1  **Measuring proteomes with long strings: A new, unconstrained**
2  **paradigm in mass spectrum interpretation**

3

4  Arun Devabhaktuni[1], Niclas Olsson[1], Carlos Gonzales[1], Keith Rawson[1], Kavya Swaminathan[1],
5  Joshua E. Elias[1*]

6

7  [1]Department of Chemical and Systems Biology, Stanford School of Medicine, Stanford
8  University, Stanford, CA 94025, USA

9  *Corresponding Author: Dr. Joshua Elias, Clark Center W300C, 318 Campus Drive, Stanford,
10  CA 94305 (josh.elias@stanford.edu) (Phone: 650-724-3422) (Fax: 650-724-5791)

11

## Summary

Thousands of protein post-translational modifications (PTMs) dynamically impact nearly all cellular functions. Mass spectrometry is well suited to PTM identification, but proteome-scale analyses are biased towards PTMs with existing enrichment methods. To measure the full landscape of PTM regulation, software must overcome two fundamental challenges: intractably large search spaces and difficulty distinguishing correct from incorrect identifications. Here, we describe TagGraph, software that overcomes both challenges with a string-based search method orders of magnitude faster than current approaches, and probabilistic validation model optimized for PTM assignments. When applied to a human proteome map, TagGraph tripled confident identifications while revealing thousands of modification types on nearly one million sites spanning the proteome. We expand known sites by orders of magnitude for highly abundant yet understudied PTMs such as proline hydroxylation, and derive tissue-specific insight into these PTMs' roles. TagGraph expands our ability to survey the full landscape of PTM function and regulation.

## Introduction

Post translational modifications (PTMs) dynamically modulate the activity, conformation states, localization, interactions, abundance, and degradation of almost all proteins encoded by the human genome[1–4], yet most remain poorly understood.  PTM dysregulation has been linked to heart[5], neurodegenerative[6], and autoimmune[7] diseases, cancer[8], and countless other major health challenges[9].  Thus, characterizing PTMs' identities, abundances, and regulation is an essential dimension for understanding overall protein function and disease etiology.  However, mapping the full breadth of PTM identities and locations across the entire human proteome has remained intractable[10].

Mass spectrometry is arguably the most robust technology capable of direct, unambiguous, and large-scale PTM measurement.  It has provided transformative insight into the roles phosphorylation[11], acetylation[12], ubiquitylation[13] – both singly and in combination[14] -- have on cell biology. However, most global PTM studies have focused on modifications with optimized enrichment workflows[15].  Consequently, our current view of PTMs' collective impact on the human proteome is heavily skewed towards a small fraction of the possible PTM landscape[16].

Even without experimental enrichment, PTM-containing peptides are readily detected by routine tandem mass spectrometry (MS/MS) experiments[17,18], and are believed to comprise much of the "dark matter" in proteome datasets that consistently evades confident identification[19].

2

44    Conventional sequence database search tools cannot identify modified peptides unless they are

45    first anticipated by the researcher[20–22]. Search parameters including the number, kind, and

46    frequency of PTMs are usually chosen to strike a difficult compromise: considering larger

47    numbers of PTMs and other sequence variants is necessary for their identification, but doing so

48    exponentially increases the time needed to interpret MS/MS datasets, and  decreases the ability

49    to distinguish correct from incorrect assignments[23]. To partially address this compromise,

50    strategies have been proposed to constrain the number of proteins being searched, protease

51    specificity rules, or the allowable types and numbers of PTMs[17,18,24–26]. In practice, these

52    approaches only marginally decrease search times without clearly distinguishing correct from

53    incorrect PTM assignments[27]. Therefore, most have not been demonstrated on large, proteome-

54    scale datasets[23]

55    Here, we describe *TagGraph*, a powerful computational tool that addresses two principle

56    challenges of searching very large sequence spaces.  First, *TagGraph* leverages accurate *de*

57    *novo* mass spectrum interpretations[28,29] to rapidly search millions of possible sequences for a

58    match with an FM-index[30] data structure. This highly efficient search method makes modern

59    next-generation genome sequencing possible[31], but has not been adapted to proteomics. By

60    combining it with a graph-based string reconciliation algorithm, TagGraph rapidly searches

61    MS/MS datasets without restrictions on number of proteins, PTMs, or protease specificity. This

62    strategy achieves speeds orders of magnitude faster than prior algorithms because it considers

63    exponentially more sequence possibilities without having to explicitly test each one against input

64    spectra. Second, by replacing conventional "target-decoy" error estimation[32] with a PTM-

65    optimized probabilistic model, TagGraph accurately discovers and discriminates high-

66    confidence peptide identifications even from such large search spaces. We demonstrate 40-fold

67    more accurate false discovery rate estimation relative to target-decoy-dependent software.

68    Combined, these advances make large-scale, untargeted PTM proteomics possible. We

69    demonstrate this new search capability on a recently published human proteome draft[33]. Our

70    analysis reveals 2,576 modification types and 936,886 total modification sites spanning 13,791

71    proteins across 30 adult and fetal tissues (44,232 likely PTM sites across 5,576 proteins), and

72    with accurately estimated false discovery rates commensurate with current proteomics

73    standards.  Simultaneously assaying such a large number of modifications substantially

74    expanded the number of known PTM sites by orders of magnitude, particularly for those lacking

75    biochemical enrichment techniques. Our analysis reveals quantitative, functionally relevant,

76    differences in PTM stoichiometry between human tissues. We focus particularly on

77  hydroxylation, demonstrating this un-enrichable PTM's prevalence in the human proteome, its

78  potential role in cancer, and its association with candidate enzymes. Through this analysis, we

79  establish TagGraph as a paradigm shift for rapid proteome characterization, which promotes

80  simultaneous, unbiased PTM identification.


## Results

81  **Results**

## TagGraph: A new paradigm in fast, unrestricted proteome analysis

82  **TagGraph: A new paradigm in fast, unrestricted proteome analysis**

83  We developed TagGraph to address the compromise between search accuracy, depth, and

84  speed proteomics researchers commonly face when searching large tandem mass spectra

85  (MS/MS) datasets. With it, we can now accurately assign peptides bearing multiple unspecified

86  post-translational modifications (PTMs) or amino acid substitutions. Conventional database

87  search algorithms perform exponentially more comparisons between MS/MS spectra and

88  peptide candidates as they consider new modification types. In contrast, TagGraph rapidly

89  selects a very small number of candidate peptides from a sequence database through an

90  efficient string matching and reconciliation procedure (**Fig. 1a, Supplementary Fig. 1**). In so

91  doing, TagGraph effectively surveys very large sequence spaces that would be impractical to

92  query using traditional database search engines.

93  TagGraph leverages the speed of indexed string matching algorithms[34] by first transforming

94  complex, numeric mass spectra into discrete, unambiguous query strings using *de novo* peptide

95  sequencing. *De novo* sequencing algorithms produce long, reasonably accurate sequence

96  predictions from high resolution MS/MS spectra[29,35]: we found these predictions were over 50%

97  correct for nearly all interpretable MS/MS spectra (**Fig. 1b, Supplementary Fig. 2**)[29].

98  Consequently, we reasoned that many *de novo* peptides should contain a sub-string that

99  perfectly matches the true protein source of the observed MS/MS spectrum. The FM-index data

100  structure[36] was  developed to facilitate this kind of search.  TagGraph uses it to rapidly

101  assemble a small number of candidate peptide matches from an arbitrarily large, pre-indexed

102  sequence database with no restrictions on protease specificity, post-translational modifications,

103  or sequence variants. These candidates are then reconciled against the input *de novo*

104  sequence using a graph-based alignment algorithm that can discover and localize multiple

105  PTMs and other sequence alterations that co-occur on a single peptide sequence without

106  anticipating them *a priori* (**Supplementary Note 1**). Modification masses localized to individual

107  amino acid positions within a peptide are cross-referenced with the Unimod resource[10] to

108  suggest the modification's most likely identity based on mass and amino acid specificity. In this

4

109    way, TagGraph effectively searches all possible sequence alterations on time scales

110    commensurate with conventional database search tools.

111    Our strategy contrasts with prior approaches that extract many (>100) short sequence

112    fragments ("tags") from each input MS/MS spectrum to restrict protein candidates[37–39]. Although

113    they consider fewer peptides per mass spectrum than conventional database search algorithms,

114    they are subject to similar speed limitations if they consider large numbers of amino acid

115    modifications and variants.  Similarly, recent iterative approaches towards refining candidate

116    proteins and modifications[40], or non-specific database searches using very wide mass

117    tolerances[18,37] are subject to these limitations: their ability to identify modifications requires

118    comparison between spectra and large numbers of modified peptide candidates. As a result,

119    these approaches are prone to infeasibly long search times. We measured the advantage of

120    shifting this computational burden to TagGraph's string matching and reconciliation procedure

121    by comparing TagGraph's execution time to four algorithms designed to consider greatly

122    expanded search spaces[18,37,40–42]. We found that none could execute on both the entire data set

123    and the search space TagGraph considered in this comparison. Even by providing them with

124    reduced number of spectra, search spaces, or both, TagGraph's analysis speed was over an

125    order of magnitude greater than the next fastest algorithm (**Fig. 1c, Table 1**).

126    **Effective error estimation for modified peptides: beyond target-decoy**

127    Indexed string searches, as implemented by TagGraph, solve the long-standing conflict

128    between search speed and search depth, but present a second challenge: estimating reliable

129    false discovery rates (FDRs). The standard target-decoy estimation method we previously

130    developed[32] is unsuitable to TagGraph search results, since it loses discrimination accuracy as

131    more peptides and PTMs are considered (**Supplementary Note 2**)[27,43,44]. Consequently, we

132    developed a probabilistic validation strategy using a hierarchical Bayes Model optimized by

133    expectation maximization (EM)[45]. Our robust model is universal, deducing the likelihood that any

134    individual peptide-spectrum match is correctly interpreted, conditioned on fourteen quantitative

135    and categorical attributes (**Supplementary Note 3, Supplementary Fig. 3, Supplementary**

136    **Fig. 4**). Of these, half relate specifically to PTM-containing peptides, enabling discrimination

137    between correctly and incorrectly interpreted spectra, regardless of the deduced peptides'

138    modification states.

139    We first evaluated TagGraph's error model by comparing it to traditional target-decoy database

140    search (SEQUEST), using the cell line dataset described in **Fig. 1c**. We found that EM-

141    generated FDR estimates tended to be more conservative than those inferred from target-decoy

142     searches (**Fig. 1d**), as expected for a model that discriminates between correct and incorrect

143     PTM assignments. Furthermore, we found that the extent to which SEQUEST and TagGraph

144     disagreed was consistent with the estimated 1% FDR threshold we applied to both

145     (**Supplementary Fig. 5a**). For the majority of these disagreements however, TagGraph-

146     generated peptide-spectrum matches were far more consistent with correct identifications based

147     on protease specificity, algorithm-assigned scores, and ion assignment (**Supplementary Fig.**

148     **5b, Supplementary Fig. 6**).

149     To further evaluate TagGraph's error model, we sought to measure how often modifications are

150     miss-assigned to specific amino acid sites in the proteome. To accomplish this, we replaced all

151     tyrosine residues with phenylalanines (mass difference of one oxygen, 15.9995 Da) in an

152     altered human proteome sequence database (**Fig. 1e**). We reasoned that an accurate

153     expanded search algorithm should return phenylalanine-containing peptides with an additional

154     oxygen localized to converted phenylalanines; peptides containing converted phenylalanines

155     without the oxygen addition are incorrect. This approach should therefore serve to benchmark

156     modification assignments, which, similar to target-decoy's use without modifications, delivers an

157     expected known result, and could be generalized across multiple search engines.

158     We benchmarked four search methods against TagGraph with this validation tool. Each

159     algorithm's results were filtered based on target-decoy-based criteria (either the algorithm's own

160     implementation or a linear discriminant analysis[11]) or, for TagGraph, the hierarchical Bayes

161     model. The proportion of peptide-spectrum matches containing unmodified phenylalanines at

162     tyrosine positions was used to estimate the modification-specific FDR relative to the 1%

163     predicted FDR. Due to their reliance on target-decoy based statistics (**Supplementary Note 2**),

164     no algorithm besides TagGraph reliably discriminated phenylalanine-containing peptides:

165     TagGraph's error model was nearly an order of magnitude closer to the expected 1% than the

166     next-best flexible search method (**Fig. 1e**), increasing sensitivity by more than four-fold

167     (**Supplementary Table 1**) at 100 times the speed (**Fig. 1c**).

168     Even when searching the conventional human proteome sequence database, the four flexible

169     search methods above produced 'confident' results with readily identifiable errors at severely

170     underestimated FDR rates (**Supplementary Table 1, Supplementary Fig. 7**). In contrast,

171     TagGraph doubled the number of unique peptide identifications relative to SEQUEST (**Fig. 1f**)

172     by enabling accurate identification of peptides with any protease specificity and modification

173     state. Once reconciled with the Unimod resource[10], we found that unanticipated post-isolation

174     modifications accounted for the majority of this increase (83%), followed by biologically-

175    regulated modifications (12%) and those with no previous association (5%) (**Supplementary**

176    **Table 1**). Our analysis of this cell line demonstrated TagGraph's unique ability to sensitively

177    characterize modified peptides with speeds and accuracies that are compatible with current,

178    large-scale proteomic workflows.

## Unrestricted analysis of the human proteome reveals a broad modification landscape

181    To further investigate TagGraph's utility for deep PTM characterization, we next applied our

182    approach to a recently described draft human proteome[33], approximately 150 times larger than

183    our initial test data set. Due to their long computation times and underpowered validation

184    techniques, performing this analysis with preexisting database search methods would not have

185    been feasible.  We interpreted 25 million tandem mass spectra derived from 30 adult and fetal

186    tissues and over 2,000 raw data files[33] with TagGraph. Once *de novo* sequencing (PEAKS ver.

187    7[46]) was complete, searching these data with TagGraph collectively took just six days on a

188    single desktop computer.  These data yielded over 1.1 million unique peptides, tripling the

189    number originally reported using traditional database searching (**Fig. 2a, Supplementary Table**

190    **2**). This analysis identified proteins not found in the initial report, ranging from 100 (Adult CD8+

191    T Cells) to over 600 (Adult Gallbladder) additional proteins per tissue (**Supplementary Fig. 8a,**

192    **Supplementary Table 3**). Several of these were supported by histological staining

193    (**Supplementary Fig. 8b**).

194    As with our cell line analysis (**Fig. 1**), TagGraph predominantly rescued peptides bearing at

195    least one modification that was not considered in the original search (**Fig. 2b**). A small number

196    of post-isolation modifications (methionine oxidation; N-terminal carbamylation,

197    carbamidomethylation, and formylation) collectively accounted for 38% of modified spectra (**Fig.**

198    **2c, Table 2**), consistent with previous findings[17,18,37,47]. TagGraph rescued other commonly

199    disregarded peptide classes, including semi-specific and non-specific trypsin cleavage, and mis-

200    assigned monoisotopic precursor masses (**Fig. 2b**).

201    In comparison to the handful of abundant yet biologically irrelevant post-isolation modifications,

202    this extremely deep proteome analysis revealed a much wider array of lesser-abundant PTMs

203    (**Fig. 2c-d, Supplementary Table 4**). For example, we found N-terminal myristoylation, lysine

204    hydroxylation, and arginine dimethylation hundreds to thousands of times in the proteome

205    without requiring the kind of targeted, sample-intensive enrichment procedures that have

206    previously been essential to PTM analysis.  This study confirmed 4,278 modifications previously

207    reported in the Uniprot proteomics resource, while extending it by an additional 39,954 (**Fig. 2e,**

7

208 **Supplementary Table 5).** Comparing MS/MS spectra from this human proteome dataset to

209 spectra derived from synthetic peptides (**Supplementary Fig. 9**) served to validate several

210 unexpected, yet confidently identified peptides.

211 Many PTMs act as reversible switches on protein function. Their enzymatic addition and

212 removal regulates signaling networks, protein binding, and other cellular processes[1,3]. Although

213 more than 90% of TagGraph-identified PTMs were previously unreported, we found several

214 PTM-flanking sequence motifs[14,48,49] enriched in this dataset (e.g. proline-directed

215 phosphorylation[11,50] and glycine-directed arginine methylation[51]), supporting their validity (**Fig.**

216 **3a, Supplementary Fig. 10**). Furthermore, we identified over 200 gene ontologies[52] that were

217 significantly enriched among proteins bearing 22 noteworthy PTMs, giving additional support to

218 their validity and functional significance (**Fig. 3b, Supplementary Fig. 11, Supplementary**

219 **Table 6**). This unbiased analysis confirmed biological processes known to be regulated by

220 multiple PTMs (e.g., acetyl Lys, methyl Lys, phosphorylated Ser regulating chromatin

221 function[53]).  Other processes, such as the cell cycle, were associated with a much more

222 restricted set of PTMs (phosphorylated Ser)[54].

