

# PTM-Shepherd: Analysis and Summarization of Post-Translational and Chemical Modifications From Open Search Results

## Authors

Daniel J. Geiszler, Andy T. Kong, Dmitry M. Avtonomov, Fengchao Yu, Felipe da Veiga Leprevost, and Alexey I. Nesvizhskii

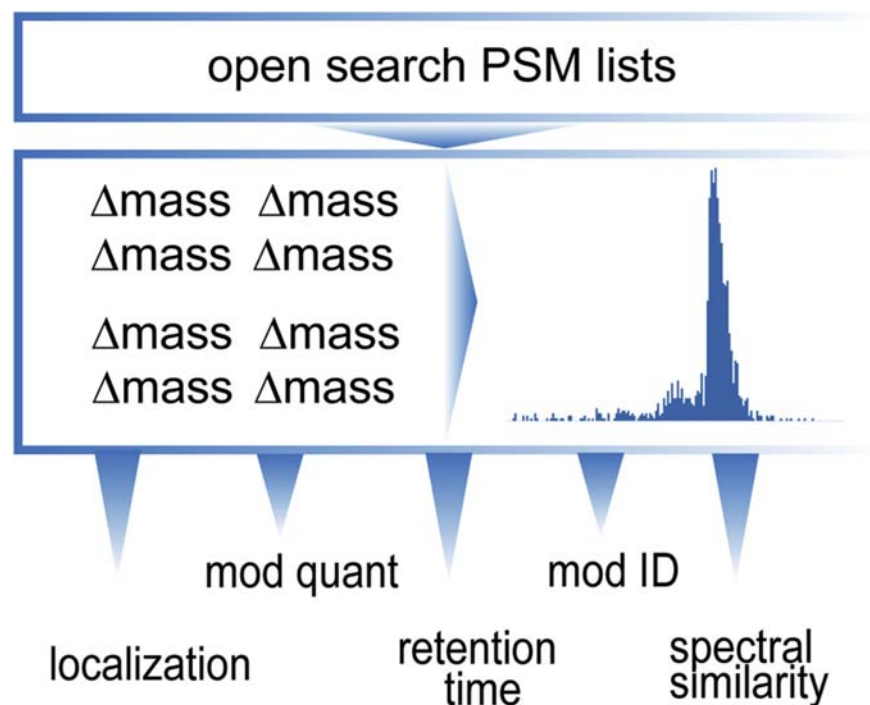
## Correspondence

[nesvi@med.umich.edu](mailto:nesvi@med.umich.edu)

## In Brief

Proteomics open searches require extensive post hoc analysis to be useful. PTM-Shepherd provides comprehensive open search annotation from peptide-spectrum match lists, including modification ID, localization, and effects on spectral similarity (SS) and retention time (RT). These features can be used to find batch effect and new post-translational modifications, and inform subsequent closed searches. It is available as a standalone JAR and as part of the FragPipe suite.

## Graphical Abstract



## Highlights

- Comprehensive open-search annotation.
- Sensitive modification detection.
- Identification of batch effects.
- Novel post-translational modification discovery.



# PTM-Shepherd: Analysis and Summarization of Post-Translational and Chemical Modifications From Open Search Results

Daniel J. Geiszler<sup>1</sup>, Andy T. Kong<sup>2</sup>, Dmitry M. Avtonomov<sup>2</sup>, Fengchao Yu<sup>2</sup>, Felipe da Veiga Leprevost<sup>2</sup>, and Alexey I. Nesvizhskii<sup>1,2,\*</sup>

**Open searching has proven to be an effective strategy for identifying both known and unknown modifications in shotgun proteomics experiments. Rather than being limited to a small set of user-specified modifications, open searches identify peptides with any mass shift that may correspond to a single modification or a combination of several modifications. Here we present PTM-Shepherd, a bioinformatics tool that automates characterization of post-translational modification profiles detected in open searches based on attributes, such as amino acid localization, fragmentation spectra similarity, retention time shifts, and relative modification rates. PTM-Shepherd can also perform multiexperiment comparisons for studying changes in modification profiles, e.g., in data generated in different laboratories or under different conditions. We demonstrate how PTM-Shepherd improves the analysis of data from formalin-fixed and paraffin-embedded samples, detects extreme underalkylation of cysteine in some data sets, discovers an artifactual modification introduced during peptide synthesis, and uncovers site-specific biases in sample preparation artifacts in a multicenter proteomics profiling study.**

Database searching of shotgun proteomics data is a commonly used strategy for identification of peptides and proteins from complex protein mixtures (1, 2). Peptide identification in this strategy most commonly relies on matching tandem mass spectrometry (MS/MS)-derived peptide spectra to their theoretical counterparts using MS/MS database search tools, which require prior knowledge of the potential modifications that might be present in a sample. This is problematic, as proteins can exist in myriad forms outside their canonical sequences. For example, protein function is commonly modulated by post-translational modifications (PTMs), and additional chemical modifications from sample processing can hinder identification. Because the search space of all potential peptides including their modifications is so large, when using conventional database search strategies,

researchers are forced to limit the modifications considered by their searches, leading to large number of unexplained spectra (3–6).

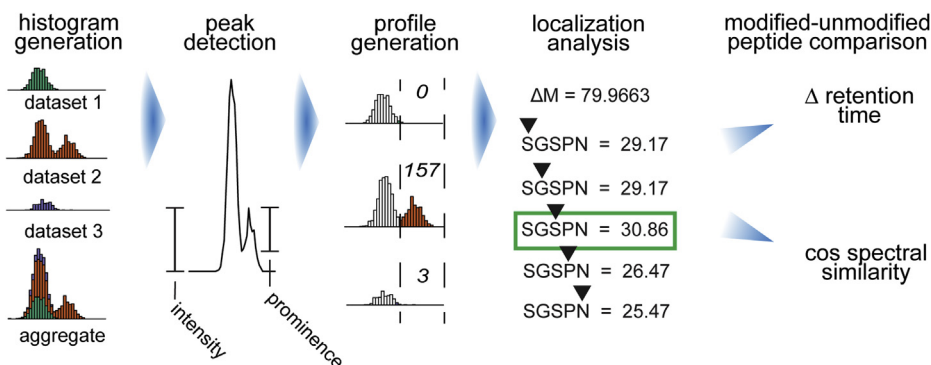
Open searching, or mass-tolerant searching, is one strategy that allows researchers to expand their search space and reduce the number of unexplained MS/MS spectra. It has proven to be an effective strategy for identifying both known and unknown modifications in shotgun proteomics experiments (3, 4, 7–9). Rather than being limited to user-specified modifications, open searches identify peptides with mass shifts corresponding to potential modifications or sequence variants. These mass shifts do not, however, contain the same information present in closed searches, most importantly the identity of the modification and what amino acids within the peptide sequence may contain it. Deciphering open search results thus requires subsequent computational characterization to recover this information (4, 10–12).

Existing tools for open search postprocessing perform a limited set of analyses on a spectrum-level basis. PTM-Prophet (13), e.g., is limited to localizing mass differences for each peptide-spectrum match (PSM) but neither does provide data summaries that can inform subsequent searches nor does it provide identities for mass differences. Philosopher (14) only provides mappings of mass differences to UniMod and generates a basic mass shift histogram. Here we present PTM-Shepherd, an automated tool that calls modifications from open search PSM lists and characterizes them based on attributes, such as amino acid localization, fragmentation spectra similarity, effect on RT, and relative modification rates. PTM-Shepherd can also perform multiexperiment comparisons for studying changes in modification profiles under differing conditions. We use these profiles in a wide array of situations to show how additional metrics, interexperiment comparisons, and bulk analytical profiles can be helpful in PTM analysis. Overall, we expect that PTM profiles produced by PTM-Shepherd will greatly enhance understanding of the

From the <sup>1</sup>Department of Computational Medicine and Bioinformatics and <sup>2</sup>Department of Pathology, University of Michigan, Ann Arbor, Michigan, USA

This article contains [supporting information](#).

\*For correspondence: Alexey I. Nesvizhskii, [nesvi@med.umich.edu](mailto:nesvi@med.umich.edu).



**FIG. 1. PTM-Shepherd workflow.** Data processing begins by aggregating the mass shifts across all data sets into a common histogram. Peaks are determined based on their prominence. The 500 most intense peaks in aggregate are then quantified for each data set and normalized to size. Peptides with each mass shift are iteratively rescored with the modification at each position, producing localization scores for each peptide and an aggregate localization enrichment for each mass shift. Finally, modified peptides and their unmodified counterparts are analyzed to have their pairwise cosine spectral similarity and change in retention time calculated.

data at both the macro level for quality control (QC) and the micro level for specific PTM identification.

#### EXPERIMENTAL PROCEDURES

##### *PTM-Shepherd: Mass Shift Histogram Construction*

A histogram of identified mass shifts is constructed using all PSMs from the PSM.tsv file (or multiple PSM.tsv files in the case of multi-experiment analysis) generated by the MSFragger/Philosopher pipeline (Fig. 1). These PSM.tsv files are typically (by default) filtered to 1% PSM level and 1% protein level false discovery rate (FDR) using target-decoy counts, as determined by the Philosopher filter command. The width of each bin in the histogram is 0.0002 Da (by default). This histogram is extended by 5 Da on either side of the most extreme values to prevent peaks at the maximum and minimum of the histogram from being truncated after smoothing. Random noise between  $-0.005$  and  $0.005$  Da is added to break ties occurring between bin boundaries and mass shifts. After bin assignment, the histogram is smoothed to make peaks more monotonic. Bin weight is distributed across five bins (by default), with the weights assigned to each bin being determined by a Gaussian distribution centered at the bin to be smoothed such that 95% of the bin's weight is distributed between them. Peaks, representing mass shifts of observed modifications, are called from this histogram.

