in PTM abundance, the proposed approach improved statistical power with small sample sizes in almost all the considered scenarios. The advantage by using the proposed approach over t-test became more profound in the presence of missing data [figref (Fig. 4d-e)]. In cases with small sample sizes (e.g., 2 replicates) as shown in [figref Fig. 4d], performance by t-test decreased dramatically with one missing value, and increasing the log-fold change to 2 did not effectively reduce the negative impact. Two-sample t-test only used data within the groups of interest while ignoring the rest of the data. Consequently, it gave similar performance across cases with different number of conditions. In contrast, the proposed

approach leveraged all available information, which resulted in improved power with increased number of conditions (see for example, the improved performance in [\[figref Fig. 4e\]](#) [with 4 conditions] over [\[figref Fig. 4d\]](#) [with 2 conditions]). Further results are provided in Supplementary [\[link to Sec. S4.1\]](#).

## Spike-in experiment

[\[Need to write results of this experiment\]](#)

We evaluated our approach using a custom designed spike-in benchmark experiment where 50 heavy-labeled KGG motif peptides from 20 human proteins were used as spike-in peptides. Quantitative changes in protein and site abundance changes of these 20 proteins were the target of the benchmark. Unmodified peptides from Human Lysate were used as the estimate of global protein abundance changes. All comparisons with respect to human lysate serve as a null model, there was no change in global protein nor PTM abundance between conditions. Additionally, E coli Lysate was used to normalize total protein levels prior to enrichment or global protein profiling. Four mixes of spike in peptides and Human Lysate were created to create conditions with known fold change. [\[Add figure of mixtures\]](#) Two sets of data were acquired for each mixture: modified peptide data, including the spike in KGG enriched peptides, and unmodified peptides from Human Lysate. These datasets were used as the PTM and global protein data.

[\[figref volcano plot of spike in peptides\]](#)

In *\*FigXX\** we can clearly see the red labeled spike-in peptides do not follow the expected fold change before adjusting for changes in global protein level. After adjustment the estimated fold change is much close to expectation. Additionally, background peptides that should not show a significant change, colored grey, show many false positives before adjustment is made. Again after adjustment the results improve, and the number of false positives decrease significantly. In this experiment adjusting to remove confounding with changes in global protein abundance leads to a more accurate conclusion about the effect of modified peptides between conditions.

## Dilution experiment [What was this experiment?]

We evaluated our approach using a custom designed dilution benchmark experiment where two sets of liver samples from five Atg16L1 deficient mice (one with KGG enrichment followed by trypsin digestion, one with trypsin digestion alone) were prepared. Following a randomized order, each sample was injected at three different concentrations 4ul, 2ul, and 1ul. The experiment resulted in 15 KGG and 15 global profiling runs. To correct for run-to-run variation, the AQUA peptide mixture was spiked into each sample. There are 4696 ubiquitinated proteins identified in this dataset. Among these, 3173 have corresponding measurements from global data. To simulate the effects of changes in protein abundance, we created three groups of data by matching one concentration for the KGG data with one for the global data: K1P1, K2P2, and K4P4. Systematic changes in the abundance of ubiquitinated sites between the groups are present in the KGG data. These changes are considered as artifacts driven by the changes in the abundance of their corresponding proteins. Using the comparison of K2P2 vs. K1P1 as example, there are 19355 possible comparisons with KGG data, and 3860 changes are detected. Out of the 19355 sites, 15036 have corresponding measurements from global data and are eligible for protein-level correction. The correction significantly reduced the number of false positives to 79. Similar observations were made in the other two comparisons K4P4 vs. K2P2 and K4P4 vs. K1P1, as shown in Table XX. [\[Create or find this table\]](#)

## Re-analysis of published dataset: TMT experiment

To evaluate our approach on TMT experiments we reanalyzed an experiment targeting primary murine macrophages infected with *Shigella flexneri* (*S.flexneri*). [10] Tandem mass tagging was used to quantify changes in total proteins, phosphorylation, and ubiquitination in Wild type (WT) controls and ATG16L1-deficient (cKO) bone marrow macrophages (BMDMs) either infected or uninfected with *Shigella flexneri*. 11-plex isobaric multiplex was used with two separate mixtures. Cell lysates were prepared over three different time periods, either uninfected, infected at early (45-60 min) or infected at late (3-3.5 hour) time

points. There were 103,700 peptides mapping to 9,430 proteins in the global profiling run. Additionally, there were 25,600 unique phosphorylation sites and 12,400 KGG ubiquitination sites. About 90% of the identified modified peptides derived from proteins that were also quantified in the global profiling run.

The major results are summarized below.