223 We found that most PTMs were enriched in multiple biological process or cellular compartment

224 were also implicated in multiple others.  For example, reversible arginine methylation

225 dynamically regulates proteins involved in RNA splicing and stabilization[50], as confirmed by our

226 ontology analysis (**Fig. 3b**). We observed a relative increase in the mono- and di-methylation

227 site abundances on RNA splicing proteins such as HNRNPA3 and SFPQ in reproductive tissues

228 and lymphocytes (**Fig. 3c, Supplementary Table 7**), suggesting that these modifications have

229 specific roles in these contexts. Chemically similar modifications like arginine mono- and

230 dimethylation showed stark contrasts: heat shock proteins were highly and consistently

231 methylated in this dataset, but were not readily identified in dimethylated states (**Fig. 3c,**

232 **Supplementary Table 7**).

233 **Quantifying PTM abundance and stoichiometry without requiring biochemical**
234 **enrichment**

235 We found that many PTMs' abundances across the 30 tissues examined here mirrored those of

236 the protein on which they were found, as exemplified by MBP R167 and HSPA8 R446

237 methylation (**Fig. 3c)**. This degree of congruity suggests consistent stoichiometry across

238 tissues. Conversely, other PTMs, including SFPQ R681 and R695 mono- and dimethylation,

239 were largely restricted to specific tissues, despite the proteins' uniform expression across the

240 entire dataset (**Fig. 3c**). Since PTMs and their host proteins can be simultaneously quantified by

241 mass spectrometry (e.g., **Fig. 3c**), we accordingly estimated each PTM's stoichiometry

242 (**Supplementary Methods, Supplementary Table 8**). This notion contrasts with previous PTM

243 stoichiometry assays which required metabolic labeling[55,56], or enzymatic removal of a single

244 target PTM class[56]. While such experimental interventions may estimate stoichiometries for a

245 single PTM class (e.g., phosphorylation), it is difficult to use them to compare multiple,

246 overlapping PTMs. Considering that a PTM's stoichiometry can have important implications for

247 its substrate protein's activity and function[57], deeply sequenced proteome datasets like this

248 stand to illuminate a wide range of protein regulation.

249 In support of our flexible stoichiometry estimation approach, we found that protein N-terminal

250 acetylation demonstrated the most consistently high stoichiometry (95.5%; stdev = 16.7%, **Fig.**

251 **3d**). This is expected, considering the broad and irreversible acetyl group addition, co-

252 translationally catalyzed by N-terminal acetyltransferases[58]. Conversely, we found that lysine

253 acetylation demonstrated consistently low and variable stoichiometry (15.2%; stdev = 22.7%,

254 **Fig. 3d**), consistent with its heterogeneous representation on histone proteins[59], and its possible

255 non-enzymatic origins on abundant cytosolic and mitochondrial proteins[60,61]. Over the entire

256 dataset, we found that neither PTM abundance nor stoichiometry correlated with substrate

257 protein abundance (**Supplementary Fig. 12**), supporting the complementary use of both

258 measurements in proteome characterization.

## TagGraph simultaneously characterizes multiple PTM types on highly modified proteins

261 TagGraph identified multiple PTMs that intersect on individual proteins, and on individual

262 residues (e.g., SFPQ R681, R695) (**Fig 3c, Supplementary Table 7**). Extreme examples

263 include Albumin (921 PTMs) and actin (514 PTMs) (**Fig. 3e**). Histones are also well understood

264 to undergo extensive and combinatorial modifications to encode epigenetic information[1].

265 However, deciphering these modifications has required individual histone isoform[62] or specific

266 modification[63] enrichment. Using TagGraph, we identified 277 PTMs across the major histone

267 proteins, 132 of which were not previously reported (**Supplementary Table 9**). While we found

268 modifications such as K28 dimethylation and K80 methylation on Histone H3 were both

269 abundant and ubiquitous across the tissues examined here (**Fig. 3f**), we note several tissue-

270 specific PTM combinations such as a 25-fold higher abundance of Histone H4 R56

271 dimethylation in fetal than adult tissues. Twenty-six PTMs, comprised of eleven PTM types,

272 showed similarly higher abundance in fetal tissues, suggesting specific roles in developmental

273 contexts (**Supplementary Table 9**). Our unbiased evaluation of these modifications, performed

274  in conjunction with the rest of the proteome and without targeted enrichment techniques, opens

275  new avenues to exploring tissue-specific epigenetic control.

## Enrichment-free PTM discovery identifies new roles for protein hydroxylation

277  We found that hydroxylation of prolines, tyrosines and lysines comprised a large (16%)

278  proportion of newly identified histone PTMs (**Supplementary Table 9**), yet only hydroxylated

279  tyrosine was previously described[64]. This coincides with our broader observation that several

280  modification classes remain uncharted across the human proteome, despite being highly

281  prevalent.  Proline hydroxylation, for example, is the most abundant modification in the human

282  body[65], yet just 171 sites have been recorded[66].  Unlike more widely studied modifications, no

283  enrichment tools exist to facilitate targeted hydroxylation analysis. Furthermore, of 11 amino

284  acids capable of becoming hydroxylated[10], four (Met, Trp, Phe, His) are often hydroxylated by

285  standard proteomics sample preparation protocols. Thus, true post-translational proline

286  hydroxylation must be distinguished from mis-localized artifacts[67].  Armed with TagGraph's

287  modification-focused error model, we confidently identified and localized 18-fold more hydroxyl

288  proline residues than were previously known in humans (**Table 2**).

289  Proline hydroxylation is best understood in the context of collagen proteins, as it is essential to

290  their role in maintaining extracellular matrix stability[65].  Despite hydroxyl proline comprising over

291  13% of mammalian collagen by weight[68], only 128 sites across all collagens were previously

292  assigned in humans (75% of all charted hydroxyl prolines in the human proteome).  TagGraph

293  identified 166 proline hydroxylation sites on COL1A2 alone, just three of which were previously

294  described, identified by Edman degradation[69](**Fig 4a**). While most proline hydroxylation sites

295  were highly represented across most solid tissues examined here (e.g., P330, P642), several

296  displayed tissue-specific abundance (e.g., P408, restricted to colon, bladder, liver, gallbladder,

297  and pancreas) (**Fig. 4b**).  TagGraph identified 25 other types of PTMs from this single protein,

298  suggesting multiple routes by which PTMs cooperatively regulate collagen structure and

299  function.

300  Although hydroxyl proline's role in maintaining collagen structure is well understood[65], its

301  prevalence and roles on other proteins has remained sparse[70–72]. Just 26 proteins were

302  previously reported to bear hydroxyl proline modifications besides collagens and collagen

303  domain-containing proteins[66], and thus, it has widely been considered a relatively specialized

304  PTM.  Our analysis extends known proline hydroxylation by nearly 3,000 sites spanning nearly

305  1,000 substrate proteins (**Supplementary Table 5**). These proteins were significantly enriched

10

306  for 113 biological processes (**Supplementary Table 6).** Thus, proline hydroxylation likely

307  shapes a diverse range of cellular processes beyond matrix homeostasis (**Fig. 3b**).

308  Noting that tumors also exploit multiple cellular processes during oncogenesis, we hypothesized

309  that proline hydroxylation could play a role in cancer.  Significant associations were previously

310  shown between specific phosphorylation sites and cancer-associated mutations[8]. Taking a

311  similar approach, we examined whether proline hydroxylation significantly intersected with

312  missense somatic cancer mutations catalogued in the COSMIC database[73]. We found that

313  hydroxylated prolines were 25% more likely to be associated with cancer mutations than

314  expected (p<6e-11, Fisher's exact test, **Fig. 4c**). Enrichment persisted even after excluding

315  collagen domain-containing proteins (22%, p<4e-6, **Fig. 4c**). Methionine oxidation, a common

316  post-isolation modification, was not enriched (p=0.49), nor were other post-isolation proline

317  modifications (**Fig. 4c**). Similar to phosphorylation, further study of hydroxylation could

318  substantially increase insight into cancer pathogenesis and reveal new therapeutic targets. The

319  identification of mutated hydroxylation sites on proteins which are hubs of post-translational

320  signaling (i.e., 20 such sites on histones) further supports this hypothesis (**Supplementary**

321  **Table 10**).

## Identifying candidate PTM interactors and regulators

323  As with proline hydroxylation, TagGraph significantly expanded the number of known sites for

324  lysine hydroxylation and asparagine hydroxylation by 18-fold and 14-fold, respectively. To

325  further elucidate protein-PTM interactions beyond ontological groupings, we reasoned that

326  PTMs should co-occur with the specific proteins with which they interact. We tested this notion

327  by screening all PTM and protein quantifications for significant correlations across the 30

328  tissues examined here. We found several protein-PTM correlations that confirmed known

329  functional associations (**Fig. 4d, Supplementary Fig. 13)**. Generally, we found that proteins

330  that were highly correlated with specific modifications did not bear those modifications

331  themselves (**Supplementary Fig. 13a**). However, they tended to be enriched for the same

332  functional ontologies as the PTMs' substrates (**Supplementary Fig. 13 c**). Such highly

333  correlated proteins are functionally associated with these PTMs and may be candidate PTM-

334  altering enzymes or indirect regulators.

335  We found over 70 proteins with abundances that correlated highly with lysine hydroxylation

336  abundance across all tissues (**Fig. 4d**). Many of these proteins, such as PXDN and CYGB have

337  known roles in oxygen transport or oxidoreductase activity (**Supplementary Fig. 13c**),

338    supporting their role in regulating hydroxylation PTMs.  Of note, one enzyme known to catalyze

339    this PTM, PLOD1[74], was the second most highly correlated (**Fig. 4d**).

340    Surprisingly, many proteins that correlated with lysine hydroxylation also correlated with

341    asparagine hydroxylation (**Fig. 4d**), despite no previous evidence linking these PTMs to the

342    same biological context. The primary substrates for lysine hydroxylation are collagens, which

343    are a major component of the extracellular matrix (ECM)[75]. Though asparagine hydroxylation

344    was not previously characterized in the ECM, TagGraph revealed 45 novel sites on Fibrillin-1

345    and Fibrillin-2 (**Supplementary Table 5**), both of which are ECM constituents. As opposed to

346    specific correlates (such as PLOD1 for lysine hydroxylation), proteins that correlate highly with

347    both PTMs may function as general positive regulators of ECM homeostasis through

348    hydroxylation (**Fig. 4d**). It is plausible that asparagine hydroxylation might stabilize fibrillins,

349    analogous to the function of lysine and proline hydroxylation in stabilizing collagen fibrils[65] and

350    asparagine hydroxylation in stabilizing ankyrins[76].

## Discussion

351    **Discussion**

352    Characterizing the identity, abundance, and function of post-translational modifications (PTMs)

353    is arguably the single-most important contribution mass spectrometry-based proteomics can

354    make to cell biology[16].  However, computational and experimental limitations have reinforced a

355    myopic view of a diverse and dynamic PTM landscape. TagGraph overcomes these obstacles

356    with two major innovations, making it possible to identify essentially any modified peptide

357    sequence from high-quality tandem mass spectra.

358    First, we circumvent the slow task of performing thousands or tens of thousands of peptide-

359    spectrum comparisons for each observed tandem mass spectrum. Instead, we devised an

360    extremely rapid string matching approach that produces only a handful of candidate sequence

361    matches that are then scored against the observed spectrum. A graph-based reconciliation

362    algorithm lets TagGraph consider any combination of modifications to a peptide.  Importantly, it

363    performs this task at speeds commensurate with conventional database search algorithms. Just

364    as similar string matching algorithms revolutionized next-generation DNA sequencing[31], we

365    expect this capability will become increasingly important as the latest generation of high-

366    resolution and high-volume mass spectrometers[77,78] become more widely available.

367    Second, we optimized a probabilistic model that simultaneously evaluates the likelihood that a

368    peptide's sequence and any modifications it bears are correct.  This component of our approach

369　was essential, since the standard "target-decoy" error estimation method we previously

370　developed is inherently blind to amino acid modifications[67,79]. We demonstrate the accuracy of

371　our error estimations using synthetic peptide validation, direct comparison to target-decoy, and

372　re-identification of known sites, motifs, protein-PTM relationships, and functional roles for

373　various PTMs. We expect the ease with which the model can be modified will offer superior

374　FDR estimation to target-decoy in several other specific applications[80,81].

375　Our analysis of an unprecedentedly rich PTM landscape could only be accomplished through

376　these two advances. This kind of high-throughput, unbiased PTM discovery is compatible with

377　any proteomics experiment using high-resolution tandem-mass spectrometry. We demonstrate

378　this ability in several ways, including a focused analysis on proline hydroxylation, a pervasive

379　PTM previously only known to occur on a small number of proteins. We expand the number of

380　known sites by 18-fold, linking proline hydroxylation to a much wider array of biological contexts

381　than originally thought, and demonstrate its significant association with somatic cancer

382　mutations. Furthermore, simultaneous identification of unmodified and modified peptides from

383　the same dataset enables high throughput quantification of PTM stoichiometries. This metric

384　holds great potential for illuminating functional relationships between PTMs, their substrates and

385　candidate regulatory proteins.

386　In addition to enabling routine, enrichment-free PTM discovery and characterization, we

387　envision many other applications for TagGraph. By searching MS/MS spectra in an enzyme-

388　independent manner, for example, TagGraph could automatically detect endogenous peptides

389　[82] and alternate start site utilization[83]. Furthermore, TagGraph's speed holds tremendous

390　advantages for refining gene predictions when applied in a proteogenomic context: by rapidly

391　evaluating multiple gene assemblies in all six frames while permitting amino acid substitutions,

392　insertions, and deletions it can validate translation products that would confound conventional

393　database search methods. This capacity could have direct application to systems that have

394　previously been intractable to large-scale proteome analysis such as the gut microbiome and

395　other complex microbial communities. Finally, by learning essential experimental details (e.g.,

396　PTMs, enzymatic digestion, and mass accuracy) directly from input data, TagGraph can help

397　standardize proteomic analyses. This capability stands to enable direct comparisons between

398　data sets collected by multiple laboratories, thereby fostering the kind of large-scale

399　collaborations that have transformed the genomics field.