##### *PTM-Shepherd: Peak Picking*

PTM-Shepherd picks peaks based on a mixture of peak prominence and signal-to-noise remainder (SNR) as measures of quality and quantification, respectively. A peak's prominence is calculated as the ratio of its apex to the more intense of either its left or right shoulder, found by following a peak downward monotonically (Fig. 1). To improve monotonicity for this procedure, adjacent histogram bins are temporarily grouped into small sets and flattened to the minimum bin height within the set, with set size internally calculated based on the total number of histogram bins. Peaks are called when their prominence exceeds 0.3 (by default). A peak's SNR is calculated with a 0.004 Da sliding window (by default) against a background of 0.005 Da on either side (scaling linearly with peak picking width). The average height per histogram bin is computed for the signal and noise regions, and then the signal remainder is calculated by subtracting off the noise. From this list of peaks, the top 500 by SNR (by default) are sent to downstream processing. Peak boundaries are considered to be

either the observed peak boundary or the defined precursor tolerance, whichever is closer to the apex. PSMs are assigned to the peak if their mass shift falls within the reported peak boundary.

##### *PTM-Shepherd: Mass Shift Annotation*

Detected peaks are iteratively annotated using entries from the Unimod (retrieved: October 2, 2019) (15) modification database (including single-residue insertions and deletions and isotopic error peaks), supplemented with a user-specified list of mass shifts. Each peak is allowed to be decomposed into at most two modifications. Some exceptionally rare or protocol-based modifications (e.g., O18 labeling, N15 labeling) that regularly confounded annotation were removed. Mass differences within 0.01 Da (by default) of a known mass shift are annotated immediately. If a mass shift does not meet this condition, it is then tested against combinations of user-defined mass shifts and known annotations before being checked against combinations of two modifications identified at the previous step. Failing both these assignments, mass differences are marked as "unannotated" and appended to the list of potential modification combinations.

##### *PTM-Shepherd: Mass Shift Localization*

PTM-Shepherd constructs localization profiles for each mass shift peak. Localization profiles are constructed for each experiment, reporting an N-terminal localization rate and a normalized amino acid propensity for each peak. The localization step is performed for every PSM by placing the mass shift at each amino acid in turn and rescoring the PSM (with the original spectrum) using the same scoring function as in MSFragger. PSMs are considered localizable if there is a position(s) within the peptide sequence that, when the mass shift is placed there, results in more matched fragment ions than using unshifted fragment ions only (i.e., without adding the mass shift anywhere). Localizable PSMs corresponding to the same peak in the mass shift histogram are aggregated, and their characteristics are analyzed. The localization rate for a peak is calculated by counting the number of instances a mass shift was localized to a particular amino acid. If the localization is ambiguous (i.e., several sites scored equally high), the weight of the localization is distributed among all localized residues. Counts are normalized to the rate of localization for a given residue, and then divided by each residue's background content. Background residue content is computed by counting the number of occurrences of each residue in every localizable PSM in the entire data set (by default). Options for experiment-level normalization at the

unique peptide level and bin-wise normalization at the PSM and unique peptide level are also available.

#### Modified–Unmodified Comparisons

Cosine SS between modified and unmodified peptides is used to determine how mass shifts affect MS/MS spectra. Unmodified PSMs, *i.e.*, PSMs with a mass shift less than 0.001 Da (by default), are aggregated based on their identified peptide sequence and charge state. If there are more than 50 unmodified spectra for a peptide, 50 are randomly selected for downstream comparisons. Then, for every mass-shifted PSM at a given charge state, the average cosine similarity score between this PSM and its corresponding unmodified PSMs at the same charge state is recorded. These similarity scores are aggregated for all PSMs for each mass shift peak, then averaged and reported as that peak's SS profile. RT effects are also examined. Peptide RTs are extracted from Philosopher's PSM.tsv output. For every PSM with a mass shift, the average difference in RT between that PSM and all its corresponding unmodified PSM is calculated. These average RT differences are aggregated for each mass shift peak, then averaged across all peptides in that peak and reported as that peak's RT difference profile.

#### Experimental Data Sets

Four formalin-fixed and paraffin-embedded (FFPE) data sets were used for identifying modifications associated with the fixing process and storage as selected by Tabb *et al.* (16) for their study. Two of these data sets, titled “Nielsen” (PXD000743) and “Buthelezi” (PXD013107), were acquired on SCIEX TripleTOFs (17). The Nielsen data set was acquired on a TripleTOF 5600+ and consists of 218,449 scans across 20 SCIEX.wiff files, and the Buthelezi data set was acquired on a TripleTOF 6600 and consists of 474,726 scans across 12 SCIEX.wiff files. Two other data sets, titled “Zimmerman” (PXD001651) and “Nair” (PXD013528), were acquired on Thermo Q-Exactive instruments (18). The Zimmerman data set consists of 79,803 scans across five .raw files, and the Nair data set consists of 245,589 scans across 10 .raw files. Files were acquired as .raw or .wiff files and converted to mzML using ProteoWizard's MSConvert, version 3.0.18208.

The synthetic peptide data set was obtained from ProteomeXchange (PXD004732) in mzML format (19). Only MS runs with the 3xHCD label were included in our analysis. Peptide pools labeled as SRM were also excluded. This synthetic peptide data set consists of unmodified proteotypic human peptides fragmented on a Thermo Fisher Orbitrap Fusion Lumos instrument. Cysteines were incorporated as alkylated cysteines during synthesis.

Additional data sets used in this work were obtained from the Clinical Proteomics Tumor Analysis Consortium (CPTAC) data portal in mzML format (20). These were limited to MS runs generated from the CompRef samples, a CPTAC reference material created using breast cancer xenograft pools for QC and data harmonization purposes. The samples were analyzed using tandem mass tag 10 (TMT-10) labeling-based technology. The first cohort of six experiments (10-plex TMTs) consists of samples processed at three sites (two experiments from each site)—the Broad Institute (BI), Johns Hopkins University (JHU), and Pacific Northwest National Laboratory (PNNL)—acquired on an Orbitrap Fusion Lumos as part of the CPTAC harmonization study (21). The second cohort consists of the same CompRef samples processed as longitudinal QC samples as part of three CPTAC data sets: the clear cell renal cell carcinoma data set generated at JHU (three experiments), the lung adenocarcinoma data set generated at BI (four experiments), and the uterine corpus endometrial carcinoma data set generated at PNNL (four experiments). All these data, at all sites, were generated using the CPTAC harmonized data generation protocol (21). These data were processed together

using PTM-Shepherd's multiexperiment setting to generate a single report.

#### Database Search and Statistical Validation

All analyses were performed using a database constructed from all human entries in the UniProtKB protein database (retrieved July 29, 2016). Reversed protein sequences were added as decoys, and common contaminants were appended (total targets and decoys: 141,585). Unless specified otherwise, all data sets were processed with the following parameters. Data were searched using MSFragger, version 2.1(4) with a precursor mass tolerance of  $\pm 500$  Da. Isotope error correction was disabled, and one missed tryptic cleavage was allowed for peptides of 7 to 50 residues in length. Oxidation of methionine was included as a variable modification, and cysteine carbamidomethylation was included as a fixed modification. MSFragger mass calibration and parameter optimization was performed for all data sets (22), including fragment ion tolerance. Shifted ions were not used in scoring.

FFPE data were processed using a  $-200$  to  $500$  Da mass range to match that used in the original publication (16). CPTAC data were processed with protein N-terminal acetylation and peptide N-terminal TMT mass of 229.1629 Da as variable modifications (TMT was also specified as fixed modification on Lys). In addition, CPTAC data were searched against combined UniProtKB mouse plus human protein database (retrieved February 10, 2020), with its respective reversed decoys appended to the database, resulting in 252,401 total target and decoy proteins.

PSMs identified using MSFragger were processed using PeptideProphet (23) via the Philosopher v2.0.0 toolkit (14). All processing and filtering were performed on a per-experiment basis. The four FFPE data sets were processed as four experiments. The chloroacetamide (CAA)-labeled HeLa cells data set was processed as one experiment containing all 39 fractions. CPTAC samples were grouped experiment wise, with each experiment containing all 24 fractions. Because of large size, the ProteomeTools synthetic peptide data set was processed as 11 subsets split based on the five-digit identifier at the beginning of each filename. PeptideProphet parameters for all analyses were default open search parameters: semiparametric modeling, clevel value set to  $-2$ , high accuracy mass mode disabled, mass width of 1000, and using expectation value for modeling. Resulting PSM matches were filtered to 1% FDR using target-decoy strategy with the help of Philosopher filter command.

## RESULTS AND DISCUSSIONS

### Overview of PTM-Shepherd

The overview of PTM-Shepherd computational workflow is shown in Figure 1. The process starts with PTM-Shepherd reading FDR-filtered PSM lists (produced by MSFragger and Philosopher, optionally with open search artifacts removed using Crystal-C (24)), and the mass shift for each PSM, to construct a mass shift histogram (11). After smoothing the histogram, PTM-Shepherd picks peaks based on a mixture of peak prominence and signal-to-noise ratio. From this list of detected peaks, the top 500 (by default) are selected for downstream processing. Basic abundance statistics are then calculated for this list of detected peaks. PSMs are assigned to a particular peak if their mass shift falls within the reported peak boundary, and abundance of the peak is calculated based on spectral counts. PTM-Shepherd can also operate in a multiexperiment mode. In this mode, peak detection is



performed on an aggregate mass shift histogram from all experiments, generated from the mass shifts of each experiment weighted according to the proportion of the total PSMs they comprise. The use of a combined histogram for peak detection can greatly simplify comparisons between modifications detected in different conditions and experiments. In this multiexperiment mode, the summary attributes for each detected peak are generated separately for each experiment and for all data combined.