## Conclusion and discussion

We proposed a general statistical modeling framework for PTM characterization. The framework is designed for bottom-up MS workflows, which are characterized with variations from multiple convoluted sources, frequent missing data, and associated uncertainty in the conclusions. The framework is general and is applicable to a variety of experimental designs. It outperforms the ad-hoc methods underlying the t-test, and yields accurate results in the broad type of experimental circumstances, including the presence of missing values, changes in protein abundance, batch effects, and different acquisition methods. The framework allows us to plan for subsequent experiments, and choose the appropriate number of replicates in consideration of adjustment with respect to protein abundance.

Our results show that when measurements from multiple related conditions are available, the proposed approach for joint modeling and summarization of all the LC-MS/MS runs leads to more sensitive PTM significance analysis and more accurate and precise quantification than when separately analyzing conditions of runs. The gain is due to a more efficient use of the data, and to a more accurate understanding of the systematic and random variations. The proposed framework can be extended beyond the experimental designs with multiple batches and conditions discussed above. For example, it can represent experimental designs with even more complex structures, such as time series or factorial investigations.

A potential limitation of the proposed framework is the assumption that all the peptides are correctly mapped to the underlying proteins and PTM sites, and the features are informative of the abundances of underlying protein and PTM. Also, characterizing PTMs with current data-dependent acquisition workflows is prone to being under sampled, leading to a sparse dataset with a large number of missing values for the analysis. Statistical methods accounting for effects due to experimental units and missing values introduced in this manuscript help interpret the data in a more objective manner. The latest development of targeted acquisition and data-independent acquisition methods are expected to further alleviate these issues.

Overall, the proposed approach balances accuracy and practicality, and enables the analysis of complex experiments in high throughput. Future work is to carry out the inference and testing for not only the relative change of PTM abundance, but also the fraction of the protein that is modified at the particular site (site occupancy, or stoichiometry). We are also interested in characterizing the interplay of PTMs at multiple sites. The proposed statistical methods are implemented as an R package MSstatsPTM available on Bioconductor and Github.

## References

- [1] Y. Benjamini and Y. Hochberg. “Controlling the false discovery rate: a practical and powerful approach to multiple testing”. In: *J.R. Statist. Soc. B* 57.1 (1995), pp. 289–300.
- [2] Benjamin M. Bolker et al. “Generalized linear mixed models: a practical guide for ecology and evolution”. In: *Trends in Ecology and Evolution* 24.3 (2009), pp. 127–135. ISSN: 01695347. DOI: 10.1016/j.tree.2008.10.008.
- [3] Florian P. Breitwieser and Jacques Colinge. “IsobarPTM: A software tool for the quantitative analysis of post-translationally modified proteins”. In: *Journal of Proteomics* 90 (2013), pp. 77–84. DOI: <https://doi.org/10.1016/j.jprot.2013.02.022>. URL: <https://www.sciencedirect.com/science/article/pii/S1874391913000973>.
- [4] M. Choi et al. “MSstats: an R package for statistical analysis of quantitative mass spectrometry-based proteomic experiments”. In: *Bioinformatics* 30 (2014), pp. 2524–2536.
- [5] P. Cohen. “The regulation of protein function by multisite phosphorylation—a 25 year update.” In: *Trends Biochem Sci.* 25 (2000), pp. 596–601.
- [6] Y.L. Deribe, T. Pawson, and I. Dikic. “Post-translational modifications in signal integration”. In: *Nature Structural & Molecular Biology* 17 (2010), pp. 666–672.
- [7] J. J. Faraway. *Extending the linear model with R*. 1st. Boca Raton, FL: Taylor & Francis Group, LLC, 2006.
- [8] T. Huang et al. “MSstatsTMT: Statistical Detection of Differentially Abundant Proteins in Experiments with Isobaric Labeling and Multiple Mixtures”. In: *Molecular & Cellular Proteomics* 19 (10 Oct. 2020), pp. 1706–1723.
- [9] L. Käll and O. Vitek. “Computational mass spectrometry-based proteomics”. In: *PLoS Comput. Biol.* 7 (12 Dec. 2011), e1002277.
- [10] Timurs Maculins et al. “Proteomics of autophagy deficient macrophages reveals enhanced antimicrobial immunity via the oxidative stress response”. In: *bioRxiv* (2020). DOI: 10.1101/2020.09.10.291344. eprint: <https://www.biorxiv.org/content/early/2020/09/12/2020.09.10.291344.full.pdf>. URL: <https://www.biorxiv.org/content/early/2020/09/12/2020.09.10.291344>.
- [11] M. Mann and O. Jensen. “Proteomic analysis of post-translational modifications”. In: *Nat Biotechnol* 21 (2003), pp. 255–261.
- [12] J. Olsen and M. Mann. “Status of Large-scale Analysis of Post-translational Modifications by Mass Spectrometry”. In: *Molecular & Cellular Proteomics* 12.12 (2013), pp. 3444–3452.
- [13] Roepstorff P. “Mass spectrometry in protein studies from genome to function.” In: *Curr Opin Biotechnol.* 8.1 (1997), pp. 6–13.
- [14] Franklin E Satterthwaite. “An approximate distribution of estimates of variance components”. In: *Biometrics bulletin* 2.6 (1946), pp. 110–114.
- [15] J. W. Tukey. *Exploratory data analysis*. Addison-Wesley, 1977.
- [16] Y. Zhu et al. “DEqMS: A Method for Accurate Variance Estimation in Differential Protein Expression Analysis”. In: *Molecular & cellular proteomics : MCP* 19 (2020), pp. 1047–1057.