**Acknowledgements**

**Author contributions**

A.D. designed and implemented all algorithms, designed and carried out all A375 mass spectrometry experiments, performed the human proteome analysis, and wrote the manuscript. J.E.E. designed algorithms and wrote the manuscript  N.O. performed synthetic peptide spectra analysis.  K.R. implemented TagGraph for web-based queries.  C.G. wrote the manuscript. . K.S. performed experimental validation studies and wrote the manuscript.
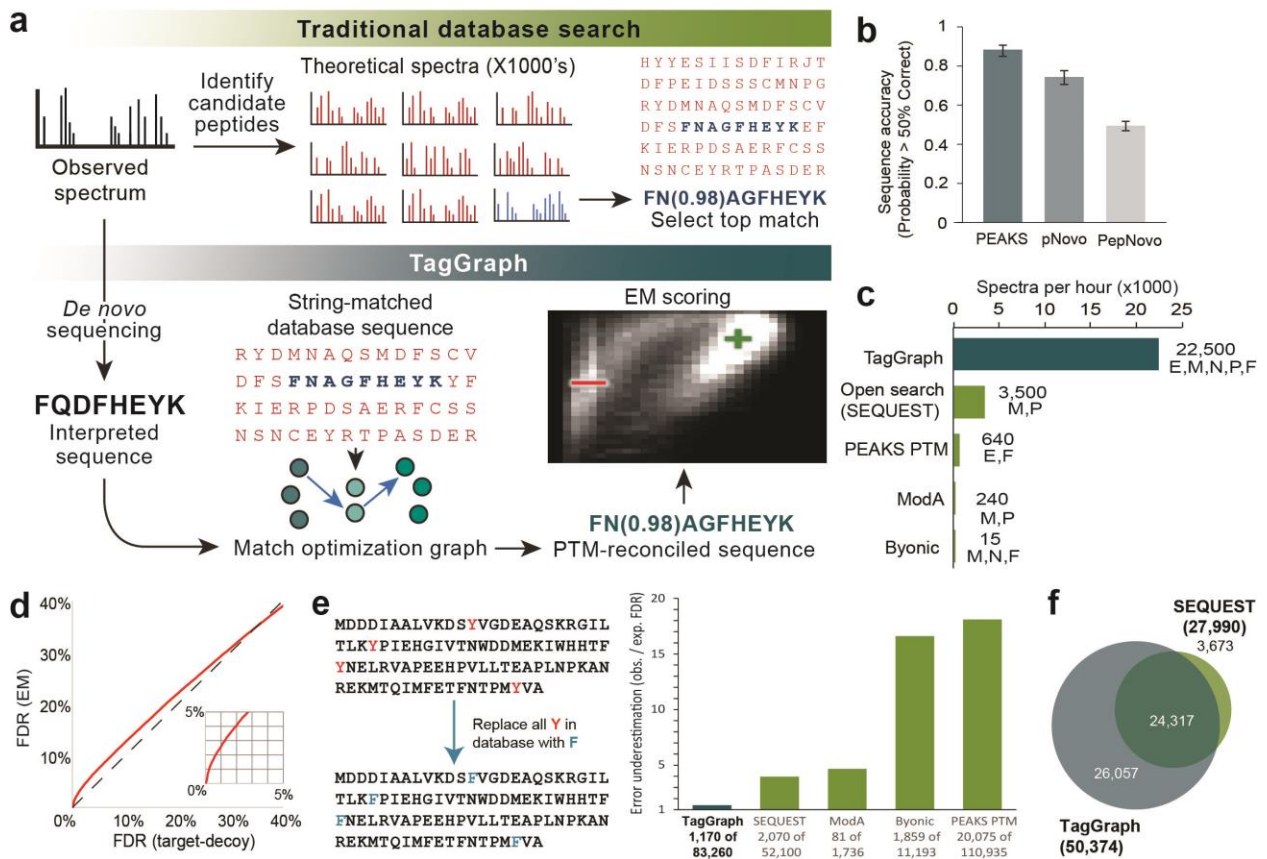
**FIGURES**

**Figure 1.**



**Figure 1. Through flexible string matching, TagGraph efficiently enables in-depth proteome characterization while controlling identification error. a)** TagGraph workflow. Traditional database search engines compare an observed MS/MS spectrum to hundreds or thousands of peptide candidates. In contrast, TagGraph first extracts a candidate peptide sequence from high resolution MS/MS spectra by *de novo* sequencing to facilitate rapid indexed protein database searching, and sequence reconciliation. This process lets TagGraph consider an unlimited number of PTMs and amino acid substitutions. A small number of top-scoring sequence candidates are ultimately scored against the input spectrum using an EM-optimized probabilistic Bayesian network. **b)** The majority of *de novo*-interpreted high-resolution MS/MS spectra are mostly correct. The proportion of analyzed spectra interpreted with over 50% sequencing accuracy by PEAKS [46], pNovo [84], and PepNovo [85] on a data set of 168,391 MS/MS spectra derived from the A375 melanoma cell line. Error bars correspond to the standard deviation in accuracy over different fractions from this data set. ROC curves corresponding with these graphs are described in **Supplementary Fig. 2b**. **c)** TagGraph search times on the A375 data set were at least an order of magnitude faster than previously described unconstrained modification and iterative search strategies, even when the latter were given comparatively reduced search spaces (**Table 1**). Letters indicate which search space expansions were compatible with the search algorithm: E, no enzyme specificity; M, any possible modification; N, any number of modifications per peptide; P, all proteins in sequence database. F indicates that the algorithm estimates a false discovery rate from its identifications. **d)** Expectation Maximization-based false discovery rate estimation is generally consistent with, but more

15

441 conservative than traditional target-decoy-based estimates when both are applied to TagGraph
442 results. This is expected considering target-decoy's inability to distinguish correct and incorrect
443 modification annotations **e)** The human proteome sequence database was modified, substituting
444 every tyrosine residue with a phenylalanine. The accuracy of PTM-specific false discovery rates
445 was estimated based on substituted phenylalanine-containing peptides reported by each
446 algorithm (**Methods**). Only TagGraph reported results with an empirically calculated FDR close
447 to 1%. Open search[18] was not included in this error rate comparison because it does not directly
448 localize modifications to specific residue positions. **f)** By allowing unrestricted modifications to
449 candidate peptides, TagGraph identified nearly twice as many unique peptide forms as
450 SEQUEST, configured with common search parameters. Analogous to the proteoform
451 concept[86], we define unique peptides by the combination of the peptide's amino acid sequence
452 and any modifications made to them.
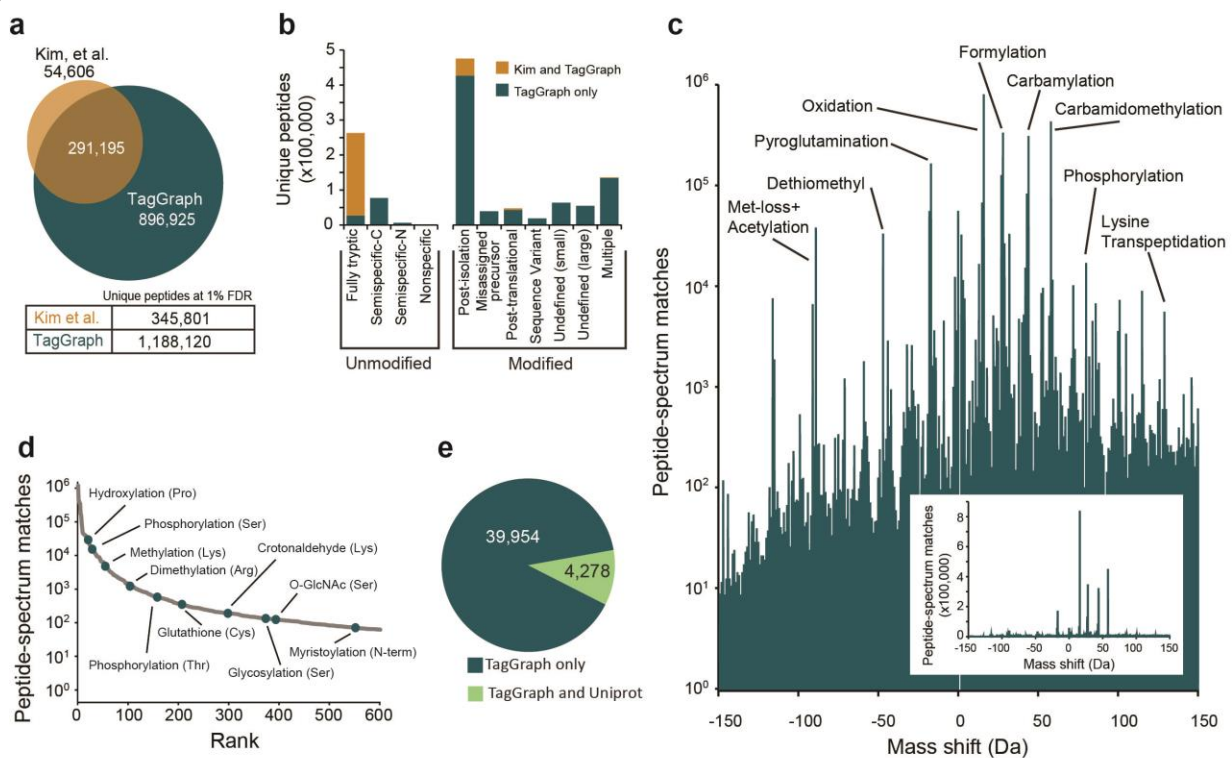
453 **Figure 2.**



454

455 **Figure 2.  TagGraph extends deep proteome characterization to post-translational
456 modifications.  a)** TagGraph confirmed the majority of identifications made by Kim et al.[33]
457 (**Supplementary Table 2**), but also expanded unique peptide identities from the human
458 proteome dataset over three-fold relative to those originally reported.  **b)** Categorical breakdown
459 of unique peptide forms (distinguishing PTMs) identified by TagGraph. As expected, the majority
460 of peptides identified by both TagGraph and Kim et al. correspond to tryptic peptides. Peptides
461 identified by TagGraph but not Kim et al. primarily originated from non-tryptic peptides and
462 peptides with unanticipated modifications. Post-isolation modifications comprised the most
463 prevalent identification category in this dataset. **c)** Mass shifts (modified amino acid mass –
464 unmodified amino acid mass) corresponding to all modifications identified by TagGraph from the
465 human proteome dataset reveal a complex modification landscape. Numbers of identifications
466 (peptide-spectrum matches) span six orders of magnitude. Despite the presence of several

467     highly abundant post-isolation modifications (e.g., formylation), the depth of the proteomic
468     profiling achieved in this dataset made it possible to characterize lower abundance post-
469     translational modifications. Inset: modification frequencies without log transformation. **d)** Ranked
470     relative abundances of 2,576 PTM-amino acid combinations, as estimated by the number of
471     spectra bearing each from the human proteome dataset.  Ten of these are highlighted; all
472     modifications are represented in **Supplementary Table 4**. **e)** TagGraph analysis of the human
473     proteome dataset identified 39,954 modification sites not present in Uniprot, 44,232 modification
474     sites total. The overlap in the sites reported by TagGraph and Uniprot is highly significant (p-
475     value < 1e-308, Fisher's exact test).
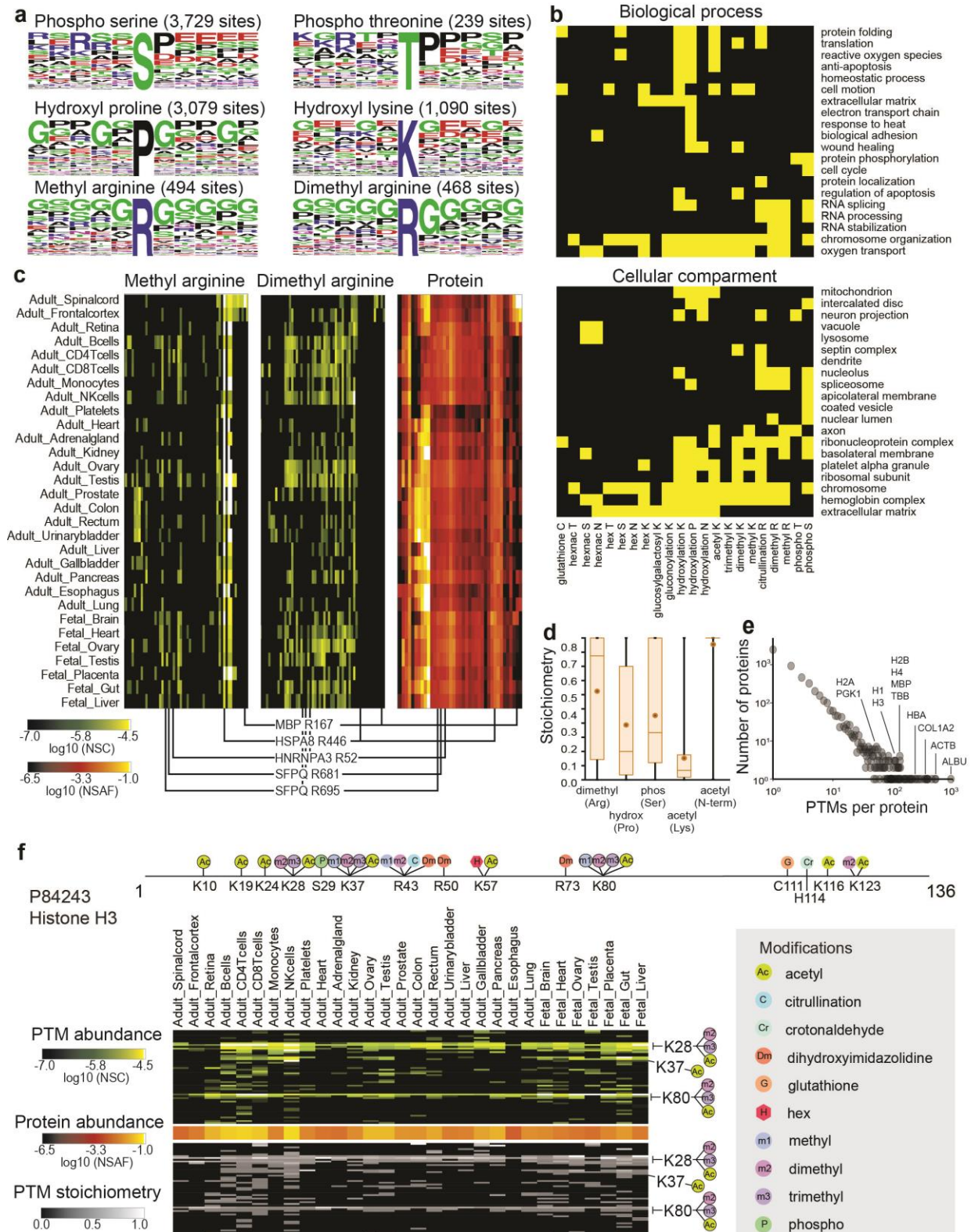
476    **Figure 3.**



477

**Figure 3. TagGraph reveals insights into PTM dynamics, function, and regulation. a)**
Sequence logos corresponding with select TagGraph-identified PTMs, as generated by

480    WebLogo[87]. Amino acids flanking the indicated PTMs were evaluated with the Motif-X[48]
481    algorithm, confirming previously characterized motifs while positing new ones (**Supplementary**
482    **Fig. 10**). **b)** Significantly enriched gene ontologies associated with prevalent post-translational
483    modifications (yellow, 1% FDR (Benjamini-Hochberg corrected)). Ontologies and significances
484    were assigned with the DAVID web tool[52]. A full list of all enriched ontologies is available in
485    **Supplemental Table 6**. Ontologies significantly enriched among post-isolation modifications
486    were excluded to correct for abundance-based biases in PTM detection (**Supplementary Fig.**
487    **11**). **c)** Arginine methylation and dimethylation distribution across proteins and tissues. The 64-
488    most abundant monomethylated or dimethylated Arg sites from the entire data set are displayed
489    across the y-axis, along with corresponding protein expression levels (49 proteins).  Three
490    modification sites on HNRNPA3 and SFPQ are highlighted for their distinct arginine
491    monomethylation and dimethylation patterns across the tissues, despite demonstrating near
492    uniform protein levels. Methyl modifications on MBP and HSPA8 are highlighted for their tissue
493    specificity and ubiquity (respectively). Proteins were ordered by hierarchical clustering.  PTMs
494    were arranged to match their substrate proteins. All methylation sites are reported in
495    **Supplementary Table 7. d)** Stoichiometry distributions vary for different PTMs, giving insight
496    into their regulation and function.  Box and whisker plots indicate the average (circle) and
497    median (horizontal bar) values, 25$^{th}$ quartile and 75$^{th}$ quartile (box), and minimum and maximum
498    (whiskers). **e)** Several proteins were found to be heavily modified in this data set. Histogram
499    shows the number of proteins identified with the indicated number of distinct PTMs (site and
500    modification). Of note, 921 distinct PTM sites were identified for human serum albumin. **f)**
501    TagGraph identified both known and novel PTM sites on Histone H3 (**Supplementary Table 9**);
502    a selection of the more abundant PTMs are shown. Site positions are numbered including the
503    initiating methionine, as is the convention in the Uniprot protein database. PTMs circled in black
504    are present in Uniprot. More detailed Histone PTM maps are presented in **Supplementary**
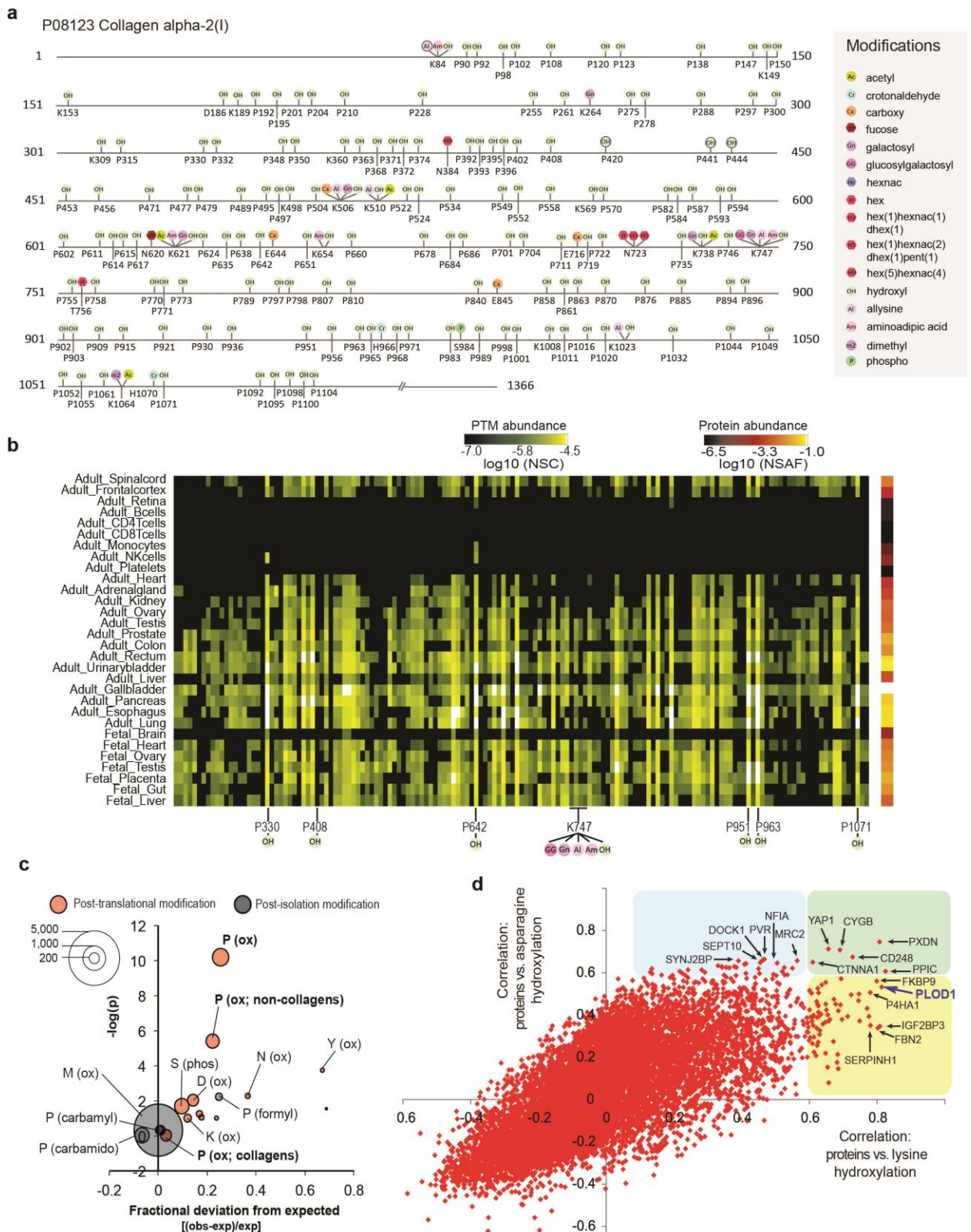505    **Table 9.**