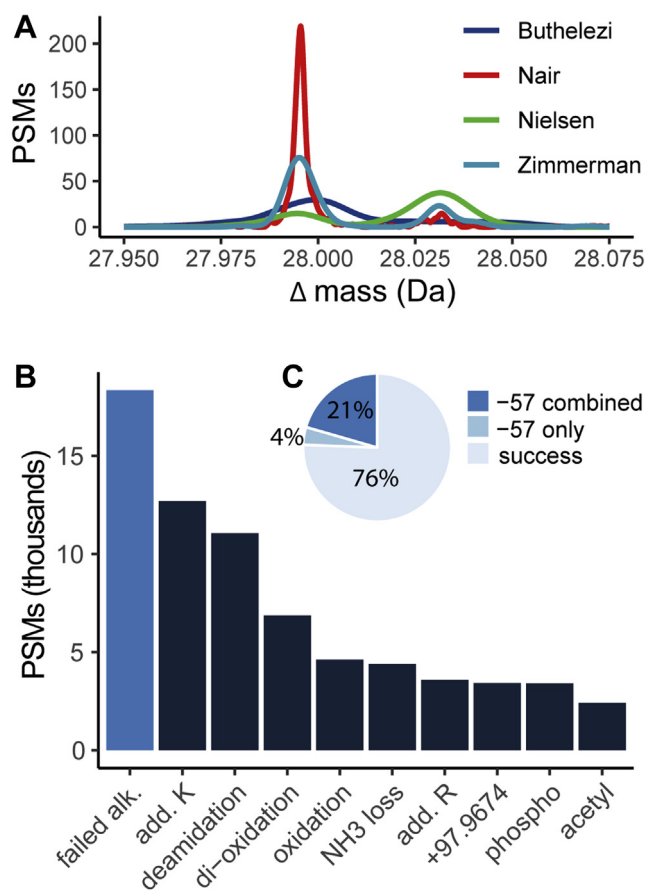
Once peaks in the mass shift histogram have been called, PTM-Shepherd attempts to determine their identities. Mass shifts are iteratively annotated using entries from the Unimod (15) modification database, isotopic error peaks, and user-specified mass shifts, allowing the mass difference to be decomposed into at most two modifications. PTM-Shepherd also constructs localization profiles for each peak. Localization profiles are constructed for each experiment, reporting an N-terminal localization rate and a normalized amino acid propensity for each modification. This analysis is performed for every PSM by placing the mass shift at each amino acid in turn and rescoring the PSM (with the original spectrum) using the scoring function presented in MSFragger (see [Methods](#) section).

PTM-Shepherd also computes several metrics that are useful for gaining a better understanding of the nature of those detected mass shifts. For each peak, PSMs containing that mass shift are compared with their unmodified counterparts if present within the same run. First, cosine SS between modified and unmodified peptides is computed, which is useful for determining how the modifications affect spectra. Then, RT effects are examined, and the average difference in RT between the peptide with and without modification is reported.

#### PTM-Palette Discovery: Analysis of FFPE Samples

Understanding how sample processing and storage affect proteins is critical to maximizing their identification. Analysis of tissue samples preserved using FFPE technique warrants the inclusion of additional modifications to reflect changes in proteins following formalin fixation. FFPE samples are also typically analyzed after long-term storage, during which they could be exposed to high temperatures and sunlight (25). Although previous studies have examined which modifications should be included when analyzing proteins from FFPE samples (25), this was revisited recently by Tabb *et al.* (16) using a two-pass search. First, an open search was used to identify prevalent mass shifts. Second, they performed a traditional search and informed the localization of their mass shifts with chemical knowledge. We sought to investigate how PTM-Shepherd could be used to validate their findings and streamline this analysis for other data sets and sample preparation protocols.

After their first-pass open search, Tabb *et al.* (16) found five modifications that were consistently present across the four data sets analyzed: methylation, dimethylation, single oxidation, double oxidation, and variable carbamidomethylation.



**Fig. 2. Basic PTM-Shepherd applications.** A, PTM-Shepherd identifies two peaks in close proximity for the four data sets of Tabb *et al.* All four data sets (Zimmerman, Nair, Nielsen, and Buthelezi) show a mixture of two Gaussian peaks about 28 Da. The consistently more intense peak is at 27.9949 formylation. Only in the Nielsen data set does dimethylation (28.0313) approach formylation's intensity. B, PTM-Shepherd identifies more failed alkylation than other common modifications such as deamidation and not-Met oxidation. C, PTM-Shepherd modification decomposition identifies six times as much failed alkylation as is identifiable based on the  $-57$  Da mass shift alone, in total accounting for a quarter of all Cys-containing peptides. PSMs, post-translational modifications.

Automated processing with PTM-Shepherd replicates most of these findings. Based on PSM counts and using the same criteria, we find mass shifts of methylation, mono-oxidation, and dioxidation within the top 10 mass shifts (excluding isotopic error peaks) for every data set ([supplemental Table S1A](#)). Interestingly, PTM-Shepherd also finds a notable discrepancy with respect to dimethylation levels. PTM-Shepherd identifies two peaks in close proximity: 27.9954 Da (corresponding to formylation) and 28.0320 (corresponding to dimethylation). Dimethylation is only higher than formylation in one data set (Nielsen; [Fig. 2A](#)), whereas others have formylation between three-fold and nine-fold higher than dimethylation. To confirm that this was not an artifact of PTM-Shepherd's signal-to-noise peak picking, we reanalyzed these results with the

DeltaMass software that implements an alternative (Gaussian mixture modeling) strategy for peak picking (11). For all these four data sets, DeltaMass found that the region of mass shifts from 27.90 to 28.10 contained two peaks (supplemental Fig. S1). For Nielsen, Nair, and Zimmerman, these are easily visible. Even the Buthezezi data set, although not exhibiting as clear a separation as the others, places the more abundant peak apex closer to the mass shift value corresponding to formylation. The presence of formylation within a list of most abundant PTMs also makes logical sense given the nature of preservation method.

Tabb and colleagues relied on chemical knowledge and other tools (7, 10) to arrive at the final search configuration that included oxidation of Met to methionine sulfone. We chose to investigate this further using PTM-Shepherd. Because a single oxidation of Met was already included as a variable modification in our open search, a Met oxidation to methionine sulfone may appear as either a variable modification and a +15.9949 Da mass shift localized to Met or a +31.9898 Da mass shift localized to Met. However, we do not observe enrichment of Met (supplemental Table S1A) localization in either of these instances. In contrast, the enrichment scores for Pro were 9.3 and 5.6 for mono-oxidation and dioxidation, respectively. Tabb and colleagues' gain in the number of PSMs when adding Met sulfone and dihydroxy Pro in the search may be explained by the diffuse nature of oxidation. When using a closed search strategy with a dynamic +31.9898 Da modification on Met, any occurrence of two +15.9949 Da events—e.g., on two alternative oxidation sites—might be interpreted as a Met sulfone because peptide ions will converge downstream of the theoretical and experimental oxidation sites. The same phenomenon can occur with multiple instances of hydroxyproline. Collagen is known to contain massive amounts of hydroxyproline and as such is likely to produce peptides with multiple hydroxyprolines (26). To determine whether the +31.9898 Da mass shift was attributable to multiple hydroxyprolines co-occurring on the same peptide, we checked whether PSMs containing it were more likely to map to collagen proteins than noncollagen proteins. This mass shift was overwhelmingly more likely to occur on collagen proteins (odds ratio = 43.5;  $p < 10^{-5}$  by Fisher's exact test), confirming that it is a combination of multiple hydroxyprolines and can be captured via including hydroxyproline in a PTM palette. Overall, in our experience in this and other data sets, PTM-Shepherd provides a very reasonable estimate of the most likely modification sites for a particular mass shift.

"Formaldehyde adduct" (+30.0106 Da, annotated as methylole in Unimod) is a known modification observed on peptides from FFPE samples, and it was detected at high levels in the Nair and Zimmerman data sets (in the top 10). According to PTM-Shepherd analysis, this mass shift lacks any significant localization characteristics (localized less than 10% of the time), indicative of a noncovalent adduct. In general, identification of labile modifications is one of the advantages of open database searching with MSFragger compared with other

PTM-focused tools or closed searches with variable modifications, all of which are less effective at finding labile modifications that cannot be localized. Summarizing our observations, PTM-Shepherd suggests a slightly modified version of the PTM palette from FFPE samples proposed by Tabb *et al.*: oxidation of Met (+15.9949 Da), hydroxyproline (+15.9949 Da on Pro; ideally, specified for collagens only), formylation on Lys and N termini (+27.9949 Da), and methylation of Lys and N termini (+14.0157 Da). It may also be beneficial to include the methylole (+30.0106) adduct but using the mass offset search option of MSFragger that allows both shifted and nonshifted fragment ions in scoring rather than as variable modification (22, 27).

#### Detection of Cysteine Artifacts Following Underalkylation

Cysteine is an extremely reactive amino acid, frequently picking up an array of chemical modifications when exposed (28). Unspecified mass shifts, such as those resulting from chemical modifications of Cys, confound peptide identifications and lead to lower recovery rates. Cys alkylation restricts the number of chemical derivatives it can form and prevents interference from disulfide bonds, and as such has been a mainstay of proteomics processing for decades (29). CAA and iodoacetamide (IAA) are the two most common alkylating reagents used in proteomics workflows. Previous comparisons of these reagents found that IAA generally has a higher rate of cysteine alkylation than CAA when applied at the same concentration, but with the caveat of higher rates of off-target effects (30). Here, we have tested the ability of PTM-Shepherd to uncover cysteine artifacts in proteomic data sets. Bekker-Jensen *et al.* (31) rigorously tested shotgun proteomics protocols to determine an optimal strategy for rapidly generating a comprehensive profile of human proteomes, ultimately producing a valuable repository of high-quality and deep proteomics data. Their protocol also included a 10 mM treatment with the alkylating agent CAA, which, per Schnatbaum *et al.* (30), only achieved two-thirds the alkylating efficiency of 10 mM IAA in complex mixtures. Unlike the other samples we analyzed for this article, this protocol also did not denature protein samples before adding the alkylating agent. This likely contributes heavily to underalkylation as well. As such, it presents an exceptional opportunity to examine these cysteine artifacts.