## Figure

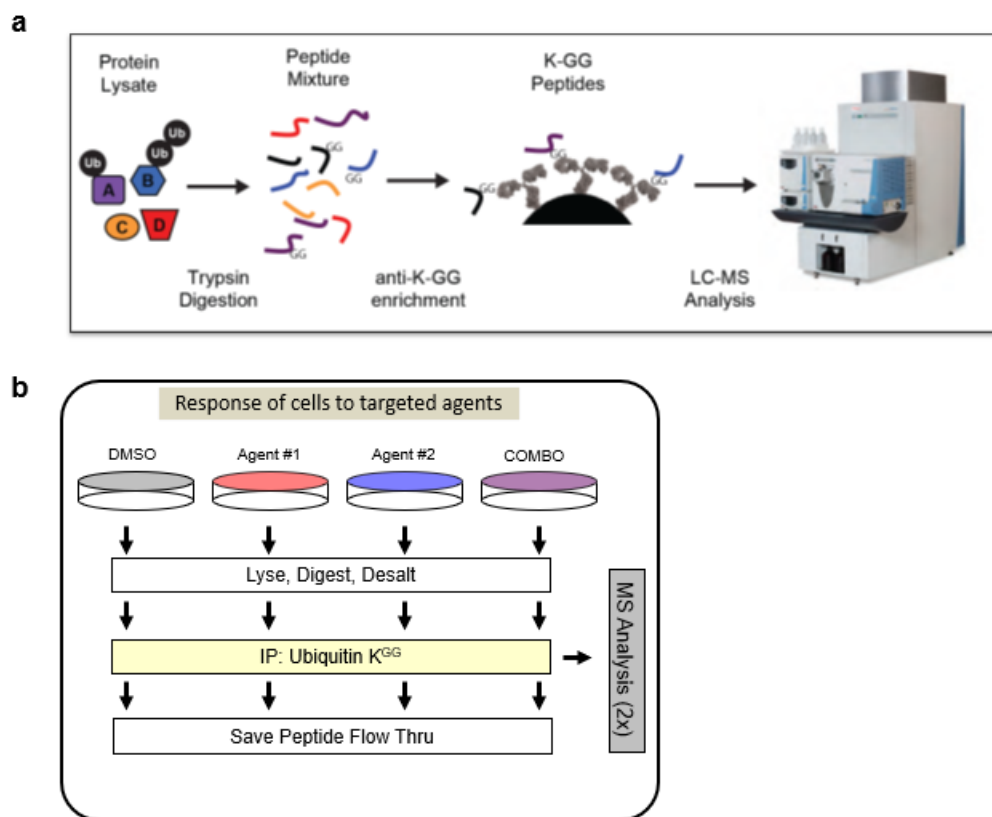


Figure 6.1: [todo]

**Figure 4** False positive rate and statistical power in PTM significance analysis by the proposed approach, t-test with protein-level adjustment, and t-test without protein-level adjustment. (a) When there was no change in protein abundance, all considered methods well calibrated the Type I error rate. (b) When the PTM changes were entirely due to changes in protein abundance across conditions, analysis accounting for the protein-level changes resulted in off-target, high false positive rates. (c) Statistical powers by the proposed approach and t-test with protein-level adjustment, where the data consisted of 2 conditions and the SD corresponding unmodified feature intensities was 0.2. (d) Statistical powers by the proposed approach and t-test with protein-level adjustment, where the data consisted of 2 conditions, the SD corresponding unmodified feature intensities was 0.2, and the PTM was missing in Run 1 of Condition 1. (e) Same as in (d), but the data consisted of 4 conditions.