19

**Figure 4.**

508       **Figure 4. Characterization of hydroxylation, an un-enrichable PTM, enabled by**
509   **TagGraph. a)** TagGraph extensively expanded known hydroxylation sites across the human
510   proteome; a selection of the most abundant PTMs on COL1A2 are shown as an example.
511   TagGraph expanded proline hydroxylation from three previously known sites (P420, P441,
512   P444)[69] to 166 while identifying 25 other types of PTMs on this protein (**Supplementary Table**
513   **5)**. **b)** PTMs identified in (a) were found to vary in abundance across tissues. Many
514   hydroxylations displayed uniform abundance across solid tissues (i.e., P330, P642), whereas
515   others displayed tissue-specific abundance variations (P408). **c)** Comparison between
516   modification and cancer mutation sites (COSMIC). The size of each bubble indicates the
517   number TagGraph-identified modification sites that were also found to be mutated in sequenced
518   tumors. The expected value and significance (Fisher's exact test) of this overlap was
519   determined from the background of all peptides confidently identified by TagGraph (**Methods**).
520   Proline hydroxylation sites significantly overlapped with mutation sites both overall ($p < 6e-11$)
521   and when restricted to non-collagen domain containing proteins ($p < 4e-6$), suggesting that
522   mutating these sites' PTM capacity plays a role in cancer pathogenesis. **d)** Correlations
523   between protein abundance and total PTM profiles across tissues (**Supplementary Fig. 13)**
524   suggest candidate regulatory enzymes and functional associations.  Proteins that highly
525   correlated with lysine hydroxylation (x-axis), asparagine hydroxylation (y-axis) or both are
526   highlighted (yellow, blue, or green, respectively). PLOD1 (in purple), the enzyme responsible for
527   lysine hydroxylation in collagen emerged among the proteins most correlated with this
528   modification. Protein expression levels were correlated with PTM stoichiometry across all
529   tissues (**Methods**).

530

531     **TABLES**

532     **Table 1.** **Search spaces considered by conventional and expanded**

533     **database search algorithms**

| Algorithm | Enzyme | Mods Considered | Mods/peptide | Protein Restriction |
|---|---|---|---|---|
| SEQUEST | LysC | methionine oxidation | three | **No** |
| Open search (SEQUEST) | LysC | Any mass between -500 and 500 Daltons | one (unlocalized) | **No** |
| PEAKS PTM | **None** | 435 previously known modifications (Unimod) | three | Yes |
| Byonic | LysC | Any mass between -40 and 100 Daltons | one | **No** |
| ModA | **None** | Any mass between -200 and 200 Daltons | **unlimited** | **No** |
| TagGraph | **None** | **Any mass** | **unlimited** | **No** |

534

**Table 2. Top 10 Post-isolation, post-translational, and previously uncharacterized amino acid modifications identified from the Kim et al. dataset.**

| Category[a] | Modification[b] | Specificity[c] | Mass Shift[d] | Unique Peptides[e] | Peptide-Spectrum Matches[f] | Sites[g] |
|---|---|---|---|---|---|---|
| **Post-isolation** | Carbamylation | Peptide N-term | +43.01 | 63,302 | 254,382 | 70,207 |
| | Formylation | Peptide N-term | +27.99 | 47,329 | 293,863 | 57,371 |
| | Carbamidomethylation | Peptide N-term | +57.02 | 42,310 | 304,898 | 52,899 |
| | Oxidation | Met | +15.99 | 42,193 | 784,001 | 41,115 |
| | Deamidation | Asn | +0.98 | 21,064 | 252,713 | 22,609 |
| | Acetaldehyde | Peptide N-term | +26.02 | 15,016 | 132,086 | 20,402 |
| | Deamidation | Gln | +0.98 | 8,911 | 36,897 | 13,015 |
| | Gln->pyroglutamate | Peptide N-term Gln | -17.03 | 8,372 | 123,277 | 7,966 |
| | Carbamidomethylation | Lys | +57.02 | 7,575 | 40,422 | 12,355 |
| | Carbamylation | Lys | +43.01 | 6,932 | 30,245 | 11,409 |
| **Post translational** | Phosphorylation | Ser | +79.97 | 3,466 | 15,798 | 3,729 |
| | Met-loss + Acetyl | Protein N-term Met | -89.03 | 2,391 | 39,799 | 1,705 |
| | Hydroxylation | Pro | +15.99 | 1,901 | 26,348 | 3,079 |
| | Citrullination | Arg | +0.98 | 1,551 | 5,328 | 2,020 |
| | GlyGly | Lys | +114.04 | 871 | 2,403 | 1,613 |
| | Allysine | Lys | -1.03 | 861 | 2,801 | 1,670 |
| | Hydroxylation | Lys | +15.99 | 766 | 3,316 | 1,090 |
| | Cyano | Cys | -32.03 | 601 | 2,591 | 567 |
| | Carboxylation | Glu | +43.98 | 467 | 726 | 913 |
| | Acetylation | Lys | +42.01 | 455 | 1,323 | 1,196 |
| **Previously undefined** | Unknown | Peptide N-term | +12.00 | 3,834 | 14,743 | 3,692 |
| | Unknown | Peptide N-term | +51.01 | 2,569 | 8,834 | 2,698 |
| | Iron(III)[h] | Asp, Glu | +52.92 | 2,552 | 9,843 | 1,468 |
| | carbamidomethyl and formyl on same residue[h] | Peptide N-term | +85.02 | 1,632 | 4,848 | 2,072 |
| | disulfide bondh | Cys with nearby Cys | -116.06 | 1,529 | 9,788 | 1,840 |
| | carbamidomethyl on C-terminus or Arg[h] | Peptide C-term Arg | +57.02 | 1,372 | 3,914 | 3,150 |
| | Unknown | Peptide N-term | +83.04 | 1,280 | 4,632 | 1,948 |
| | carbamidomethyl and pyro-glutamination on same residue[h] | Peptide N-term Glu | +39.01 | 1,099 | 3,622 | 908 |
| | Unknown | Peptide N-term | +23.98 | 1,059 | 2,919 | 1,381 |
| | reaction of N-terminal carbamidomethyl with internal Met[h] | Peptide N-term, co-occurs with dethiomethyl modification of internal Met | +105.02 | 957 | 4,907 | 991 |

a) Top ten modifications of the three indicated categories are shown, ordered by the number of unique peptides identified in the Kim *et al.* human proteome dataset. Categories assigned based on the likely modification identity, as determined by TagGraph.

b) Modification identities assigned by TagGraph, based on observed mass shifts, modification specificity, and evidence in the Unimod resource.

c) Specificity determined from the sites within modified peptdies to which observed mass shifts were assigned.

d) Mass shift measured from the difference between an observed amino acid residue's monoisotopic mass and the expected value. Negative values indicate a net mass loss.

e) Number of unique peptide sequences bearing the annotated modification. Does not include peptide sequences identified with multiple modifications.

f) Total number of spectra in which indicated modification was confidently identified.

g) Total number of distinct amino acid residue sites bearing indicated modification.

h) Hypothetical identity of mass shift

## References

538

539   1.   Jenuwein, T. & Allis, C. D. Translating the histone code. *Science* **293,** 1074–80
540         (2001).

541   2.   Eisenhaber, B. & Eisenhaber, F. Posttranslational modifications and subcellular
542         localization signals: indicators of sequence regions without inherent 3D structure?
543         *Curr. Protein Pept. Sci.* **8,** 197–203 (2007).

544   3.   Nussinov, R., Tsai, C.-J., Xin, F. & Radivojac, P. Allosteric post-translational
545         modification codes. *Trends Biochem. Sci.* **37,** 447–55 (2012).

546   4.   Grillari, J., Grillari-Voglauer, R. & Jansen-Dürr, P. Post-Translational Modification of
547         Cellular Proteins by Ubiquitin and Ubiquitin-Like Molecules: Role in Cellular
548         Senescence and Ageing. (2013).

549   5.   Smith, L. E. & White, M. Y. The role of post-translational modifications in acute
550         and chronic cardiovascular disease. *Proteomics. Clin. Appl.* **8,** 506–21 (2014).

551   6.   Rakhit, R. *et al.* Oxidation-induced Misfolding and Aggregation of Superoxide
552         Dismutase and Its Implications for Amyotrophic Lateral Sclerosis. *J. Biol. Chem.*
553         **277,** 47551–47556 (2002).

554   7.   Ryan, B. J., Nissim, A. & Winyard, P. G. Oxidative post-translational modifications
555         and their involvement in the pathogenesis of autoimmune diseases. *Redox Biol.* **2,**
556         715–24 (2014).

557   8.   Reimand, J., Wagih, O. & Bader, G. D. The mutational landscape of
558         phosphorylation signaling in cancer. *Sci. Rep.* **3,** 2651 (2013).

559   9.   Karve, T. M., Cheema, A. K., Karve, T. M. & Cheema, A. K. Small Changes Huge
560         Impact: The Role of Protein Posttranslational Modifications in Cellular
561         Homeostasis and Disease. *J. Amino Acids* **2011,** 1–13 (2011).

562   10.  Creasy, D. M. & Cottrell, J. S. Unimod: Protein modifications for mass
563         spectrometry. *Proteomics* **4,** 1534–6 (2004).

564   11.  Huttlin, E. L. *et al.* A tissue-specific atlas of mouse protein phosphorylation and
565         expression. *Cell* **143,** 1174–1189 (2010).

566   12.  Choudhary, C. *et al.* Lysine acetylation targets protein complexes and co-regulates
567         major cellular functions. *Science* **325,** 834–40 (2009).

568   13.  Kim, W. *et al.* Systematic and quantitative assessment of the ubiquitin-modified
569         proteome. *Mol. Cell* **44,** 325–40 (2011).

570   14.   Beltrao, P. *et al.* Systematic functional prioritization of protein posttranslational
571         modifications. *Cell* **150,** 413–25 (2012).

572   15.   Zhao, Y. & Jensen, O. N. Modification-specific proteomics: strategies for
573         characterization of post-translational modifications using enrichment techniques.
574         *Proteomics* **9,** 4632–41 (2009).

575   16.   Prabakaran, S., Lippens, G., Steen, H. & Gunawardena, J. Post-translational
576         modification: Nature's escape from genetic imprisonment and the basis for
577         dynamic information encoding. *Wiley Interdiscip. Rev. Syst. Biol. Med.* **4,** 565–583
578         (2012).

579   17.   Savitski, M. M., Nielsen, M. L. & Zubarev, R. A. ModifiComb, a new proteomic tool
580         for mapping substoichiometric post-translational modifications, finding novel
581         types of modifications, and fingerprinting complex protein mixtures. *Mol. Cell.*
582         *Proteomics* **5,** 935–48 (2006).

583   18.   Chick, J. M. *et al.* A mass-tolerant database search identifies a large proportion of
584         unassigned spectra in shotgun proteomics as modified peptides. *Nat. Biotechnol.*
585         **33,** 743–9 (2015).

586   19.   Skinner, O. S. & Kelleher, N. L. Illuminating the dark matter of shotgun proteomics.
587         *Nat. Biotechnol.* **33,** 717–718 (2015).

588   20.   Craig, R. & Beavis, R. C. TANDEM: matching proteins with tandem mass spectra.
589         *Bioinformatics* **20,** 1466–7 (2004).

590   21.   Eng, J. K., McCormack, A. L. & Yates, J. R. An approach to correlate tandem mass
591         spectral data of peptides with amino acid sequences in a protein database. *J. Am.*
592         *Soc. Mass Spectrom.* **5,** 976–989 (1994).

593   22.   Perkins, D. N., Pappin, D. J., Creasy, D. M. & Cottrell, J. S. Probability-based protein
594         identification by searching sequence databases using mass spectrometry data.
595         *Electrophoresis* **20,** 3551–67 (1999).

596   23.   Ahrné, E., Müller, M. & Lisacek, F. Unrestricted identification of modified proteins
597         using MS/MS. *Proteomics* **10,** 671–86 (2010).

598   24.   Nesvizhskii, A. I. A survey of computational methods and error rate estimation
599         procedures for peptide and protein identification in shotgun proteomics. *J.*
600         *Proteomics* **73,** 2092–123 (2010).

601   25.   Mann, M. & Wilm, M. Error-tolerant identification of peptides in sequence
602         databases by peptide sequence tags. *Anal. Chem.* **66,** 4390–9 (1994).

603   26.   Bern, M., Cai, Y. & Goldberg, D. Lookup peaks: a hybrid of de novo sequencing and
604         database search for protein identification by tandem mass spectrometry. *Anal.*
605         *Chem.* **79,** 1393–400 (2007).

606   27.   Gupta, N., Bandeira, N., Keich, U. & Pevzner, P. A. Target-decoy approach and
607         false discovery rate: when things may go wrong. *J. Am. Soc. Mass Spectrom.* **22,**
608         1111–20 (2011).

609   28.   Guthals, A., Clauser, K. R., Frank, A. M. & Bandeira, N. Sequencing-Grade De novo
610         Analysis of MS/MS Triplets (CID/HCD/ETD) From Overlapping Peptides. *J.*
611         *Proteome Res.* **12,** 2846–57 (2013).

612   29.   Devabhaktuni, A. & Elias, J. E. Application of de novo sequencing to large-scale
613         complex proteomics datasets. *Journal of Proteome Research* (2016).

614   30.   Ferragina, P. & Manzini, G. Opportunistic data structures with applications. in
615         *Proceedings 41st Annual Symposium on Foundations of Computer Science* 390–
616         398 (IEEE Comput. Soc, 2000). doi:10.1109/SFCS.2000.892127

617   31.   Langmead, B., Trapnell, C., Pop, M. & Salzberg, S. L. Ultrafast and memory-
618         efficient alignment of short DNA sequences to the human genome. *Genome Biol.*
619         **10,** R25 (2009).

620   32.   Elias, J. E. & Gygi, S. P. Target-decoy search strategy for increased confidence in
621         large-scale protein identifications by mass spectrometry. *Nat. Methods* **4,** 207–14
622         (2007).

623   33.   Kim, M.-S. *et al.* A draft map of the human proteome. *Nature* **509,** 575–81 (2014).

624   34.   Langmead, B., Trapnell, C., Pop, M. & Salzberg, S. L. Ultrafast and memory-
625         efficient alignment of short DNA sequences to the human genome. *Genome Biol.*
626         **10,** R25 (2009).