Open search analysis followed by PTM-Shepherd shows a number of prevalent mass shifts enriched on Cys, consistent with what we expect from underalkylated samples. Note that because Cys alkylation was searched as a fixed modification to elucidate the identity of modifications occurring on unalkylated Cys residues, the mass shift must be decomposed into two components: a failed alkylation event ( $\Delta m = -57.0215$  Da from the theoretical mass of the identified peptide) and the modification itself. Consider the mass shift  $-9.03680$  Da detected in this data set. PTM-Shepherd decomposes this mass shift into a failed alkylation event

(−57.0215 Da) and a triple oxidation of Cys to cysteic acid ( $\Delta m = +47.9847$  Da). This becomes particularly important when trying to directly assess the number of failed alkylations in a sample. Strictly counting the number of −57.0215 Da mass shifts will severely undercount their total occurrences because it ignores cases where it is found in conjunction with another modification, which are very likely given the reactivity of Cys. We implemented an additional parameter in PTM-Shepherd to account for this that prioritizes user-defined modifications and allows them to identify mass shifts that do not directly correspond to entries in Unimod (15). On its own, failed alkylation was the sixth most abundant mass shift and was more prevalent than other common events that are often accounted for in closed searches, such as pyroglutamate formation, and accounted for 3.9% (3450 PSMs) of the 89,281 total Cys-containing PSMs (supplemental Table S2A). However, because it frequently occurs with other mass shifts as demonstrated previously, we also pooled the instances in which it was annotated as one of two mass shifts on a peptide. Remarkably, the total number of failed alkylation events jumps to 20.5% (18,343 PSMs) when considering all instances of failed alkylation annotations (Fig. 2C). When considering decomposed mass shifts, failed alkylation is nearly twice as common as deamidation and four times as common as non-Met oxidation events (Fig. 2B and supplemental Table S2B).

After applying an abundance cutoff of 0.01% of total spectra, we detected 10 mass shifts exhibiting strong Cys localization (> tenfold enrichment) that were also annotated with a failed alkylation of Cys. These were a large portion of the broadly occurring Cys-enriched PTMs in the samples (supplemental Table S2A). The most abundant of these modifications correlate with what would be expected in a poorly alkylated sample. The +1 and +2 isotopic error peaks in conjunction with unalkylated Cys were particularly abundant, accounting for 1601 combined spectra. The aforementioned −9.0368 Da—triple oxidation on Cys without alkylation—was most prevalent aside from these. Its heavily enriched localization to Cys (48.5-fold) lends credence to this compound identification that would be missed by other annotation tools and, consequently, a count of total failed alkylation events. Surprisingly, failed alkylation combined with a formaldehyde adduct (+12.0000 Da) was also common. The combined mass shift of −45.0216 Da was heavily localized to Cys and had a 96% N-terminal localization rate, pointing to potential thiazoladine formation *via* N-terminal Cys cyclization. These are known to occur in formaldehyde-treated data, but the authors did not report the use of formaldehyde (32). A lone formaldehyde adduct mass shift accounting for 305 PSMs and heavily localized to Trp (42.5-fold enrichment) was also detected in the data set, however. Taken together, these indicate that the thiazoladine was probably an artifact of formaldehyde exposure rather than underalkylation, although the latter may be a necessary condition for Cys to react with formaldehyde. Failed alkylation of Cys and the subsequent

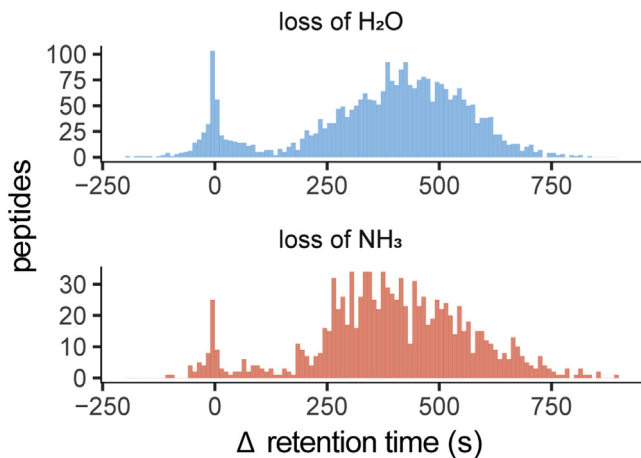
triple and double oxidations conform to our chemical knowledge of Cys artifacts and, along with glutathione disulfide as a biological modification, comprise 8.7% of all Cys-containing PSMs. Including these modifications should increase Cys-peptide recovery in underalkylated samples.

#### *PTM-Shepherd Computed Metrics Facilitate Granular PTM Identification*

Open searches are inherently limited in the information they provide, providing only peptide lists and their associated mass shifts (4). Data interpretation efforts are further complicated by the ambiguity of mass shifts. Two methylation events and an ethylation event, for instance, would be indistinguishable from each other based on mass. However, more granular identities can be discerned by incorporating additional metrics: changes in RT, SS, and localization. To demonstrate that these additional metrics improve open search result comprehension, we analyzed the synthetic unmodified tryptic peptide data set generated as part of the ProteomeTools project (19). This data set allows us to examine and characterize instrumental artifacts apart from confounding biological factors.

In-source losses from peptides result from low-energy fragmentation pathways that can occur during tandem MS as well as during ionization and transmission, resulting in artifactual changes to the observed precursor mass (33). Because in-source losses occur after column elution and consequent RT recording, they have no effect on peptide RT. This property can be used to distinguish them from sample modifications (34). We used PTM-Shepherd to elucidate the origins of two of this data set's most common mass shifts attributable to both in-source losses and real modifications: loss of H<sub>2</sub>O and loss of NH<sub>3</sub> (supplemental Table S3). Peptides with multiple spectra corresponding to each loss had their RT shifts pooled and collapsed to their median.

Interestingly, both losses of H<sub>2</sub>O and NH<sub>3</sub> exhibited bimodal changes in RT (Fig. 3). For peptides presenting losses of H<sub>2</sub>O ( $\Delta m = -18.0104$  Da; 2961 peptides), both composite RT distributions were approximately Gaussian with approximate means of 0 and 450 s. As anticipated, many peptides (16.7%) fall within the mean zero distribution, indicating that they do not experience increases in column RT despite the loss of a highly polar group. This is characteristic of in-source losses and indicates that peptides within this distribution are exhibiting in-source loss of H<sub>2</sub>O in the mass spectrometer prior to precursor selection. Peptides presenting losses of NH<sub>3</sub> ( $\Delta m = -17.0270$  Da;  $n = 1094$ ) showed a similar pattern. Note that H<sub>2</sub>O and NH<sub>3</sub> are only two examples of in-source loss and, in some cases, entire residues can be lost *via* this mechanism. As a significant source of instrumental bias, it is important to be able to classify in-source losses properly and remove them from experimental sample pools. In fact, for researchers studying the isobaric biological forms of these mass shifts, it is critical to exclude these.



**FIG. 3. Retention time (RT) profiles for peptides with losses of H<sub>2</sub>O and NH<sub>3</sub>.** Modified peptides are compared with their homologous unmodified peptides, with multiple RT changes being collapsed to their median. The effect on RT for losses of H<sub>2</sub>O (*top*) and NH<sub>3</sub> (*bottom*) is distributed bimodally. These mass shifts are known to correspond to both in-source losses and spontaneous conversions. In-source losses should not have an effect on RT, and as such are suspected to fall within a Gaussian distribution centered at zero.

The second population of peptides with H<sub>2</sub>O and NH<sub>3</sub> loss exhibited a RT-shift consistent with a pre-elution modification—the loss of a polar group increased RT by an average of 450 s. H<sub>2</sub>O losses are known to manifest as a conversion to pyroglutamic acid from Glu (35) as well as on Asn, Gln, Ser, Thr, Tyr, Asp, and Cys as sample-derived modifications (15). NH<sub>3</sub> losses are known to manifest as a conversion to pyroglutamic acid, but from Gln rather than Glu (36). Other losses of NH<sub>3</sub> are known to occur on some N termini occupied by Thr, Ser, and Cys, and on any Asn (15). While RT shifts alone do not contain enough information to fully identify these modifications, additional metrics calculated by PTM-Shepherd—localization propensity and modified-to-unmodified peptide similarity—allowed us to delineate the primary sources of these mass shifts.

In the case of a loss of NH<sub>3</sub>, two primary sources (in addition to the in-source losses) were identified: a spontaneous conversion of Gln to pyroglutamate and a cyclization of Cys. Cys cyclization is expected to be present in peptides being synthesized with carbamidomethylated Cys, as the reaction is known to occur after alkylation (37). Localization analysis showed that Cys (enrichment score = 10.4) and Gln (enrichment score = 4.7) were the two most enriched residues for this mass shift (Fig. 4A). In-source losses localized to Cys are rare and in-source losses localized to Gln are relatively common (38), which is reflected in the number of peptides each residue produces with  $\Delta RT = 0$ . To illustrate, Cys has an RT profile very different from other residues in aggregate, whereas Gln has a similar distribution (Fig. 4C); none of those containing Cys localized NH<sub>3</sub> losses were in-source losses, as compared with 19.6% of other residues in aggregate.