**Figure 5** Results corresponding to estimation error and false positive rate and statistical power of the PTM significance analysis, where the data were acquired in two batches. The following parameters were considered to represent the batch effects: no batch-condition interaction, mean intensity level in Batch 2 was higher than Batch 1 by 2 on log scale, and the SDs in Batch 1 and Batch 2 were 0.2 and 0.3, respectively. (a) Estimation based on the most statistically significant batch with t-test was more variable than other methods and frequently biased. (b) The proposed approach better calibrated Type I error rate. (c) The proposed approach improved statistical power in all the considered cases with different number of conditions and replicates. This was achieved by properly characterizing batch effects and leveraging all available information. Ignoring batch effects by t-test (no batch) lost power dramatically. The negative impact was only partially reduced by increasing the sample size to 5.



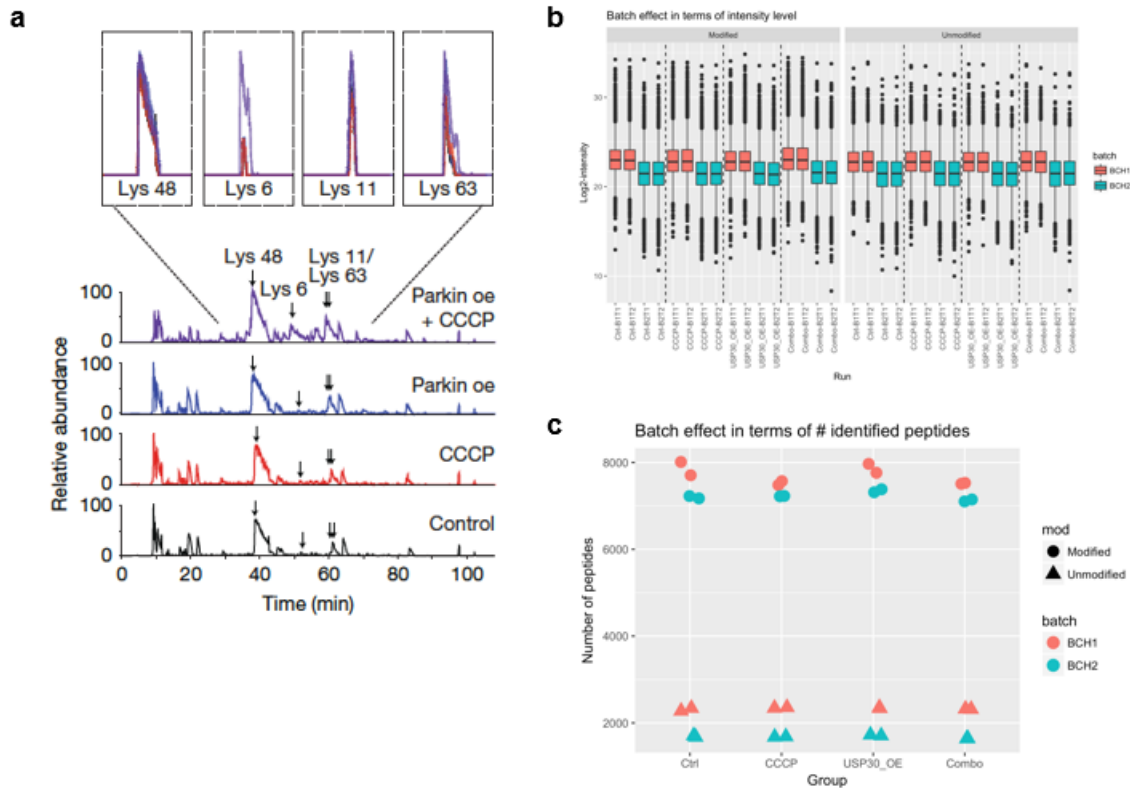


Figure 6.2: [todo]

**Figure 6** Design of future PTM experiments in terms of sample size calculations and power analysis. (a) Protein-level adjustment relies on the inference of protein abundance, which introduces additional uncertainty in the estimate of PTM difference. Therefore, the required sample size to detect a systematic change is higher than as expected for standard differential analysis without adjustment. Sample size calculations without accounting for the uncertainty would lead to over-optimistic, under-powered studies. (b) In complex designs, simultaneously analyzing all the conditions effectively increases the degrees of freedom and requires fewer replicates. (c) Increasing the sample size and analyzing multiple conditions together both result in improved statistical power.