627   35.   Frank, A. M., Savitski, M. M., Nielsen, M. L., Zubarev, R. A. & Pevzner, P. A. De
628         Novo Peptide Sequencing and Identification with Precision Mass Spectrometry
629         research articles. *J. Proteome Res.* 114–123 (2007).

630   36.   Li, H. & Durbin, R. Fast and accurate long-read alignment with Burrows-Wheeler
631         transform. *Bioinformatics* **26,** 589–95 (2010).

632   37.   Na, S., Bandeira, N. & Paek, E. Fast multi-blind modification search through
633         tandem mass spectrometry. *Mol. Cell. Proteomics* **11,** M111.010199 (2012).

634   38.   Tabb, D. L., Saraf, A. & Yates, J. R. GutenTag: high-throughput sequence tagging
635         via an empirically derived fragmentation model. *Anal. Chem.* **75,** 6415–21 (2003).

636 39. Tsur, D., Tanner, S., Zandi, E., Bafna, V. & Pevzner, P. A. Identification of post-
637   translational modifications by blind search of mass spectra. *Nat. Biotechnol.* **23,**
638   1562–7 (2005).

639 40. Bern, M., Kil, Y. J. & Becker, C. Byonic: advanced peptide and protein identification
640   software. *Curr. Protoc. Bioinformatics* **Chapter 13,** Unit13.20 (2012).

641 41. Bern, M., Cai, Y. & Goldberg, D. Lookup peaks: a hybrid of de novo sequencing and
642   database search for protein identification by tandem mass spectrometry. *Anal.*
643   *Chem.* **79,** 1393–400 (2007).

644 42. Han, X., He, L., Xin, L., Shan, B. & Ma, B. PeaksPTM: Mass spectrometry-based
645   identification of peptides with unspecified modifications. *J. Proteome Res.* **10,**
646   2930–6 (2011).

647 43. Beausoleil, S. A., Villén, J., Gerber, S. A., Rush, J. & Gygi, S. P. A probability-based
648   approach for high-throughput protein phosphorylation analysis and site
649   localization. *Nat. Biotechnol.* **24,** 1285–92 (2006).

650 44. Nesvizhskii, A. I., Vitek, O. & Aebersold, R. Analysis and validation of proteomic
651   data generated by tandem mass spectrometry. *Nat. Methods* **4,** 787–97 (2007).

652 45. Keller, A., Nesvizhskii, A. I., Kolker, E. & Aebersold, R. Empirical Statistical Model
653   To Estimate the Accuracy of Peptide Identifications Made by MS/MS and Database
654   Search. *Anal. Chem.* **74,** 5383–5392 (2002).

655 46. Ma, B. *et al.* PEAKS: powerful software for peptide de novo sequencing by tandem
656   mass spectrometry. *Rapid Commun. Mass Spectrom.* **17,** 2337–42 (2003).

657 47. Creasy, D. M. & Cottrell, J. S. Error tolerant searching of uninterpreted tandem
658   mass spectrometry data. *Proteomics* **2,** 1426–34 (2002).

659 48. Schwartz, D. & Gygi, S. P. An iterative statistical approach to the identification of
660   protein phosphorylation motifs from large-scale data sets. *Nat. Biotechnol.* **23,**
661   1391–8 (2005).

662 49. Pearlman, S. M., Serber, Z. & Ferrell, J. E. A mechanism for the evolution of
663   phosphorylation sites. *Cell* **147,** 934–46 (2011).

664 50. Guo, A. *et al.* Immunoaffinity enrichment and mass spectrometry analysis of
665   protein methylation. *Mol. Cell. Proteomics* **13,** 372–87 (2014).

666 51. Thandapani, P., O'Connor, T. R., Bailey, T. L. & Richard, S. Defining the RGG/RG
667   Motif. *Molecular Cell* **50,** 613–623 (2013).

668  52.  Huang, D. W., Sherman, B. T. & Lempicki, R. A. Systematic and integrative analysis
669      of large gene lists using DAVID bioinformatics resources. *Nat. Protoc.* **4,** 44–57
670      (2009).

671  53.  Tan, M. *et al.* Identification of 67 histone marks and histone lysine crotonylation
672      as a new type of histone modification. *Cell* **146,** 1016–28 (2011).

673  54.  Fisher, D., Krasinska, L., Coudreuse, D. & Novák, B. Phosphorylation network
674      dynamics in the control of cell cycle transitions. *J. Cell Sci.* **125,** 4703–11 (2012).

675  55.  Olsen, J. V *et al.* Quantitative phosphoproteomics reveals widespread full
676      phosphorylation site occupancy during mitosis. *Sci. Signal.* **3,** ra3 (2010).

677  56.  Wu, R. *et al.* A large-scale method to measure absolute protein phosphorylation
678      stoichiometries. *Nat. Methods* **8,** 677–83 (2011).

679  57.  Khoury, G. A., Baliban, R. C. & Floudas, C. A. Proteome-wide post-translational
680      modification statistics: frequency analysis and curation of the swiss-prot database.
681      *Sci. Rep.* **1,** (2011).

682  58.  Starheim, K. K., Gevaert, K. & Arnesen, T. Protein N-terminal acetyltransferases:
683      When the start matters. *Trends in Biochemical Sciences* **37,** 152–161 (2012).

684  59.  Grunstein, M. Histone acetylation in chromatin structure and transcription.
685      *Nature* **389,** 349–352 (1997).

686  60.  Weinert, B. T. *et al.* Acetylation dynamics and stoichiometry in Saccharomyces
687      cerevisiae. *Mol. Syst. Biol.* **10,** 716 (2014).

688  61.  Wagner, G. & Hirschey, M. D. Nonenzymatic Protein Acylation as a Carbon Stress
689      Regulated by Sirtuin Deacylases. *Molecular Cell* **54,** 5–16 (2014).

690  62.  Garcia, B. a, Pesavento, J. J., Mizzen, C. a & Kelleher, N. L. Pervasive combinatorial
691      modification of histone H3 in human cells. *Nat. Methods* **4,** 487–489 (2007).

692  63.  Xie, Z. *et al.* Lysine succinylation and lysine malonylation in histones. *Mol Cell
693      Proteomics* **11,** 100–107 (2012).

694  64.  Huang, H., Sabari, B. R., Garcia, B. A., David Allis, C. & Zhao, Y. SnapShot: Histone
695      modifications. *Cell* **159,** (2014).

696  65.  Shoulders, M. D. & Raines, R. T. Collagen structure and stability. *Annu. Rev.
697      Biochem.* **78,** 929–58 (2009).

698  66.  The UniProt Consortium. UniProt: a hub for protein information. *Nucleic Acids Res.*
699      **43,** D204–12 (2014).

700   67.   Beausoleil, S. A., Villén, J., Gerber, S. A., Rush, J. & Gygi, S. P. A probability-based
701         approach for high-throughput protein phosphorylation analysis and site
702         localization. *Nat. Biotechnol.* **24,** 1285–92 (2006).

703   68.   NEUMAN, R. E. & LOGAN, M. A. The determination of hydroxyproline. *J. Biol.*
704         *Chem.* **184,** 299–306 (1950).

705   69.   FIETZEK, P. P., KUHN, K. & FURTHMAYR, H. Comparative Sequence Studies on
706         alpha2-CB2 from Calf, Human, Rabbit and Pig-Skin Collagen. *Eur. J. Biochem.* **47,**
707         257–261 (1974).

708   70.   Xie, L. *et al.* Oxygen-regulated beta(2)-adrenergic receptor hydroxylation by
709         EGLN3 and ubiquitylation by pVHL. *Sci. Signal.* **2,** ra33 (2009).

710   71.   Qi, H. H. *et al.* Prolyl 4-hydroxylation regulates Argonaute 2 stability. *Nature* **455,**
711         421–4 (2008).

712   72.   Masson, N., Willam, C., Maxwell, P. H., Pugh, C. W. & Ratcliffe, P. J. Independent
713         function of two destruction domains in hypoxia-inducible factor-alpha chains
714         activated by prolyl hydroxylation. *EMBO J.* **20,** 5197–206 (2001).

715   73.   Forbes, S. A. *et al.* COSMIC: Exploring the world's knowledge of somatic mutations
716         in human cancer. *Nucleic Acids Res.* **43,** D805–D811 (2015).

717   74.   Yamauchi, M. & Sricholpech, M. Lysine post-translational modifications of
718         collagen. *Essays Biochem.* **52,** 113–133 (2012).

719   75.   Yamauchi, M. & Sricholpech, M. Lysine post-translational modifications of
720         collagen. *Essays Biochem.* **52,** 113–133 (2012).

721   76.   Yang, M. *et al.* Asparagine and aspartate hydroxylation of the cytoskeletal ankyrin
722         family is catalyzed by factor-inhibiting hypoxia-inducible factor. *J. Biol. Chem.* **286,**
723         7648–60 (2011).

724   77.   Senko, M. W. *et al.* Novel parallelized quadrupole/linear ion trap/orbitrap tribrid
725         mass spectrometer improving proteome coverage and peptide identification
726         rates. *Anal. Chem.* **85,** 11710–11714 (2013).

727   78.   Andrews, G. L., Simons, B. L., Young, J. B., Hawkridge, A. M. & Muddiman, D. C.
728         Performance characteristics of a new hybrid quadrupole time-of-flight tandem
729         mass spectrometer (TripleTOF 5600). *Anal. Chem.* **83,** 5442–5446 (2011).

730   79.   Fermin, D., Walmsley, S. J., Gingras, A.-C., Choi, H. & Nesvizhskii, A. I. LuciPHOr:
731         algorithm for phosphorylation site localization with false localization rate
732         estimation using modified target-decoy approach. *Mol. Cell. Proteomics* **12,** 3409–

29

733        19 (2013).

734    80.    Woo, S. *et al.* Proteogenomic Database Construction Driven from Large Scale RNA-
735           seq Data. *J. Proteome Res.* **13,** 21–8 (2014).

736    81.    Verberkmoes, N. C. *et al.* Shotgun metaproteomics of the human distal gut
737           microbiota. *ISME J.* **3,** 179–89 (2009).

738    82.    Slavoff, S. A. *et al.* Peptidomic discovery of short open reading frame-encoded
739           peptides in human cells. *Nat. Chem. Biol.* (2012). doi:10.1038/nchembio.1120

740    83.    Ingolia, N. T., Ghaemmaghami, S., Newman, J. R. S. & Weissman, J. S. Genome-
741           wide analysis in vivo of translation with nucleotide resolution using ribosome
742           profiling. *Science* **324,** 218–23 (2009).

743    84.    Chi, H. *et al.* pNovo: de novo peptide sequencing and identification using HCD
744           spectra. *J. Proteome Res.* **9,** 2713–24 (2010).

745    85.    Frank, A. & Pevzner, P. PepNovo: de novo peptide sequencing via probabilistic
746           network modeling. *Anal. Chem.* **77,** 964–73 (2005).

747    86.    Smith, L. M. & Kelleher, N. L. Proteoform: a single term describing protein
748           complexity. *Nat. Methods* **10,** 186–7 (2013).

749    87.    Crooks, G. E., Hon, G., Chandonia, J.-M. & Brenner, S. E. WebLogo: a sequence
750           logo generator. *Genome Res.* **14,** 1188–90 (2004).

751

1    **Supplementary material for the manuscript:**

2

3    **Measuring proteomes with long strings: A new, unconstrained**
4    **paradigm in mass spectrum interpretation**

5

6    Arun Devabhaktuni[1], Niclas Olsson[1], Carlos Gonzales[1], Keith Rawson[1], Kavya Swaminathan[1],
7    Joshua E. Elias[1*]

8

9    [1]Department of Chemical and Systems Biology, Stanford School of Medicine, Stanford
10    University, Stanford, CA  94025, USA

11    *Corresponding Author: Dr. Joshua Elias, Clark Center W300C, 318 Campus Drive, Stanford,
12    CA 94305 (josh.elias@stanford.edu) (Phone: 650-724-3422) (Fax: 650-724-5791)

13

**SUPPLEMENTARY METHODS**

**1. Datasets**

*a. A375 data set*

<u>1.a.i. Sample processing</u>

A375 melanoma cells (ATCC) [1] were cultured in DMEM supplemented with 10% FCS and antibiotics. Cells were detached by trypsinization, pelleted, washed with PBS and flash frozen in liquid nitrogen. $5 \times 10^7$ flash-frozen A375 cells were thawed on ice and lysed by tip sonication in Urea lysis buffer (8 M Urea, 100 mM NaCl, 50 mM Tris, 1 mM PMSF, 10 μM E-64, 100 nM bestatin, pH 8.2). The cell lysate was reduced (5 mM DTT, 55 ºC, 30 min), alkylated (12.5 mM iodoacetamide, room temperature, 1 hr in the dark), and digested overnight with LysC at an enzyme:substrate ratio of 1:100 (37 ºC). The resulting peptide mixture was desalted using C-18 Sep-Pak cartridges (Waters), dried using vacuum centrifugation, and resuspended in 10 mM ammonium formate, pH 10 prior to high pH reverse phase (HPRP) separation. HPRP was performed using an Agilent 1100 binary HPLC, delivering a gradient (0%-5% B over 10 min, 5%-35% B over 60 min, 35%-70% B over 15 min, 70% B for 10 min) across an Agilent C-18 Zorbax Extend column. Buffer A was 10 mM ammonium formate, pH 10 and buffer B was 10 mM ammonium formate, 90% acetonitrile, pH 10. Sixty one-minute fractions were collected and concatenated into twelve fractions as described previously [2].

<u>1.a.ii. Mass Spectrometry</u>

All HPRP fractions were desalted using C-18 Sep-pak cartridges (Waters), vacuum centrifuged, and resuspended in 5% ACN, 5% formic acid at approximately 1 ug/ul. One microgram of each fraction was analyzed by microcapillary liquid chromatography electrospray ionization tandem mass spectrometry (LC-MS/MS) on an LTQ Orbitrap Velos mass spectrometer (Thermo Fisher Scientific, San Jose, CA) equipped with an in-house built nanospray source, an Agilent 1200 Series binary HPLC pump, and a MicroAS autosampler (Thermo Fisher Scientific). Peptides were separated on a 125 um ID x 18 cm fused silica microcapillary column with an in-house pulled emitter tip with an internal diameter of approximately 5 um. The column was packed with ProntoSIL $C_{18}$ AQ reversed phase resin (3 um particles, 200Å pore size; MAC-MOD, Chadds Ford, PA). Each sample was separated by applying a two-step gradient: 7% -25% buffer B over 2h; 25-45% B over 30 min. Buffer A was 0.1% formic acid, 2.5% ACN and buffer B was 0.1% formic acid, 97.5% ACN. The mass spectrometer was operated in a data dependent mode in which a full MS scan was acquired in the Orbitrap (AGC: $5 \times 10^5$; resolution: $6 \times 10^4$; m/z range: 360-1600; maximum ion injection time, 500 ms), followed by up to 10 HCD MS/MS spectra,

2

47    collected from the most abundant ions from the full MS scan. MS/MS spectra were collected in

48    the Orbitrap (AGC: $2x10^5$; resolution: $7.5x10^3$; minimum m/z: 100; maximum ion injection time,

49    1000 ms; isolation width: 2 Da; normalized collision energy: 30; default charge state: 2;

50    activation time: 30 ms; dynamic exclusion time: 60 sec; exclude singly-charged ions and ions for

51    which no charge-state could be determined).The mass calibration of the Orbitrap analyzer was

52    maintained to deliver mass accuracies of ±5 ppm without an external calibrant. All raw mass

53    spectrometry data were uploaded to the PRIDE repository[3] and assigned the reference ID

54    PXD######.

55        *1.b. Synthetic peptide confirmation data set*

56    In total, 86 synthetic peptides (SpikeTides from JPT Peptide Technologies GmbH) validating

57    various modifications, semi- or non-specific peptide assignments were evaluated. A final pool of

58    all peptides dissolved in 0.1% formic acid was generated and the concentration of each

59    individual peptide was roughly 250 fmol/µL. Two LC-MS/MS runs were performed and the auto-

60    sampler injected 1 µl of the synthetic peptide pool. An ESI-Orbitrap Elite mass spectrometer

61    (Thermo Electron, Bremen, Germany) interfaced with an Eksigent ekspert nanoLC 425 system

62    (Eksigent technologies, Dublin, CA, USA) was used. Peptides were introduced into the mass

63    spectrometer via a fused silica microcapillary column (100 µm inner diameter) ending in an in-

64    house pulled needle tip. The columns were packed in-house to a length of 18 cm with a C18

65    reversed-phase resin (with Reprosil-Pur C18-AQ resin (3 µm Dr. Maisch, GmbH, Germany). For

66    elution a two-step gradient of 4-25% buffer B (5 % DMSO, 0.2% formic acid and 94.8 %

67    acetonitrile (v/v)) in buffer A (5 % DMSO, 0.2% formic acid in water (v/v)) over 60 min followed

68    by a second phase of 25-45% buffer B over 20 min was used. The LTQ-Orbitrap was operated

69    in data-dependent mode to automatically switch between Orbitrap-MS (from *m/z* 340 to 1600)

70    and ten MS/MS acquisition. Each FT-MS scan was acquired at 60,000 FWHM nominal

71    resolution settings while the MS/MS spectra were acquired using HCD and at a resolution of

72    15,000. Precursor ion charge state screening was enabled (charge state 1 rejected) and the

73    normalized collision energy was set to 35%.