Unlike NH<sub>3</sub>, we expected losses of H<sub>2</sub>O to only be heavily enriched in glutamate. Glutamate is known to have two sources for loss of H<sub>2</sub>O; it is known to be a source of water in-source loss but can also spontaneously undergo N-terminal cyclization *in vitro* to produce the same mass shift (39). The localization enrichment profile for loss of H<sub>2</sub>O (Fig. 4B), however, revealed that two residues were exceptional contributors to the prevalence of PTM: Glu (enrichment score = 2.2) and Lys (enrichment score = 7.6). We identified populations of peptides with losses of H<sub>2</sub>O corresponding to both these populations based on  $\Delta RT$  as described previously (Fig. 4B, red). An intense peak near  $\Delta RT = 0$  s indicates that many unique peptides are capable of producing in-source losses of H<sub>2</sub>O (22.8%) on Glu, consistent with the conclusions of Sun *et al.* (38). Another peak near  $\Delta RT = 300$  s indicates that the remainder of the peptides had a loss of H<sub>2</sub>O that was present before column elution, and as such is likely to be an N-terminal cyclization of glutamate occurring *in vitro*.

Even more so than glutamate, lysine was the largest contributor to losses of H<sub>2</sub>O (example spectrum at supplemental Fig. S2). Puzzlingly, lysine's side chain does not have a hydroxyl group that it can readily lose, and as such any H<sub>2</sub>O losses attributable to lysine must be derived from the C-terminal hydroxyl group of tryptic peptides. It is worth noting that this phenomenon is unique to lysine, as Arg had the lowest localization enrichment score (0.2) of all 20 residues (Fig. 4B). Based on RTs, there were also no appreciable H<sub>2</sub>O in-source losses attributable to lysine—consistent with previous findings (38)—indicating that these lysines were being dehydrated prior to elution (Fig. 4D). We believe this is most likely because of a C-terminal lysine cyclization event. Although undescribed in proteomics, lysine derivative cyclization has been induced in other settings (40). This theory is supported by SS calculations between peptides with and without this lysine-localized mass shift (Fig. 4F). There is exceptionally low SS between the spectra of modified and unmodified peptides. Peptides containing nonlabile modifications and modifications near the C terminus result in low SS; nonlabile modifications are likely to be retained in the MS/MS spectra rather than being removed during MS1 analysis, and C-terminal modifications shift the intense y-ion series. A covalently bound lysine ring structure on peptide C termini fits both these criteria and may be the underlying cause of low SS.

The similarity profiles for losses on Glu, Cys, and Gln were distinct from Lys in that they were not enriched for peptides with MS/MS showing low similarity to their unmodified counterparts. This may be accounted for by the fact that Glu cyclization, like Cys and Gln, occurs at the N terminus; shifting the b-ion series has less of an effect on the MS/MS spectra than shifting the y-ion series. Unsurprisingly, all three of these have similarity profiles roughly corresponding to the proportion of their spectra experiencing *in vitro* modifications, which are the only modifications we can be sure are occurring at the N terminus.



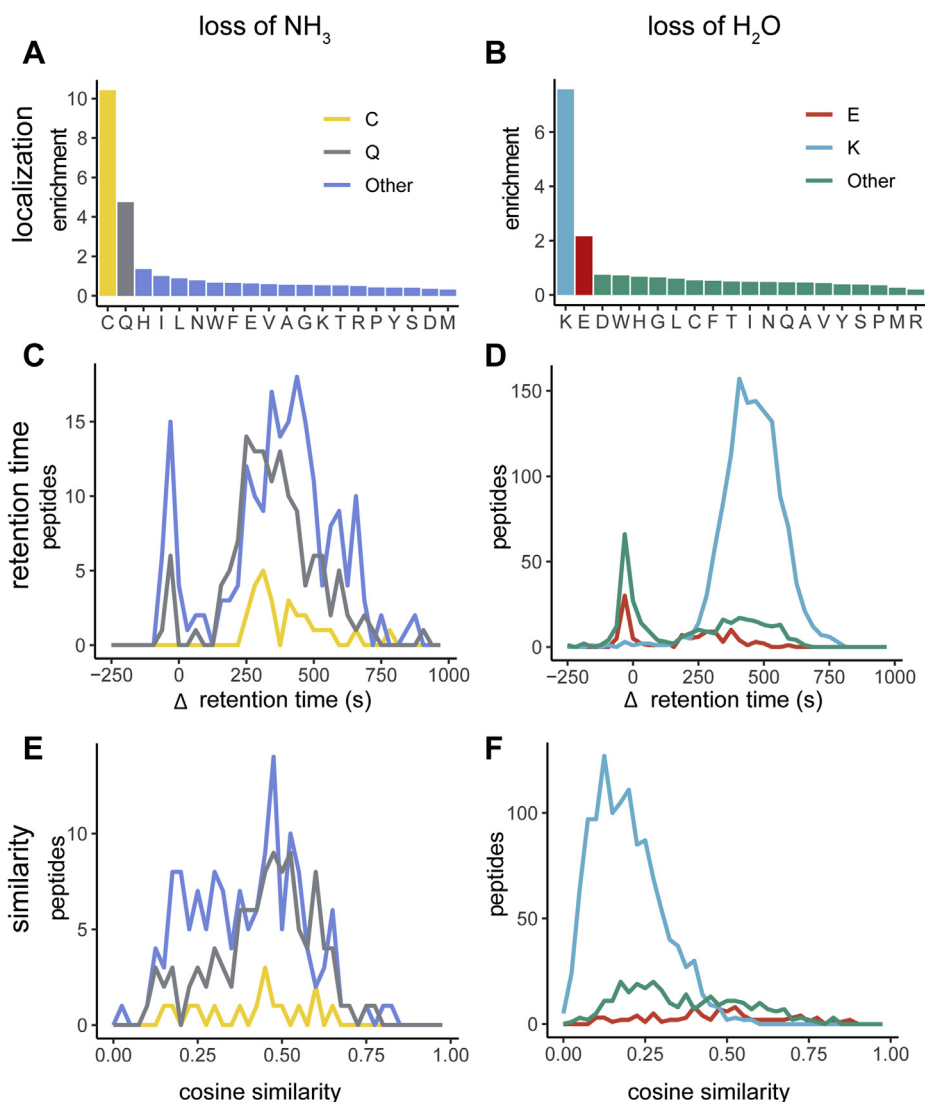


FIG. 4. **Analytical profiles for losses of  $\text{H}_2\text{O}$  and  $\text{NH}_3$ .** A and B, localization profiles reveal a nonhomogeneous landscape with specific residues showing enrichment. C and D, select modifications are distinguishable from background in-source decays in their effect on retention time. E and F, similarity scores show lower profiles for C-terminal modifications on Lys, whereas N-terminal modifications on Glu, Cys, and Gln have higher similarity.

Overall, by including metrics beyond mass shifts in PTM identification, we show that much more information beyond the chemical composition of a mass shift can be deduced. RTs can be used to discriminate between in-source losses and sample modifications, localization profiles can be used to deduce biological or chemical origins, and SS provides additional localization and lability metrics. Incorporating all these, we found, to our knowledge, a previously unknown or at least underappreciated modification (C-terminal lysine cyclization) in a deeply consequential synthetic peptide library.

#### *PTM-Shepherd in Multiexperiment Settings*

PTM-Shepherd can be run in a multiexperiment mode to analyze modification profile across a large number of

experiments. Such an analysis could be performed for visualization of interesting biological trends and in search of experiment-specific biological modifications. It could also be useful for QC and detection of batch effects—a common source of variation in high-throughput data (41). Previous efforts have been made to identify MS performance metrics (42), and some groups have shown how these can be leveraged to identify QC issues (43) and better understand intralaboratory and interlaboratory variability (44, 45). We posited that open search-derived modification profiles could be used to determine interexperiment variation while simultaneously providing insight into its origins.

To evaluate PTM profiling in multiexperiment settings, we used CPTAC CompRef reference material data (pooled tumor

TABLE 1  
Top mass shifts from CPTAC QC samples

Mass shift	PSMs	Assigned modification	% in unmodified	Similarity	Delta RT	N-terminal rate (%)	AA_1	AA_2
0.000	4,372,080	None	100.00	0.98	0	0		
1.002	913,405	+1 isotopic error	79.43	0.81	5	3		
229.163	312,295	TMT	74.52	0.38	-70	6	S (5.6)	T (2.4)
2.005	177,121	+2 isotopic error	75.95	0.74	4	1		
-0.984 (15.011*)	156,315	NH addition to M*	90.25	0.71	-173	1	M (15.8)	
230.166	141,498	+1 isotopic error + TMT	72.51	0.42	-153	6	S (3.8)	
0.017	117,922	+1 isotopic error + NH addition to M*	69.97	0.78	1	1	M (5.8)	
0.984	110,967	Deamidation	68.84	0.68	53	6	N (12.7)	R (5.3)
17.026	92,862	Deuterated methyl ester	94.96	0.76	-50	1		
15.011	87,013	Conversion of carboxylic acid to hydroxamic acid	95.22	0.54	-434	4	E (5.5)	D (3.5)
-1.002	72,465	-1 isotopic error/ +1 isotopic error + half of a disulfide bridge*	81.40	0.78	16	1	C (4.8)	W (4.5)
15.995	59,846	Oxidation	91.19	0.49	-342	21	W (24.0)	M (11.7)
100.016	59,070	Succinic anhydride labeling reagent	95.76	0.55	596	77	P (3.6)	S (2.1)
27.995	52,323	Formylation	93.22	0.61	265	61	S (2.7)	
79.967	48,379	Phosphorylation	55.30	0.64	569	4	S (6.2)	
21.981	47,471	Sodium adduct	98.72	0.25	-148	2		
115.027	47,022	Cleavage product of EGS protein crosslinks by hydroxylamine	97.81	0.58	610	86	H (2.5)	
43.010	45,450	Carbamylation	97.10	0.64	640	70		
-18.010	45,445	Dehydration/Pyro-glu from E	97.48	0.64	-175	12	D (3.2)	T (2.5)
1.987	42,988	+1 isotopic error + deamidation	69.08	0.66	45	4	N (11.9)	R (3.0)

CPTAC, Clinical Proteomics Tumor Analysis Consortium; EGS, ethylene glycolbis(succinimidylsuccinate); PSMs, peptide-spectrum matches; QC, quality control; RT, retention time; TMT, tandem mass tag.