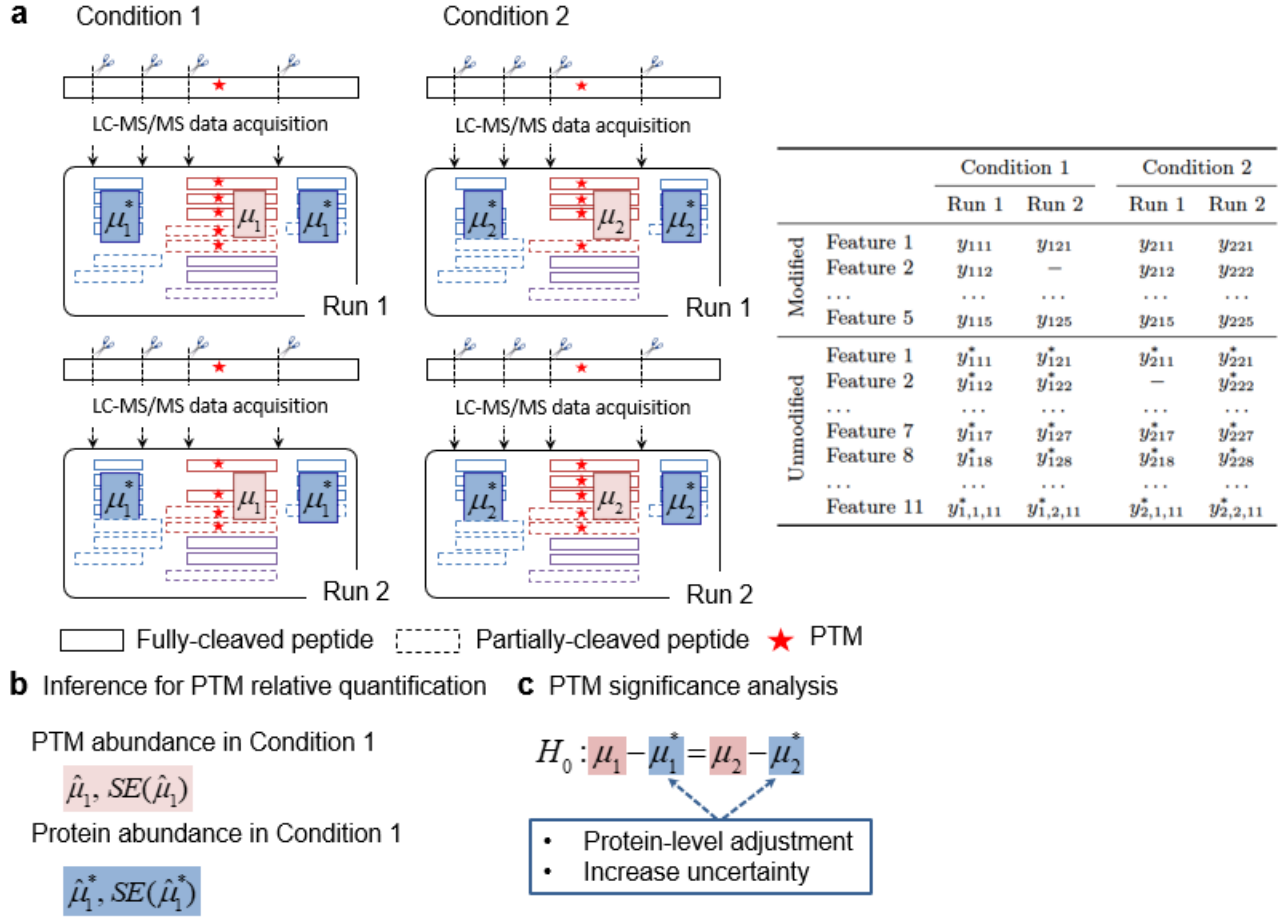


Figure 6.3: Data structure of a typical PTM experiment and goals of PTM characterization. (a) Schematic data representation, in a simplified case of two conditions and two replicate runs. Each PTM site is modeled and characterized separately, where a PTM is quantified with multiple spectral features (boxes), distinguished by different charge states of a peptide. The feature intensities are viewed as repeated measurements of the underlying abundance of the PTM, where the abundance in Condition  $i$  is denoted by  $i$ . Features corresponding to unmodified peptides are considered together to perform adjustment with respect to protein abundance, where the protein abundance in Condition  $i$  is denoted by  $i^*$ . Peptides can be fully cleaved (solid lines) and/or partially cleaved (dashed lines). Some spectral features can be missing. (b) PTM relative quantification by statistical inference, which makes use of the feature intensities to infer the underlying PTM abundance and protein abundance with an estimate of associated uncertainty. (c) Model-based testing for differential PTM abundance, which corrects for the underlying protein abundance with a cost of increased uncertainty about the estimate of difference between conditions. [there is no d.] (d) Statistical experimental design in terms of sample size calculations and power analysis.