74    The resulting data were analyzed by TagGraph to match mass spectra with their best-matching

75    synthetic peptide sequence. Synthetic-derived and experiment-derived mass spectra (e.g., from

76    the draft human proteome data set) were only compared if spectra were assigned to the same

77    peptide sequence in the same charge state. Of the 86 peptides synthesized, 75 were matched

78    in this manner and used to validate TagGraph peptide-spectrum assignments (**Supplementary**

79    **Fig. 9**). The mass tolerance used to match b- and y-ions was 0.1 Daltons.

80      *1.c. Draft human proteome data set*

81    All RAW data and database search results from the draft human proteome data set [4] were

82    downloaded from the PRIDE[5] data repository using accession PXD000561.

83    **2. *De novo* search engine comparisons**

84    We compared the performance of three *de novo* sequencing algorithms, PEAKS 7 [6], PepNovo+

85    (ver. 3.1) [7], and pNovo (ver 1.1) [8] and assayed their ability to generate *mostly* correct sequence

86    interpretations of high mass accuracy MS/MS spectra[9]. Each algorithm was used to search the

87    A375 dataset, and the resulting peptide identifications were compared against high confidence

88    peptide spectrum matches obtained from a SEQUEST search of the same dataset as previously

89    described[9](**3***,* below). Static and differential modifications PTM were set as for the SEQUEST search.

90    Mass tolerance parameters were optimized to achieve maximum sequencing accuracy for each

91    algorithm individually[9]. PEAKS was run with a precursor mass tolerance of 10 ppm and a

92    fragment mass tolerance of 0.01 Da. PepNovo+ was run a 0.01 Dalton precursor mass

93    tolerance and 0.05 Da mass tolerance on fragment ions. pNovo was run with a 6 ppm precursor

94    mass tolerance and a 25 ppm fragment ion tolerance.

95    The accuracy of each *de novo* algorithm was assessed using the **sequence accuracy** metric

96    (**Supplementary Fig. 2a**)[9]. For a given *de novo* peptide-spectrum match and its corresponding

97    high confidence SEQUEST peptide-spectrum match, sequence accuracy represents the fraction

98    of prefix residue masses[10] present in the high confidence SEQUEST match which were also

99    present in the *de novo* sequence[9].

100    **3. Database search engine comparisons**

101    We assessed several expanded database search algorithms' abilities to detect undefined

102    modifications, without constraining protease specificity, using the A375 dataset. As a baseline,

103    we searched all 168,391 MS/MS spectra in this dataset with SEQUEST [11] (ver. 28 rev 12) using

104    an indexed sequence database comprised of the human proteome (Uniprot, downloaded

105    12/9/2014) [12] plus common contaminants. The database was concatenated with a reversed

106    database for target-decoy FDR estimation. The SEQUEST search was conducted with LysC

107    protease specificity, 50 ppm precursor ion mass tolerance, and 0.5 Da fragment ion mass

108    tolerance. Cysteine carbamidomethylation (+57.021464 Da) was set as a static modification and

109    methionine oxidation (+15.994915 Da) was set as a differential modification.

110    The A375 dataset was then searched using PEAKS PTM (ver. 7) [13], Byonic (ver. 2.5.6) [14], ModA

111    (ver. 1.03) [15], and the Open search method using SEQUEST [16] – four strategies described as

4

112    being able to either consider relatively large numbers of discrete amino acid modifications, or

113    searching spectra with no *a priori* constraints on possible modifications. It was not possible to

114    search the entire A375 dataset with any of the above algorithms using completely unconstrained

115    parameters with respect to both modifications and protease specificity: either the algorithms

116    would not execute, or they did not complete within a reasonable amount of time (5 days per

117    RAW data file). Thus, CPU times were calculated using the most feasible parameters

118    approximating such a search for each algorithm, and extrapolated from a limited subset of

119    search results in cases where searching the entire dataset would be too computationally

120    intensive. As such, the search times reported in **Fig. 1d** represent a substantial underestimate

121    of the true time needed for each of these algorithms to analyze a sample in a manner equivalent

122    to TagGraph, as described below.

123    For PEAKS PTM, the A375 dataset was first *de novo* sequenced using the following settings: 10

124    ppm precursor ion tolerance and 0.01 Da fragment ion tolerance, cysteine

125    carbamidomethylation as a static modification, and methionine oxidation as a differential

126    modification. The dataset was then analyzed with PEAKS PTM using the same modification and

127    mass tolerances as for the *de novo* sequencing, LysC enzyme specificity allowing for

128    nonspecific cleavage at both ends of the peptide, and additionally searching with all 485

129    modifications curated in PEAKS's internal PTM database. For ModA, the A375 dataset was

130    analyzed with the following settings: 0.05 Da precursor mass tolerance, 0.05 Da fragment ion

131    tolerance, allowing one modification per peptide, no protease specificity, modification size

132    between -200 Da and 200 Da. For Byonic, the dataset was analyzed using the following

133    settings: 10 ppm precursor ion tolerance, 20 ppm fragment ion tolerance, LysC protease

134    specificity, cysteine carbamidomethylation as a static modification, methionine oxidation as a

135    differential modification, wildcard search enabled with a minimum mass of -200 Daltons and a

136    maximum mass of 200 Daltons. For all three algorithms, the sequence database used was the

137    same as for SEQUEST. PEAKS PTM and Byonic were allowed to use their own internal

138    methods for creating decoy sequences, whereas ModA was given a concatenated

139    forward/reversed database as input. The open search method was conducted using SEQUEST

140    and the same indexed database as above, with a 500 Da tolerance on the precursor ion mass

141    and a 0.1 Da mass tolerance on fragment ion masses.

142    For the search time comparison between all algorithms (**Fig. 1d**), CPU times were calculated as

143    the sum of the CPU time over all processes spawned by each database search algorithm to

144    analyze the data. Due to computational constraints, it was not possible to run Byonic with a

145 semi specific or nonspecific enzyme specificity, or to search the entire A375 dataset with either

146 Byonic or ModA. Thus, we conducted the Byonic search with full LysC specificity, and analyzed

147 only a single fraction of the A375 dataset for both Byonic and ModA. The estimated CPU time

148 over the entire dataset was extrapolated by multiplying the CPU time recorded from the analysis

149 of a single HPRP fraction by the ratio of the total MS/MS spectra in the dataset (168,391)

150 divided by the number of MS/MS spectra in the single fraction (16,613).

151 **4. TagGraph Parameters**

152 TagGraph was used to analyze both the A375 dataset and the draft human proteome data set.

153 In both cases, all available MS/MS were first *de novo* sequenced using PEAKS. The resulting

154 peptide sequences and raw mass spectra (mzXML-formatted [17]) were given as input to

155 TagGraph.

156 For the A375 dataset, PEAKS was run as described above. For the draft Human Proteome data

157 set, PEAKS was run with a 10 ppm precursor mass tolerance, 0.05 Dalton fragment ion

158 tolerance, cysteine carbamidomethylation as a static modification, and methionine oxidation as

159 a differential modification to maximize sequence accuracy.

160 The *de novo* sequencing results were searched with TagGraph against the human proteome

161 (Uniprot, downloaded 12/9/2014) plus common contaminants. The database was concatenated

162 with reversed decoy sequences only for searches of the A375 dataset.  This enabled direct

163 comparisons of FDR estimates derived from TagGraph with those derived from target-decoy

164 searching, and fair comparison of the CPU time of TagGraph with the other database search

165 algorithms. For the draft human proteome data set, only the forward database was searched as

166 FDR estimates were derived using the hierarchical Bayes scoring model. For both datasets,

167 mass tolerances were set to 10 ppm precursor ion tolerance and 0.1 Dalton fragment ion

168 tolerance. Enzyme specificity was set to LysC for the A375 dataset and Trypsin for the human

169 proteome data set. Although enzyme specificity was considered as a scoring attribute in the

170 hierarchical Bayes model, TagGraph is able to return high-confidence semi specific and

171 nonspecific peptide-spectrum matches regardless of the input enzyme specificity.

172 High confidence results were retrieved at a 1% FDR for the A375 dataset by ranking all returned

173 peptide-spectrum matches according to their probabilities $P(D|+)$ from highest to lowest, then

174 adding matches in order of decreasing rank to the set of high confidence results until the

175 expected FDR equaled 1% (**Supplementary Note 3, Equation 2**). The human proteome data

6

176   set was evaluated in a similar manner, considering each experiment (e.g., gel fractionation of

177   adult heart tissue) individually rather than for the entire dataset as a whole.

## 5. PTM FDR estimation with amino acid-substituted proteome

179   Target-decoy based error estimation accuracy declines when applied to peptide modifications

180   and other large search spaces [18] (**Supplementary Note 2**). Despite this, all previously

181   described unrestricted search algorithms rely on target-decoy to delineate sets of confidently

182   identified spectra. To assess the degree to which predicted FDRs produced by these algorithms

183   deviate from the actual FDR, we employed a modified human proteome sequence database in

184   which every tyrosine residue was replaced by a phenylalanine. The mass difference between

185   these residues (15.994915 Da) corresponds with an oxygen atom, and is a frequently observed

186   modification on several residues (e.g., methionine), while distinguishing other unmodified

187   residues (e.g., alanine and serine).  Thus, we reasoned that search engines capable of accurate

188   PTM assignment and discrimination should search the tyrosine-substituted database and return

189   phenylalanine-containing peptides modified by an oxygen only on those phenylalanines that

190   were previously tyrosines.  They should be able to discriminate these identifications from

191   erroneous ones in which oxidation modifications were assigned to unaltered residues.  We

192   analyzed the A375 dataset against this modified sequence database using SEQUEST, PEAKS

193   PTM, Byonic, and ModA. The results from each algorithm were then filtered to a 1% predicted

194   FDR using target-decoy based statistics. Byonic and PEAKS PTM were allowed to use their

195   own internal target-decoy based filtering algorithms. Search results provided by SEQUEST and

196   ModA were filtered using a linear discriminant method [19]. The FDR for each set of search results

197   was calculated as the proportion of peptide-spectrum matches containing phenylalanines at

198   tyrosine positions which were not annotated with a phenylalanine to tyrosine modification

199   (+15.9995 Da) or a modification of equivalent mass.

200   The SEQUEST search was conducted with phenylalanine to tyrosine, methionine hydroxylation,

201   proline hydroxylation, lysine hydroxylation, and asparagine hydroxylation as differential

202   modifications. Aside from this change, all searches were conducted with the same parameters

203   and data as described in section 3 above. Open search was not considered in this comparison

204   as it does not provide amino acid localizations for its predicted modification masses.

205   The A375 dataset was also analyzed with TagGraph against the modified human sequence

206   database. The parameters used were identical to those described in section 4 above. Results

207   were filtered to a 1% predicted FDR using the hierarchical Bayes model and the actual FDR

208   was calculated as described above.

209 **6. Abundance calculations**

210      *6b. Protein abundance calculation*

211   Protein abundances were calculated using the distributed normalized spectral abundance factor

212   (NSAF) method [20]. Briefly, the number of spectral counts originating from peptides that uniquely

213   map to single proteins were summed over all proteins identified in an experiment. Spectral

214   counts recorded from peptides that map to multiple proteins were distributed across all such

215   proteins according to the proportion of spectral counts assigned to them from uniquely mapped

216   peptides. Finally, summed spectral counts for each protein were normalized by protein length,

217   and the sum of all protein abundances for each experimental dataset was normalized to one.

218   Protein abundances per tissue were calculated as the average of the individual NSAF for that

219   protein over all experiments performed on that tissue.

220      *6b. Site abundance and stoichiometry calculations*

221   To compare modification sites between tissues, we quantified the abundance of sites using two

222   methods: normalized spectral counts (NSC) and estimated stoichiometry. For both methods, we

223   first generated a catalog of all confident peptide identifications that span a given modified amino

224   acid position of a protein, regardless of modification state. The total spectral counts

225   corresponding to all peptides containing the amino acid position ($S_T$) and just those

226   corresponding to peptides containing the exact modification on the site of interest ($S_m$) were

227   calculated for each experimental dataset.

228   The normalized spectral count of a modification site is calculated as $S_m$ divided by the number

229   of confidently identified peptide-spectrum matches in the experimental dataset. The

230   stoichiometry of the modification is calculated as $S_m$ divided by $S_T$. Modification site

231   abundances (stoichiometry or NSC) per tissue were calculated as the average of the site

232   abundances over all experiments performed on that tissue. Due to inherent difficulties in

233   accurately reporting very low abundances with spectral counting [21], experiments in which no

234   peptides overlapping the site of modification were detected were not included in the average.

235   Thus, the sum of stoichiometries of all modifications at a particular site in a particular tissue may

236   not be normalized. Finally, the abundance, stoichiometry or normalized spectral count of a

237   modification site was set to zero for a particular tissue if the corresponding protein NSAF was

238   zero in that tissue.

239 **7. Gene ontology analysis**

8

240     Gene ontology analysis was conducted using the DAVID web portal [22]. For each post-

241     translational modification of interest, proteins bearing that modification were compiled and input

242     as a gene list. The background list used was the Uniprot human proteome. The resulting gene

243     ontologies were downloaded and a global FDR threshold (Benjamini-Hochberg) of 1% was used

244     as a threshold for determining significantly enriched ontologies.

245     We observed that many ontologies were broadly enriched across all post-translational

246     modifications and hypothesized that these were simply associated with highly abundant proteins

247     and did not reflect true post-translational modification properties. As a control, we applied the

248     above enrichment analysis to fifteen post-isolation modifications and observed many ontologies

249     that were significantly enriched for all post-isolation modifications considered (**Supplementary**

250     **Fig. 11**). These ontologies were excluded from the set of enriched ontologies in the post-

251     translational modification analysis (**Fig. 3b, Supplementary Table 6**).

252     **8. COSMIC dataset comparison**

253     A database of cancer mutations was downloaded via FTP from the COSMIC website[23]. The

254     mutation list was then filtered to keep only missense mutations. To guard against slight

255     variations in protein sequence between COSMIC and Uniprot, mutations for which the amino

256     acid residue at the denoted position in the Uniprot protein sequence did not match the non-

257     mutated amino acid identity in the corresponding COSMIC entry were discarded.

258     To guard against biases in background amino acid distributions, overlap statistics were only

259     calculated for proteins on which both cancer mutations and the PTM of interest was detected

260     and only against the background of peptides confidently identified by TagGraph in the human

261     proteome dataset. Using proline hydroxylation as an example, the number of prolines, number

262     of hydroxylation prolines, number of mutated prolines, and number of mutated and hydroxylated

263     prolines were counted only on peptides confidently identified by TagGraph and on proteins

264     containing both cancer mutations and proline hydroxylation. This overlap was then tested for

265     significance via Fisher's exact test. This analysis was carried out analogously for other types of

266     hydroxylations (lysine, asparagine, methionine, etc.).

267     **9. Protein-PTM correlation analysis**

268     Reasoning that many modifications' abundances and stoichiometries will depend on specific

269     protein-modifying enzymes, we sought to discover functional relationships between post-

270     translational modifications and proteins.  We identified highly correlated subsets of modifications

271     and proteins by comparing their abundances across the tissues examined here. Modification

272    site and protein lists were first filtered to include only those identified from at least three tissues.

273    For a particular post-translational modification of interest (e.g., Lysine hydroxylation), the

274    abundance of the modification was averaged across all measured sites from all proteins within

275    each tissue, forming a vector representing the abundance of the modification across all tissues.

276    Similarly, for all identified proteins, the calculated NSAF was used to form an abundance vector

277    of that protein's expression across all tissues.  We next determined the Pearson correlation

278    coefficient between all modification and all protein vectors computed and filtered as described

279    above. The proteins with the largest magnitude correlations (positive or negative) were then

280    considered as candidates having a functional relationship with a modification of interest.