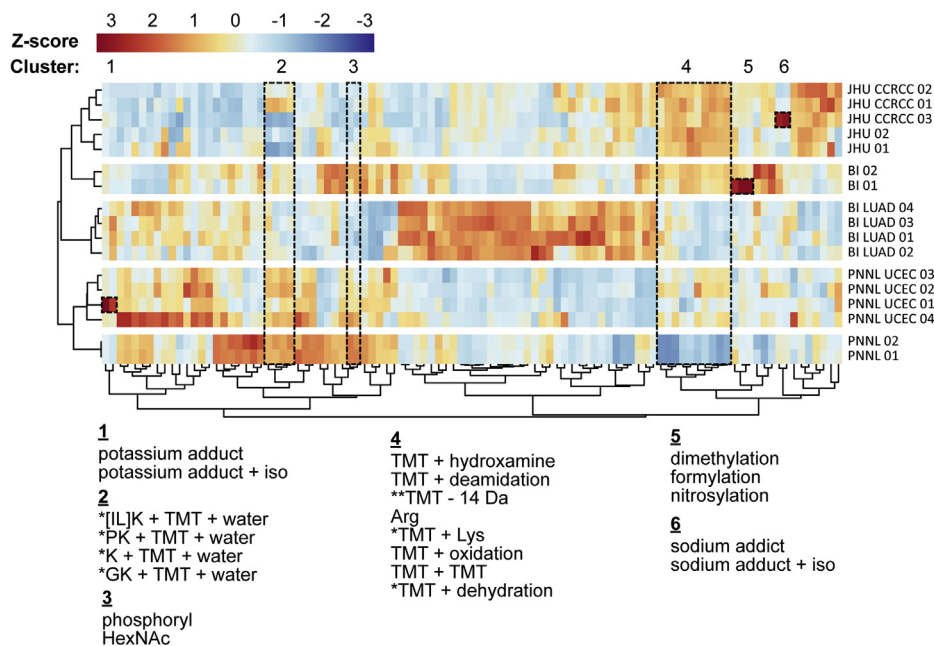
Assigned modifications correspond to automated Unimod matches, with \* indicating a partially manually reannotated mass shift. The % in unmodified column corresponds to the percent of PSMs with a matching unmodified PSM in the unmodified bin. Top two enriched amino acid localizations are shown in columns denoted AA.

xenografts comprising 10 samples each from two different breast cancer subtypes, cryopulverized and shipped to different processing locations) obtained from the CPTAC data portal (see [Methods](#) section). The samples were processed at three different locations (PNNL, JHU, and BI) and analyzed using TMT 10-plex labeling technology as part of the CPTAC3 Harmonization study (21). The same CompRef samples were also analyzed as longitudinal QC samples as part of the three large cancer profiling studies, clear cell renal cell carcinoma (46) (MS data collected at JHU; uterine corpus endometrial carcinoma (47) (MS data collected at PNNL), and lung adenocarcinoma (48) (MS data collected at BI).

We first investigated the most abundant mass shifts identified by MSFragger and PTM-Shepherd in these data (Table 1; see [supplemental Table S4A](#) for the full list), which revealed several interesting observations. In general, PTM-Shepherd accurately reconstructed expected trends including the localization profiles of the most abundant modifications. For example, carbamylation and formylation were most highly enriched on N terminus, phosphorylation on Ser, and oxidation on Trp. Not considering isotope errors, the mass shift of 229.1629 Da (TMT overlabeling) was the most

common modification, localized predominantly to Ser (enrichment factor of 5.6). Of note, only 74.5% of peptides found with TMT on Ser were also found in unmodified form (*i.e.*, with unlabeled Ser). In contrast, many other abundant modifications, such as formylation and carbamylation, were found in both modified and unmodified forms in almost all cases. Interestingly, the second most abundant modification was a mass shift of 15.0107, predominantly localized to Met (that was indistinguishable in MSFragger output from a combination of oxidation and -0.9842 Da loss on Met). This mass shift may represent the addition of an NH group to Met because of exposure to hydroxylamine, a reagent used in TMT labeling. At present, UniMod database annotates a 15.0107 Da mass shift only as conversion of carboxylic acid to hydroxamic acid, with Asp and Glu as only possible sites (which were observed in these data, but at a lower frequency than on Met, see [Table 1](#)).

The PTM profiles resulting from PTM-Shepherd analysis of these data are presented in [Figure 5](#). Sample-wise K-means clustering revealed distinct sample clusters, and mass shift-wise clustering on correlation between columns (transposed PTM-Shepherd output) revealed some highly similar



**FIG. 5. Clustered heat map representation of CPTAC3 quality control samples transposed from PTM-Shepherd output.** Values shown are column-wise z-scores of spectral counts. Column clustering shows highly related modifications, and row clustering shows experiments clustering by processing location. Mass shift clusters discussed in the text are numbered, and their corresponding mass shifts are shown left to right in the bottom of the figure. Samples processed longitudinally throughout their respective studies are indicated using tumor type label (LUAD, UCEC, or CCRCC). Samples with no tumor type label were processed as part of the CPTAC harmonization study. Mass shifts in cluster 2 correspond to negative mass shifts. \*This annotation was constructed manually. \*\*−14 Da can correspond to a large number of modifications and single-residue mutations. BI, Broad Institute; CCRCC, clear cell renal cell carcinoma; JHU, Johns Hopkins University; LUAD, lung adenocarcinoma; PNNL, Pacific Northwest National Laboratory; TMT, tandem mass tag; UCEC, uterine corpus endometrial carcinoma.

modifications. Sample clustering precisely reconstitutes sample processing location. For example, cluster 4 in [Figure 5](#) shows a series of mass shifts related to TMT overlabeling or TMT labeling that was not captured by fixed sequence expansion on Lys and dynamic sequence expansion on peptide N termini. PNNL data consistently show lower TMT overlabeling than BI and JHU for every mass shift in this cluster, and PSMs corresponding to a single additional TMT are 5 to 8 times lower than at the other two locations. BI and JHU also show enrichments of TMT labeling on Ser and, to lesser degree, Thr.

Although we expect to see differences in TMT labeling fidelity, PTM-Shepherd also allows us to explore unexpected batch effects. Lenčo *et al.* (49) recently raised concerns about the use of formic acid in sample preparation, specifically stating that an excess of Ser, Thr, and N-terminal formylation events are present in samples reconstituted with it. Within the CPTAC harmonization study, the localization profile did exactly match that described by Lenčo *et al.*: Ser enrichment of 2.7, Thr enrichment of 1.9, and a 61% potential N-terminal rate ([supplemental Table S4A](#)). Interestingly, this formylation peak appears to be disproportionately high in the first of the two BI replicates from the harmonization studies ([Fig. 5](#), cluster 5). BI01 replicate had formylation 10- to 20-fold higher than JHU01 or PNNL01 and was even fourfold higher than BI02. Ser and Thr also exhibit inflated formylation localization

for this sample, consistent with the results of Lenčo *et al.* ([supplemental Table S4, B–C](#)). Overall, a deeper analysis may be warranted in future studies to reduce batch effects caused by formic acid use. These analyses could be extended to other sample handling artifacts, *e.g.*, potassium adducts (cluster 1) and sodium adducts (cluster 6), which also exhibit marked longitudinal variability.

PTM-Shepherd also reveals how instrument parameters might be playing a role in laboratory-specific batch effects. We noted four large negative mass shifts that were differentially identified across laboratories ([Fig. 5](#), cluster 2). RT profiling showed that these mass shifts were likely to be in-source losses ([supplemental Fig. S3](#)). Their profiles were similar to loss of ammonia, a known in-source loss, and dissimilar to formylation, a known pre-elution modification ([supplemental Fig. S3](#)). Sequence analysis facilitated manual decomposition, revealing that three of the mass shifts were composed of a hydrophobic residue (Iso, Leu, Pro, or Gly), a C-terminal Lys, a TMT tag, and water. The final mass shift was the loss of a C-terminal Lys, a TMT tag, and water. It is possible that slight differences in fragmentation or ionization energies between laboratories may be manifesting as disparities in precursor charge states and proton mobility; however, determining the mechanism through which these are occurring is outside the scope of this work.

Aside from analyzing how samples cluster together, there is also useful information to be gleaned from analyzing how mass shifts cluster together. We expect that related modifications should have highly correlated abundances between experiments, with coclustering of isotopic error peaks as the most obvious example (e.g., potassium and sodium adducts and their related +1 isotopic error peaks; Fig. 5, clusters 1 and 6). Using an z-score normalization to coerce PTMs derived from related sources across experiments to cluster together and clustering on the correlation between columns, we were able to identify some unknown mass shifts based on their coclustering with other (known) modifications. In the TMT-related cluster noted previously (Fig. 5, cluster 4), we observed three mass shifts that were missed by automatic annotation. Of the six that were annotated, five of the mass shifts are directly attributable to TMT overlabeling and one to a missed tryptic cleavage or additional Arg at one end of the peptide sequence. This knowledge allowed us to explain one of the unannotated mass shifts as combinations of missed cleavages and TMT labeling. One modification (+357.2584 Da) is precisely the mass of a TMT-labeled Lys residue. This was missed by automatic annotation because both addition of Lys and TMT10 Plex would have to be more abundant than the combination of the two during PTM-Shepherd analysis. The other unexplained mass shift (+213.1680 Da) can be explained by a combination of TMT overlabeling and a dehydration event or a misattributed Met oxidation included as a variable modification. Overall, this analysis demonstrates the utility of multi-experimental analysis with PTM-Shepherd to better identify that the mass shifts are not amenable to automatic annotation.

Finally, although our aforementioned analysis focused mostly on modifications introduced because of labeling and other sample handling steps, clustering of mass shifts may also be useful for uncovering correlated biological modifications. Of note, clustering of PTM-Shepherd results shows that phosphorylation (+79.9663 Da) is most correlated with HexNAc (+203.0794 Da) (Fig. 5, cluster 3). Interestingly, coenrichment of glycopeptides was recently observed in data sets experimentally enriched for phosphopeptides (50); however, no phosphopeptide enrichment steps were applied to generate the data used in this work.

#### CONCLUSIONS

Despite advancements in computational proteomics, many MS/MS spectra remain unexplained. Open searching with tools like MSFragger has proven to be an effective way to overcome the limitations of traditional database searches by removing the requirement of having prior knowledge of the peptide modifications present in the sample. The modifications elucidated by open searches, however, lack many of the metrics necessary to make proper determinations about their identities and origins. We addressed these challenges in PTM-Shepherd, which produces comprehensive PTM profiles for

open search-derived mass shifts, including multiple Unimod annotations, RT changes, SS, and localization profiles.

We demonstrated the utility of PTM-Shepherd in four examples, providing a broadly applicable guide for others interested in using open searches for PTM analysis in their own research. First, in the development of an FFPE-treatment PTM palette, we showed how PTM-Shepherd disambiguated two overlapping peaks: formylation and demethylation. We also demonstrated how PTM palettes can be easily constructed for other sample preparation methods without extensive postprocessing. Second, we showed how PTM-Shepherd's unique ability to decompose mass shifts into multiple Unimod modifications allows us to identify and quantify the degree of failed alkylation, although this is easily extensible to other scenarios, e.g., identifying mass shifts corresponding to an absent variable modification and another co-occurring modification. We also demonstrated how incorporating additional metrics into PTM identification provides researchers with more granular and high confidence PTM identities, including the ability to distinguish between sample-derived and instrument-derived artifacts. Finally, when applied to data from a large multicenter proteomics study, PTM-Shepherd helped us to visualize batch effects and the effect of sample processing location, as well as elucidate the identities of unannotated mass shifts. We believe PTM-Shepherd will become a widely used component in our MSFragger-based pipeline for comprehensive analysis of post-translational and chemical modifications, including searches for rare and even novel modifications, across a wide range of biological applications.

#### DATA AND SOFTWARE AVAILABILITY

All raw MS data used in the article can be found from the ProteomeXchange Consortium via the PRIDE partner repository or CPTAC data from the CPTAC data portal, using specific data set identifiers cited in the text. PSM lists can be accessed at [10.5281/zenodo.4042962](https://zenodo.org/record/4042962). PTM-Shepherd is available as a standalone JAR executable (<https://github.com/Nesvilab/PTM-Shepherd>) and also fully integrated into the FragPipe graphical user interface (<http://fragpipe.nesvilab.org/>).

*Acknowledgments*—The authors thank the users of our tools for their feedback.

*Author contributions*—A. I. N. and A. T. K. conceived the project; A. T. K. developed the first version of the software, later extended and improved by D. J. G.; D. J. G. performed all analyses; D. M. A., F. Y., and F. L. contributed to the software development; D. J. G. and A. I. N. wrote the article with input from all authors, and A. I. N. supervised the entire project.

*Funding and additional information*—This work was funded in part by National Institutes of Health grants R01-GM-



094231, R01-GM-135504, and U24-CA210967. D. J. G. was supported in part by the Proteogenomics of Cancer Training Program (T32-CA140044). The content is solely the responsibility of the authors and does not necessarily represent the official views of the National Institutes of Health.

**Conflict of interest**—Authors declare no competing interests.

**Abbreviations**—The abbreviations used are: BI, Broad Institute; CAA, chloroacetamide; CCRCC, clear cell renal cell carcinoma; CPTAC, Clinical Proteomics Tumor Analysis Consortium; EGS, ethylene glycolbis(succinimidylsuccinate); FDR, false discovery rate; FFPE, formalin-fixed and paraffin-embedded; IAA, iodoacetamide; JHU, Johns Hopkins University; LUAD, lung adenocarcinoma; MS, mass spectrometry; PNNL, Pacific Northwest National Laboratory; PTMs, post-translational modifications; PSM, peptide-spectrum match; QC, quality control; RT, retention time; SNR, signal-to-noise remainder; SS, spectral similarity; TMT, tandem mass tag; UCEC, uterine corpus endometrial carcinoma.

Received July 8, 2020, and in revised form, November 13, 2020  
Published, MCPRO Papers in Press, December 1, 2020, <https://doi.org/10.1074/mcp.TIR120.002216>

## REFERENCES

- Eng, J. K., Searle, B. C., Clauser, K. R., and Tabb, D. L. (2011) A face in the crowd: recognizing peptides through database search. *Mol. Cell. Proteomics* **10**, R111–009522
- Nesvizhskii, A. I. (2010) A survey of computational methods and error rate estimation procedures for peptide and protein identification in shotgun proteomics. *J. Proteomics* **73**, 2092–2123
- Chick, J. M., Kolippakkam, D., Nusinow, D. P., Zhai, B., Rad, R., Huttlin, E. L., and Gygi, S. P. (2015) A mass-tolerant database search identifies a large proportion of unassigned spectra in shotgun proteomics as modified peptides. *Nat. Biotechnol.* **33**, 743–749
- Kong, A. T., Leprevost, F. V., Avtonomov, D. M., Mellacheruvu, D., and Nesvizhskii, A. I. (2017) MSFragger: ultrafast and comprehensive peptide identification in mass spectrometry-based proteomics. *Nat. Methods* **14**, 513–520
- Nesvizhskii, A. I., Roos, F. F., Grossmann, J., Vogelzang, M., Eddes, J. S., Grissom, W., Baginsky, S., and Aebersold, R. (2006) Dynamic spectrum quality assessment and iterative computational analysis of shotgun proteomic data: toward more efficient identification of post-translational modifications, sequence polymorphisms, and novel peptides. *Mol. Cell. Proteomics* **5**, 652–670
- Ning, K., Fermin, D., and Nesvizhskii, A. I. (2010) Computational analysis of unassigned high-quality MS/MS spectra in proteomic data sets. *Proteomics* **10**, 2712–2718
- Chi, H., Liu, C., Yang, H., Zeng, W.-F., Wu, L., Zhou, W.-J., Wang, R.-M., Niu, X.-N., Ding, Y.-H., and Zhang, Y. (2018) Comprehensive identification of peptides in tandem mass spectra using an efficient open search engine. *Nat. Biotechnol.* **36**, 1059–1061
- Solntsev, S. K., Shortreed, M. R., Frey, B. L., and Smith, L. M. (2018) Enhanced global post-translational modification discovery with MetaMorpheus. *J. Proteome Res.* **17**, 1844–1851
- Na, S., Bandeira, N., and Paek, E. (2012) Fast multi-blind modification search through tandem mass spectrometry. *Mol. Cell Proteomics* **11**, M111.010199
- Dasari, S., Chambers, M. C., Slebos, R. J., Zimmerman, L. J., Ham, A.-J. L., and Tabb, D. L. (2010) TagRecon: high-throughput mutation identification through sequence tagging. *J. Proteome Res.* **9**, 1716–1726
- Avtonomov, D. M., Kong, A., and Nesvizhskii, A. I. (2018) DeltaMass: automated detection and visualization of mass shifts in proteomic open-search results. *J. Proteome Res.* **18**, 715–720
- An, Z., Zhai, L., Ying, W., Qian, X., Gong, F., Tan, M., and Fu, Y. (2019) PTMiner: localization and quality control of protein modifications detected in an open search and its application to comprehensive post-translational modification characterization in human proteome. *Mol. Cell. Proteomics* **18**, 391–405
- Shteynberg, D. D., Deutsch, E. W., Campbell, D. S., Hoopmann, M. R., Kusebauch, U., Lee, D., Mendoza, L., Midha, M. K., Sun, Z., and Whetton, A. D. (2019) PTMProphet: fast and accurate mass modification localization for the trans-proteomic pipeline. *J. Proteome Res.* **18**, 4262–4272
- da Veiga Leprevost, F., Haynes, S. E., Avtonomov, D. M., Chang, H. Y., Shanmugam, A. K., Mellacheruvu, D., Kong, A. T., and Nesvizhskii, A. I. (2020) Philosopher: a versatile toolkit for shotgun proteomics data analysis. *Nat. Methods* **17**, 869–870
- Creasy, D. M., and Cottrell, J. S. (2004) Unimod: protein modifications for mass spectrometry. *Proteomics* **4**, 1534–1536
- Tabb, D. L., Murugan, B. D., Okendo, J., Nair, O., Blackburn, J. M., Buthelezi, S. G., and Stoychev, S. (2019) Open search unveils modification patterns in formalin-fixed, paraffin-embedded thermo HCD and SCIEX TripleTOF shotgun proteomes. *Int. J. Mass Spectrom.* **448**, 116266
- Nielsen, N. S., Poulsen, E. T., Klintworth, G. K., and Enghild, J. J. (2014) Insight into the protein composition of immunoglobulin light chain deposits of eyelid, orbital and conjunctival amyloidosis. *J. Proteomics Bioinform.* **Suppl 8**, 002
- Nair, O. (2017) *Profiling Medulloblastoma and Juvenile Pilocytic Astrocytoma Brain Tumours in a South African Paediatric Cohort*. University of Cape Town, Cape Town, South Africa
- Zolg, D. P., Wilhelm, M., Schnatbaum, K., Zerweck, J., Knaute, T., Delanghe, B., Bailey, D. J., Gessulat, S., Ehrlich, H.-C., Weininger, M., Yu, P., Schlegl, J., Kramer, K., Schmidt, T., Kusebauch, U., et al. (2017) Building ProteomeTools based on a complete synthetic human proteome. *Nat. Methods* **14**, 259–262
- Edwards, N. J., Oberli, M., Thangudu, R. R., Cai, S., McGarvey, P. B., Jacob, S., Madhavan, S., and Ketchum, K. A. (2015) The CPTAC data portal: a resource for cancer proteomics research. *J. Proteome Res.* **14**, 2707–2713
- Mertins, P., Tang, L. C., Krug, K., Clark, D. J., Gritsenko, M. A., Chen, L., Clauser, K. R., Clauss, T. R., Shah, P., Gillette, M. A., Petyuk, V. A., Thomas, S. N., Mani, D. R., Mundt, F., Moore, R. J., et al. (2018) Reproducible workflow for multiplexed deep-scale proteome and phosphoproteome analysis of tumor tissues by liquid chromatography-mass spectrometry. *Nat. Protoc.* **13**, 1632–1661
- Yu, F., Teo, G. C., Kong, A. T., Haynes, S. E., Avtonomov, D. M., Geiszler, D. J., and Nesvizhskii, A. I. (2020) Identification of modified peptides using localization-aware open search. *Nat. Commun.* **11**, 4065
- Keller, A., Nesvizhskii, A. I., Kolker, E., and Aebersold, R. (2002) Empirical statistical model to estimate the accuracy of peptide identifications made by MS/MS and database search. *Anal. Chem.* **74**, 5383–5392
- Chang, H. Y., Kong, A. T., da Veiga Leprevost, F., Avtonomov, D. M., Haynes, S. E., and Nesvizhskii, A. I. (2020) Crystal-C: a computational tool for refinement of open search results. *J. Proteome Res.* **19**, 2511–2515
- Zhang, Y., Muller, M., Xu, B., Yoshida, Y., Horlacher, O., Nikitin, F., Garessus, S., Magdeldin, S., Kinoshita, N., Fujinaka, H., Yaoita, E., Hasegawa, M., Lisacek, F., and Yamamoto, T. (2015) Unrestricted modification search reveals lysine methylation as major modification induced by tissue formalin fixation and paraffin embedding. *Proteomics* **15**, 2568–2579
- Etherington, D. J., and Sims, T. J. (1981) Detection and estimation of collagen. *J. Sci. Food Agric.* **32**, 539–546
- Polasky, D. A., Yu, F., Teo, G. C., and Nesvizhskii, A. I. (2020) Fast and comprehensive N- and O-glycoproteomics analysis with MSFragger-Glyco. *Nat. Methods* **17**, 1125–1132
- Chung, H. S., Wang, S.-B., Venkatraman, V., Murray, C. I., and Van Eyk, J. E. (2013) Cysteine oxidative posttranslational modifications: emerging regulation in the cardiovascular system. *Circ. Res.* **112**, 382–392
- Sechi, S., and Chait, B. T. (1998) Modification of cysteine residues by alkylation. A tool in peptide mapping and protein identification. *Anal. Chem.* **70**, 5150–5158
- Schnatbaum, K., Zolg, D., Wenschuh, H., and Reimer, U. (2016) *Fast and accurate determination of cysteine reduction and alkylation efficacy in proteomics workflows*. JPT Application Note, Berlin, Germany