281    Modification abundance vectors were calculated using both modification stoichiometries and

282    modification-normalized spectral counts. Both types of quantification were used in the

283    correlation analysis, often yielding different results (**Supplementary Fig. 13a**). However, in both

284    cases, our analysis revealed previously described associations between proteins and post-

285    translational modifications (e.g., arginine methylation and RNA splicing proteins), supporting the

286    validity of this analysis (**Supplementary Fig. 13c**).

287    **10. Code availability.**

288    The TagGraph algorithm is currently available via a web interface at

289    http://kronos.stanford.edu/TAG_GRAPH/.  The source code is available through github at
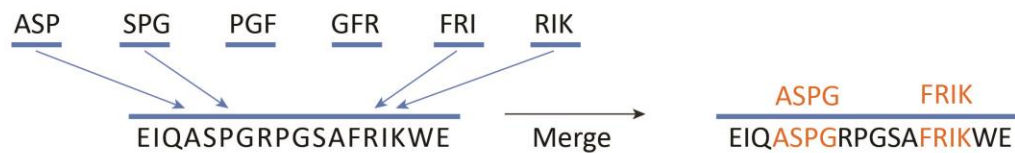
290    http://github.com/adevabhaktuni/XXXXXXXXX

291    **SUPPLEMENTARY FIGURES**
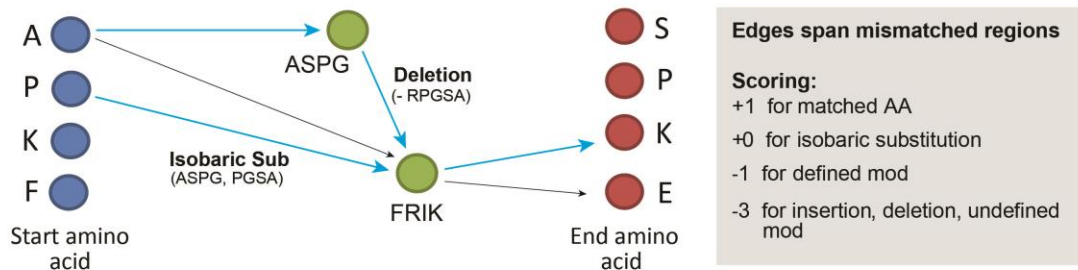
292    **Supplementary Fig. 1**

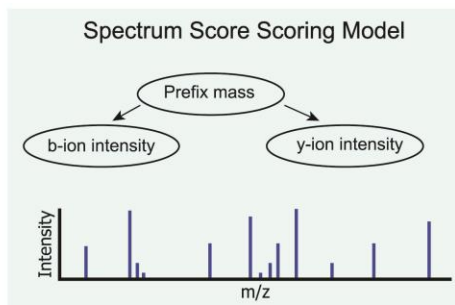## i. Identify candidates in database from *de novo* sequenced peptide



## ii. Tag Alignment



## iii. Find longest paths through graph



## iv. Score and rank matches



| | Candidates | Path Score | Spectrum Score | Total |
|---|---|---|---|---|
| 1 | R.PGSAFRIK.W | 4 | 5.6 | 9.6 |
| 2 | Q.ASPG(del RPGSA)FRIK.W | 5 | 4.2 | 9.2 |

## v. Secondary Result Population

**Mass Sum-based:**
Split mod into set of two mods with masses that sum to the original mod mass

pyrophospho                        phospho, phospho
K.ATGS(+159.93)SILYMK.E ⟶ K.ATGS(+79.97)S(+79.97)ILYMK.E ⟶ add to candidate list for spectrum

**Site based:**
Move mod site to more plausible amino acid

phospho H                          phospho S
K.H(+79.97)AVYILSPCCER.C ⟶ K.HAVYILS(+79.97)PCCER.E ⟶ add to candidate list for spectrum
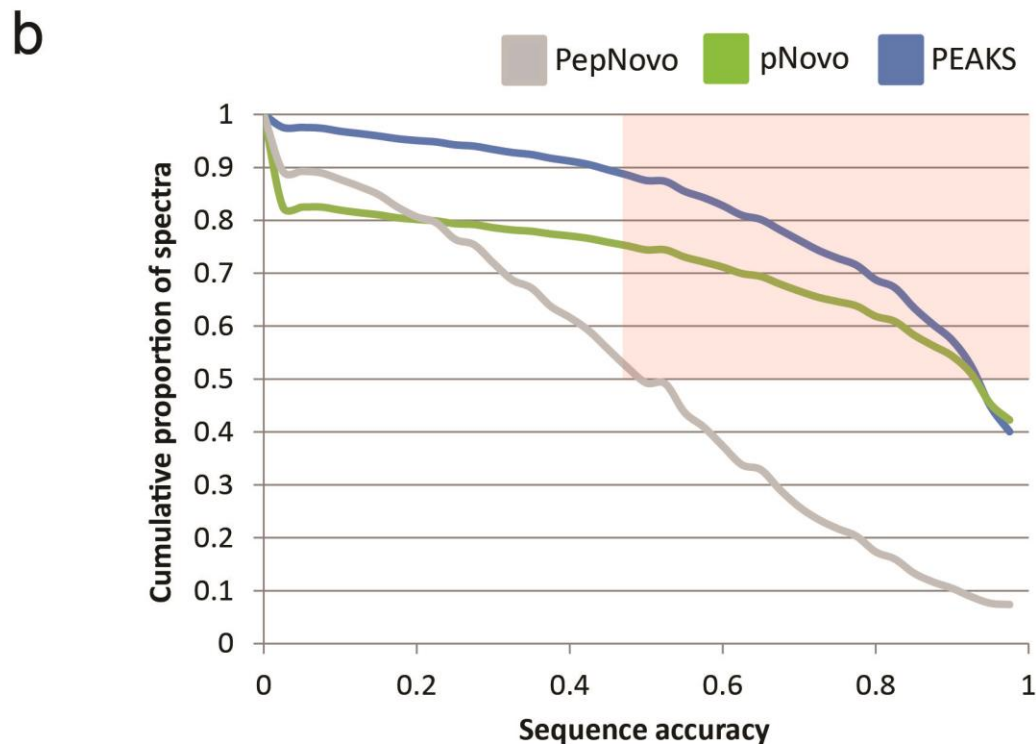
293

11

294     **Supplementary Fig. 1. TagGraph algorithm workflow overview.** TagGraph employs a five-

295     step procedure as depicted below, and detailed in **Supplementary Note 1**: **(i)** *De novo*

296     sequences are used to query an indexed sequence database. All candidate database entries

297     containing a maximum-length substring in common with the *de novo* sequence are retrieved. **(ii)**

298     The *de novo* sequence is compared against each database-derived candidate match.

299     Continuous amino acid substrings of length >2 that are identical between the query and

300     database candidate are identified as putative tags. **(iii)** Candidate matches (defined as a

301     peptide plus the set of its assigned modifications) are retrieved using a longest path algorithm

302     on a directed acyclic graph. Sequence tags defined in (ii) above are represented as nodes in the

303     graph and modifications as edges. Paths are drawn from start positions on the database peptide

304     to end positions through nodes and edges. **(iv)** Candidate matches over all database peptides

305     are collected and scored against the MS/MS spectrum using a probabilistic scoring model. **(v)**

306     After all *de novo* sequences are analyzed, additional candidate modification annotations are

307     created for select spectra if they are likely to be correct based on global dataset modification

308     abundances.

309     **Supplementary Fig. 2**



**a**     **Example sequence accuracy calculation**

Database peptide: A|S|F|R|G|K     } # Database peptide bonds: 5 (**Nd**)

} # Shared bonds: 3 (**Ns**)

*De novo* peptide: A|S|G|V|G|F|K

Sequence accuracy:     **Ns**/Nd     =     3/5     =     0.6
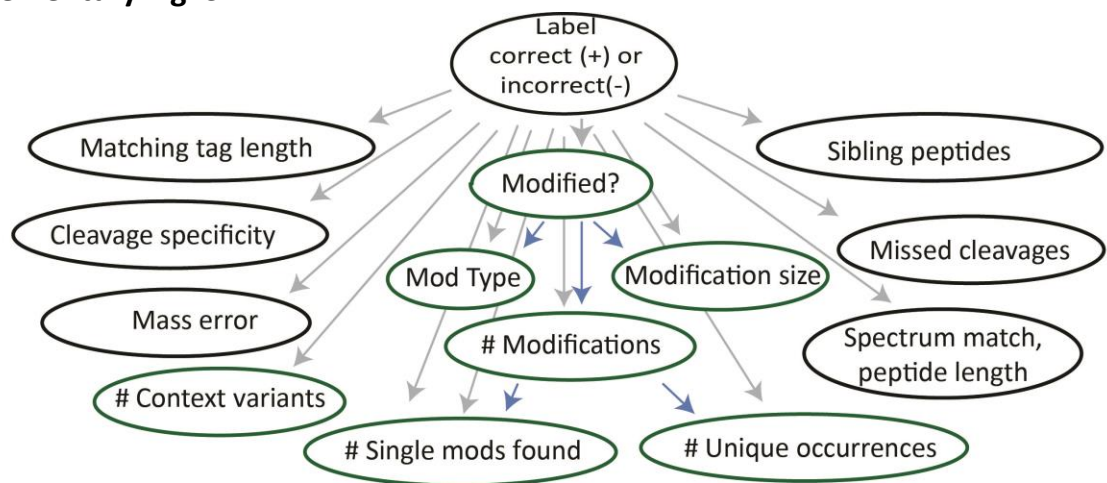
**b**

310

311     **Supplementary Fig. 2. *De novo* sequencing algorithms yield mostly correct**

312     **interpretations of most input spectra.  a)** Example calculation of sequence accuracy – the

313     proportion of peptide bonds shared between a high-confidence peptide identification and the

314     corresponding *de novo* peptide interpretation[9]. **b)** Cumulative proportion of spectra exceeding a

315     given sequencing accuracy threshold (x-axis) for three *de novo* sequencers, PEAKS, PepNovo,

316     and pNovo, as benchmarked on the A375 proteomic dataset (**Fig. 1**). PEAKS demonstrated the

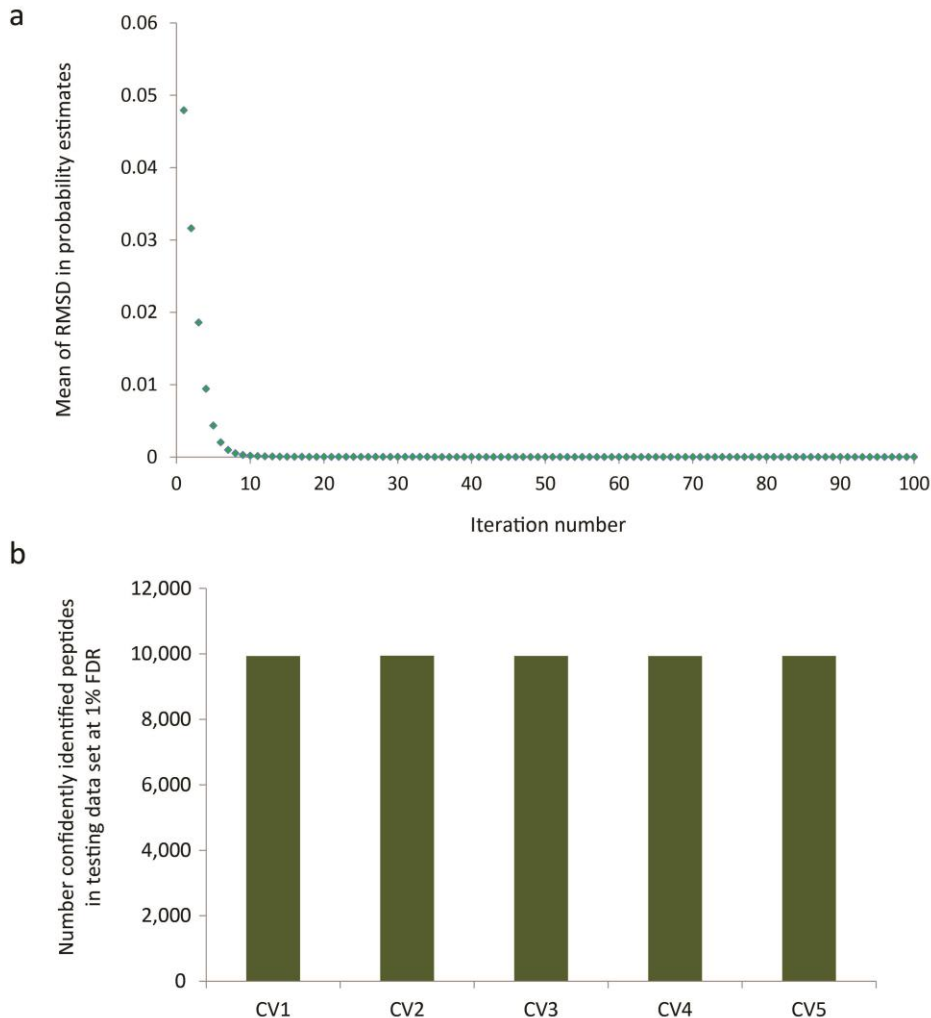317     best performance overall.

318 **Supplementary Fig. 3**



319

320 **Supplementary Fig. 3. Hierarchical Bayes model description. a)** Bayes model used for

321 fitting correct (+) and incorrect (-) peptide-spectrum match distributions. Grey arrows indicate

322 dependencies between model attributes and the distribution being trained. Blue arrows indicate

323 dependencies between model attributes. Attributes in green oval specifically pertain to

324 sequence modifications. Further details are provided in **Supplementary Note 3. b)** Example

325 distributions for several model attributes derived from the A375 dataset (**Fig. 1**). Likelihood
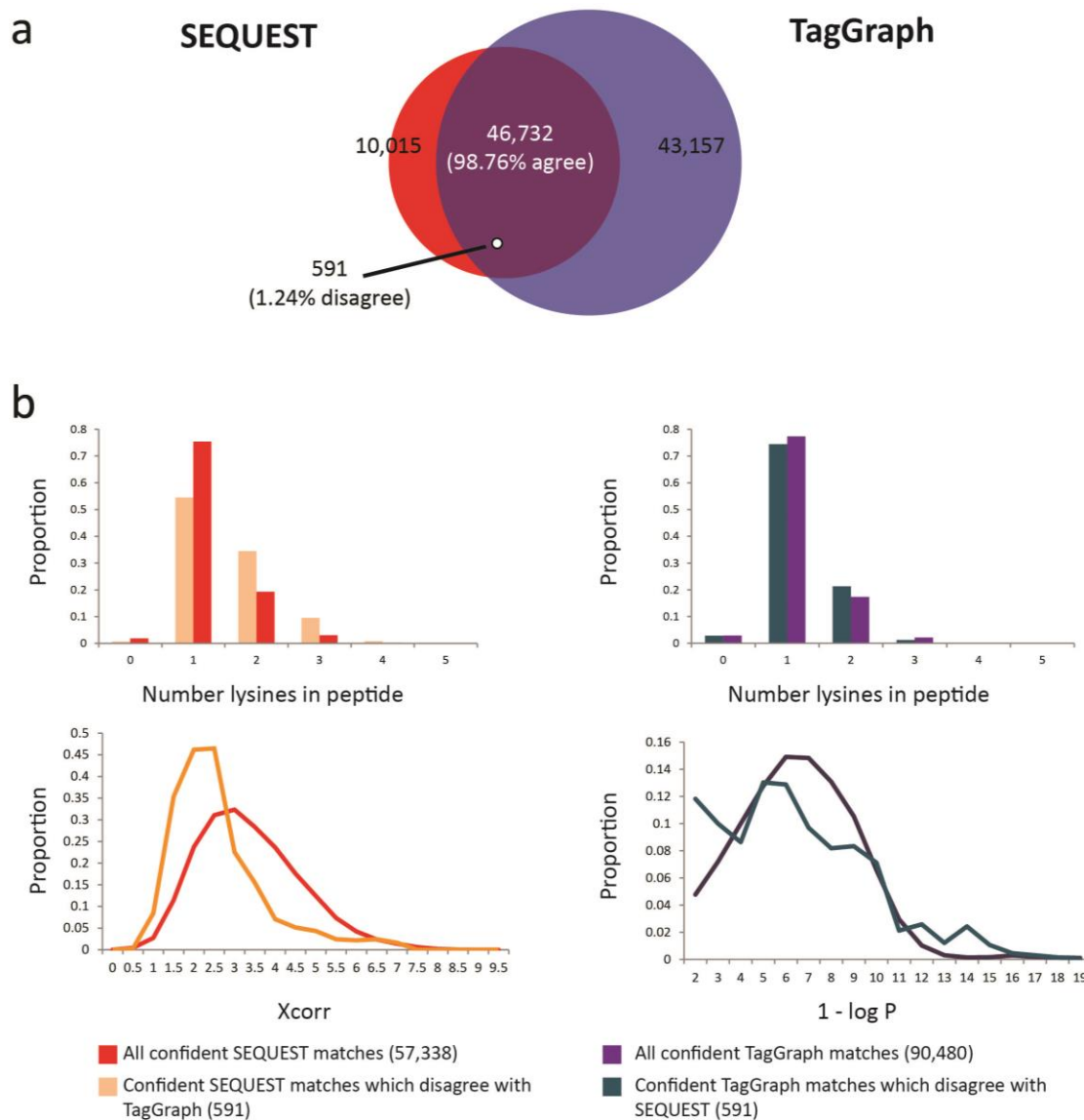
14

326    distributions were iteratively refined across multiple measurement dimensions using

327    expectation-maximization (EM).