31. Bekker-Jensen, D. B., Kelstrup, C. D., Batth, T. S., Larsen, S. C., Haldrup, C., Bramsen, J. B., Sørensen, K. D., Høyer, S., Ørntoft, T. F., Andersen, C. L., and Others. (2017) An optimized shotgun strategy for the rapid generation of comprehensive human proteomes. *Cell Syst.* **4**, 587–599
32. Metz, B., Kersten, G. F. A., Hoogerhout, P., Brugghe, H. F., Timmermans, H. A. M., De Jong, A. D., Meiring, H., ten Hove, J., Hennink, W. E., Crommelin, D. J. A., and Others. (2004) Identification of formaldehyde-induced modifications in proteins reactions with model peptides. *J. Biol. Chem.* **279**, 6235–6243
33. Cordero, M. M., Houser, J. J., and Wesdemiotis, C. (1993) The neutral products formed during backbone fragmentations of protonated peptides in tandem mass spectrometry. *Anal. Chem.* **65**, 1594–1601
34. Savitski, M. M., Kjeldsen, F., Nielsen, M. L., and Zubarev, R. A. (2007) Relative specificities of water and ammonia losses from backbone fragments in collision-activated dissociation. *J. Proteome Res.* **6**, 2669–2673
35. Kumar, A., and Bachhawat, A. K. (2012) Pyroglutamic acid: throwing light on a lightly studied metabolite. *Curr. Sci.* **102**, 288–297
36. Dick, L. W., Jr., Kim, C., Qiu, D., and Cheng, K.-C. (2007) Determination of the origin of the N-terminal pyro-glutamate variation in monoclonal antibodies using model peptides. *Biotechnol. Bioeng.* **97**, 544–553
37. Reimer, J., Shamshurin, D., Harder, M., Yamchuk, A., Spicer, V., and Krokhin, O. V. (2011) Effect of cyclization of N-terminal glutamine and carbamidomethyl-cysteine (residues) on the chromatographic behavior of peptides in reversed-phase chromatography. *J. Chromatogr. A.* **1218**, 5101–5107
38. Sun, S., Yu, C., Qiao, Y., Lin, Y., Dong, G., Liu, C., Zhang, J., Zhang, Z., Cai, J., Zhang, H., and Bu, D. (2008) Deriving the probabilities of water loss and ammonia loss for amino acids from tandem mass spectra. *J. Proteome Res.* **7**, 202–208
39. Chelius, D., Jing, K., Lueras, A., Rehder, D. S., Dillon, T. M., Vizel, A., Rajan, R. S., Li, T., Treuheit, M. J., and Bondarenko, P. V. (2006) Formation of pyroglutamic acid from N-terminal glutamic acid in immunoglobulin gamma antibodies. *Anal. Chem.* **78**, 2370–2376
40. He, W., Tao, Y., and Wang, X. (2018) Functional polyamides: a sustainable access via lysine cyclization and organocatalytic ring-opening polymerization. *Macromolecules* **51**, 8248–8257
41. Leek, J. T., Scharpf, R. B., Bravo, H. C., Simcha, D., Langmead, B., Johnson, W. E., Geman, D., Baggerly, K., and Irizarry, R. A. (2010) Tackling the widespread and critical impact of batch effects in high-throughput data. *Nat. Rev. Genet.* **11**, 733–739
42. Rudnick, P. A., Clauser, K. R., Kilpatrick, L. E., Tchekhovskoi, D. V., Neta, P., Blonder, N., Billheimer, D. D., Blackman, R. K., Bunk, D. M., Cardasis, H. L., Ham, A.-J. L., Jaffe, J. D., Kinsinger, C. R., Mesri, M., Neubert, T. A., et al. (2010) Performance metrics for liquid chromatography-tandem mass spectrometry systems in proteomics analyses. *Mol. Cell. Proteomics* **9**, 225–241
43. Wang, X., Chambers, M. C., Vega-Montoto, L. J., Bunk, D. M., Stein, S. E., and Tabb, D. L. (2014) QC metrics from CPTAC raw LC-MS/MS data interpreted through multivariate statistics. *Anal. Chem.* **86**, 2497–2509
44. Paulovich, A. G., Billheimer, D., Ham, A.-J. L., Vega-Montoto, L., Rudnick, P. A., Tabb, D. L., Wang, P., Blackman, R. K., Bunk, D. M., Cardasis, H. L., and Others. (2010) Interlaboratory study characterizing a yeast performance standard for benchmarking LC-MS platform performance. *Mol. Cell. Proteomics* **9**, 242–254
45. Tabb, D. L., Vega-Montoto, L., Rudnick, P. A., Variyath, A. M., Ham, A.-J. L., Bunk, D. M., Kilpatrick, L. E., Billheimer, D. D., Blackman, R. K., Cardasis, H. L., and Others. (2009) Repeatability and reproducibility in proteomic identifications by liquid chromatography-tandem mass spectrometry. *J. Proteome Res.* **9**, 761–776
46. Clark, D. J., Dhanasekaran, S. M., Petralia, F., Pan, J., Song, X., Hu, Y., da Veiga Leprevost, F., Reva, B., Lih, T. M., Chang, H. Y., Ma, W., Huang, C., Ricketts, C. J., Chen, L., Krek, A., et al. (2019) Integrated proteogenomic characterization of clear cell renal cell carcinoma. *Cell* **179**, 964–983.e931
47. Dou, Y., Kawaler, E. A., Cui Zhou, D., Gritsenko, M. A., Huang, C., Blumenberg, L., Karpova, A., Petyuk, V. A., Savage, S. R., Satpathy, S., Liu, W., Wu, Y., Tsai, C. F., Wen, B., Li, Z., et al. (2020) Proteogenomic characterization of endometrial carcinoma. *Cell* **180**, 729–748.e726
48. Gillette, M. A., Satpathy, S., Cao, S., Dhanasekaran, S. M., Vasaikar, S. V., Krug, K., Petralia, F., Li, Y., Liang, W. W., Reva, B., Krek, A., Ji, J., Song, X., Liu, W., Hong, R., et al. (2020) Proteogenomic characterization reveals therapeutic vulnerabilities in lung adenocarcinoma. *Cell* **182**, 200–225.e235
49. Lenčo, J., Khalikova, M. A., and Švec, F. (2020) Dissolving peptides in 0.1% formic acid brings risk of artificial formylation. *J. Proteome Res.* **19**, 993–999
50. Palmisano, G., Parker, B. L., Engholm-Keller, K., Lendal, S. E., Kulej, K., Schulz, M., Schwämmle, V., Graham, M. E., Saxtorph, H., and Cordwell, S. J. (2012) A novel method for the simultaneous enrichment, identification, and quantification of phosphopeptides and sialylated glycopeptides applied to a temporal profile of mouse brain development. *Mol. Cell. Proteomics* **11**, 1191–1202