328    **Supplementary Fig. 4**



329

330    **Supplementary Fig. 4.  Expectation Maximization-estimated false discovery rate**

331    **estimations are robust.  a)** Randomized starting model guesses for expectation-maximization-

332    based training of the hierarchical Bayes model rapidly converged, and yielded highly consistent

333    probability estimates. **b)** Five-fold cross-validation (CV) demonstrated that training the EM-

334    optimized hierarchical Bayes model did not substantially affect the returned set of confidently

335    identified spectra, when each model was tested on a dedicated test spectra set.  Further details

336    can be found in **Supplementary Note 3.**

**Supplementary Fig. 5**



338

**Supplementary Fig. 5. Conflicting high-confidence peptide-spectrum matches strongly favor TagGraph interpretations over SEQUEST. a)** 98.76% of 47,323 PSMs for which both TagGraph and SEQUEST return a high-confidence result (1% estimated FDR) agree, consistent with an estimated 1% FDR for both algorithms. PSMs wer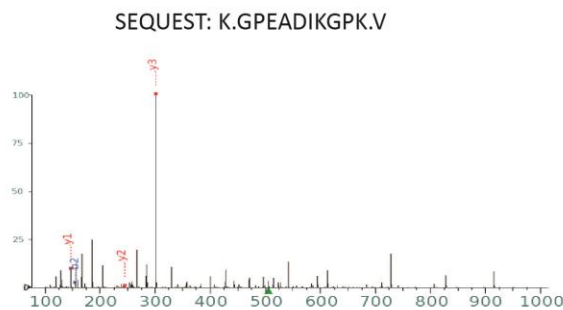e derived from the A375 dataset. **b)** Of the remaining 1.24% of PSMs for which SEQUEST and TagGraph disagree, TagGraph score and peptide missed cleavage distributions were more consistent with high-confidence identifications.
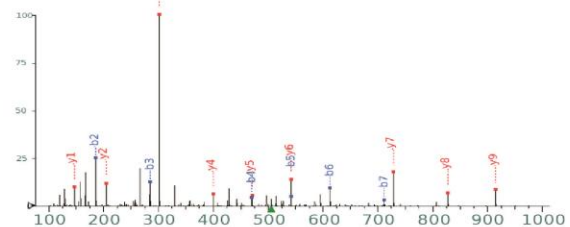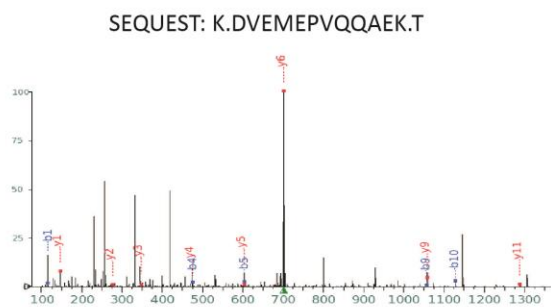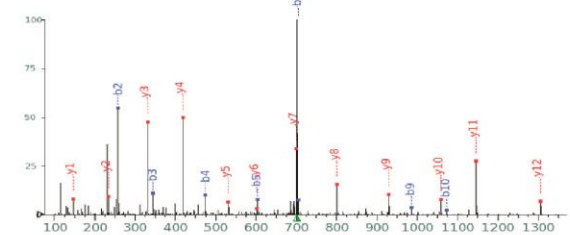
16

346 **Supplementary Fig. 6**

**Fraction 3 Scan 6679**

SEQUEST: K.GPEADIKGPK.V

TagGraph: S.PSVWAAVPGK.T

**Fraction 9 Scan 1760**

SEQUEST: K.DVEMEPVQQAEK.T

TagGraph: M.PCSEETPAISPSK.R

**Fraction 6 Scan 4745**

SEQUEST: K.NEADDYHLRNK.T

TagGraph: K.I(+43)MEIHGEGSSSGK.A
N-term carbamyl

**Fraction 8 Scan 3798**

SEQUEST: K.EDCRDEKVYSK.E

TagGraph: K.E(-18)FHINESGDPSSK.S
N-term pyroglutamination



347
348 **Supplementary Fig. 6. Examples of TagGraph-assigned peptide-spectrum matches that**

349 **conflict with high-confidence SEQUEST assignments.** Representative spectra

350 demonstrating superior fragment ion assignments made by TagGraph for peptides more

351 consistent with LysC digestion than the conflicting peptides SEQUEST assigned to the same

17

352    spectra.  Both results were assigned scores consistent with a 1% FDR on the A375 dataset with

353    respect to each set of search results.

354    **Supplementary Fig. 7**

355    *[See file TG_Figure_S07_Search_Algorithm_TypeI_II_Errors.pdf]*

356

357    **Supplementary Fig. 7. Examples spectra depicting "Case 1" (modification**

358    **mislocalization) and "Case 2" (incorrect peptide sequence) interpretation errors, as**

359    **defined Supplementary Note 2.**

**Supplementary Fig. 8**

a



b



Adult Testis
Q96NJ3 Zinc finger protein 285

Adult Kidney
Q9BYG8 Gasdermin-C

Adult Colon
Q99525 Histone H4-like protein type G

Adult Liver
Q15034 Probable E3 ubiquitin protein ligase

Adult Testis
Q96BD5 PHD finger protein 21A

Adult Frontal Cortex
P61962 DCAF7_HUMAN DDB1- and CUL4-associated factor 7

Adult Colon
Q15651 High mobility group nucle- osome-binding domain-containing protein 3

Adult Liver
Q8IUC6 TIR domain-containing adapter molecule 1

**Supplementary Fig. 8. Increased proteome coverage by TagGraph relative to Kim *et al*. a)**
The number of proteins identified by TagGraph and not Kim et. al. are shown for each tissue examined in this dataset. Identified proteins were assigned one of three categories: (i) proteins with any unmodified tryptic peptides mapped to them, (ii) proteins with unmodified non-tryptic peptides mapped to them and no unmodified tryptic peptides mapped, and (iii) proteins with only modified peptides mapped to them. Proteins were designated as identified in the Kim et. al., analysis if at least one peptide was mapped to them, and proteins were designated as present in the TagGraph analysis if their normalized spectral abundance factor (NSAF) was greater than
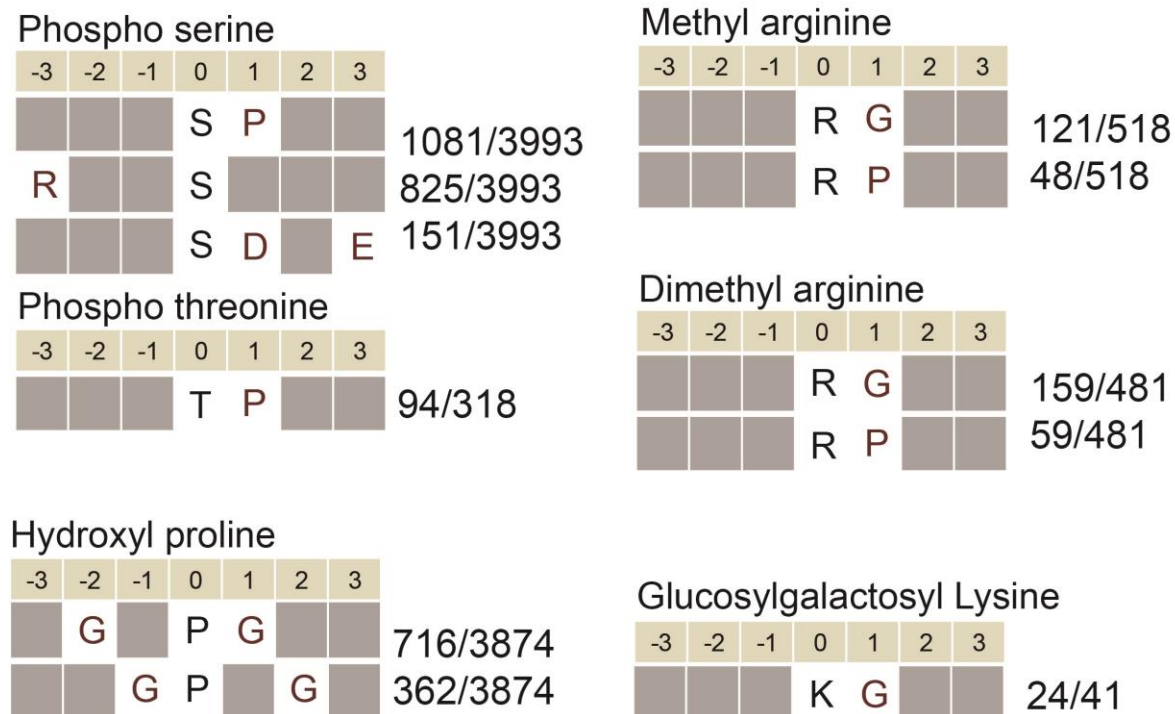
19

370     zero. We attribute the pronounced spike protein identifications from the Adult Monocytes tissue

371     to a procedural error made by the study's original authors:  We found that the pepXML-

372     formatted search result file corresponding with 'bRP_Elite' analysis, which we downloaded from

373     the PRIDE database (PXD000561) was identical to the 'bRP_Velos' pepXML file.  The raw data

374     files corresponding with these two conditions were clearly distinct, and were used as input to

375     TagGraph. This spike in identifications can only partially be attributed to TagGraph's enhanced

376     identification capabilities. **b)** Immunostaining images taken from ProteinAtlas[24] for select

377     proteins identified by TagGraph and not Kim et. al., validates TagGraph-specific protein

378     predictions.

379     **Supplementary Fig. 9**

380     *[See file TG_Figure_S09_SyntheticPeptides.pdf]*

381

382     **Supplementary Fig. 9. Synthetic peptides confirm novel post-translational modifications**

383     **and other unexpected sequence variants TagGraph measured from draft map of the**

384     **human proteome.**

385     **Supplementary Fig. 10**



386

387 **Supplementary Fig. 10. Motif-X analysis for known PTMs reveals known and novel**

388 **substrate motifs.** Motifs identified by the Motif-X algorithm[25] surrounding several abundant

389 PTMs. This analysis recovers known motifs for phosphoserine, phosphothreonine, dimethyl

390 arginine, and methyl arginine PTMs, and predicts new motifs for the less well-characterized

391 proline hydroxylation and lysine glucosylgalactosylation PTMs. Fraction indicates the number of

392 times the indicated motif was identified out of the total number of modification sites entered into

393 the Motif-X algorithm.

394    **Supplementary Fig. 11**



395

396    **Supplementary Fig. 11. Accounting for ontologies enriched among post-isolation**

397    **modifications.** Gene ontology enrichment analyses of PTM-bearing proteins may be biased by

398    mass spectrometers' tendency towards identifying modified peptides from highly abundant

399    proteins. Consequently, some ontologies could reach statistical significance based on protein

400    abundance alone, rather than PTM-specific biological phenomena. To account for this, we

22

401  identified significantly enriched ontologies (1% FDR, Benjamini-Hochberg corrected; yellow)

402  among proteins bearing any of 15 abundant post-isolation modifications. Because these

403  modifications should not have any inherent biological relevance, any ontology enriched among

404  these post-isolation modifications were deemed false (red brackets), and removed from the

405  analysis presented in **Fig. 3b**.

406  **Supplementary Fig. 12**



407

408  **Supplementary Fig. 12. Modification abundances and stoichiometries are not correlated**

409  **with protein abundances.** Scatterplots of protein normalized spectral abundances factor

410  (NSAF) with modification stoichiometry (left) or modification normalized spectral counts (NSC,

411  right). In both cases, modification abundance did not correlate with protein abundance.

412　　**Supplementary Fig. 13.**



413

414　**Supplementary Fig. 13. Proteins that correlate with PTM substrates share functional**

415　**properties. a)** Expression level (NSAF) profiles for 15,747 proteins spanning 30 tissues were

416　correlated with averaged PTM profiles across the same tissues, using either stoichiometry or

417　normalized spectral counts (NSC).  The representative scatter plot shown here for lysine

418　hydroxylation indicates the extent to which each protein's tissue profile (points) correlates with

419　lysine hydroxylation across the 30 tissues as measured by estimated stoichiometry (x-axis) or

420　total abundance (y-axis). These data show that the two PTM quantification methods are broadly

421　similar. However, protein correlation ranks may differ greatly between the two quantification

422　methods.  Thus, both can produce complementary but similar sets of highly correlated proteins.

423　**b)** Protein-PTM correlations generally did not indicate specific modified substrates.  A protein's

424　abundance could correlate with a particular PTM because it regulates or directly catalyzes the

425　PTM's formation on its substrate.  Alternatively, proteins could be correlated with a modification

426　because they are themselves heavily-modified substrates of the PTM.  Kinases, which both

427　catalyze phosphorylation events and are themselves highly phosphorylated, would be expected

428　to be examples of both conditions, for example.  By contrast, collagens would be examples of

429　the latter condition, as abundant proteins in certain tissues that carry a highly degree of

430　hydroxylated prolines. To evaluate these possibilities, we first identified the 20 proteins that

431　most highly correlated with each of the 28 PTMs shown here, as computed using either

24

432  modification NSC or stoichiometry.  Of these, we plotted the number of proteins that were also

433  modified by the indicated PTM.  For the most part, however, PTMs were not identified on the

434  same proteins to which they were most highly correlated, suggesting that they may be

435  candidate regulators of PTM transfer. **c)** Enriched gene ontologies for the top fifty most

436  correlated proteins for several PTMs suggests either enzymatic activity (i.e., oxidoreductase

437  activity is known to be required for lysine hydroxylation to occur) or common functional activity

438  (i.e., arginine dimethylation is known to be enriched in RNA splicing proteins, **Fig. 3b**). As

439  demonstrated in part b, these proteins are themselves not substrates of the PTM of interest.

440  Thus, these ontologies further suggest functional relationships between PTMs and proteins

441  which are highly correlated with them.

## SUPPLEMENTARY TABLES

443  **Supplementary Table 1**
444  **Supplementary Table 1. TagGraph, PEAKS PTM, Byonic, ModA, and Open search**
445  **results per mass spectrum from A375 cell line data set**

446  *[NB: this is a 92 Mb data file]*

447  a. Fraction number from high-pH reversed phase concatenated fractionation (see **Methods**)

448  b. Number of the MS2 fragmentation scan

449  c. Charge state of the precursor ion that gave rise to the indicated MS2 scan

450  d. Computed peptide's mass, based on its amino acid sequence and any additional
451  modifications

452  e. Inferred singly-charged ion mass, based on observed precursor ion's m/z ratio and
453  charge

454  f. Parts-per-million mass deviation between observed and theoretical peptide masses

455  g. Probability of indicated peptide identification being correct, as computed by search
456  algorithm.

457  h. Log-transformed, inverse probability of indicated peptide identification being correct, as
458  computed by the expectation maximization-optimized Bayesian network: $-\log_{10}(1-p)$

459  i. Peptide sequence from the input FASTA sequence database, noting the amino acids
460  flanking the peptide, appearing outside the periods.

461  j. TagGraph-resolved peptide sequencing, noting deviations from the database sequence
462  with a "-"

463  k. Modifications assigned to TagGraph-resolved peptide: Nested series are of the format:
464  (('*Mod1 name from Unimod if exists*', *Mod1 delta mass from Unimod if it exists*, *Mod1 delta*

465 *mass vs. Unimod if exists*), (*Mod1 target amino acid from Unimod if exists*, *Mod1 target amino*
466 *acid location on peptide from Unimod if exists*), *indexed location of Mod1 on peptide sequence*
467 *counting from zero*), ((*'Mod2 name from Unimod if exists'…, indexed location of Mod2 on*
468 *peptide sequence counting from zero*),(...)]

469     l. List of proteins from FASTA sequence database containing indicated peptide

470     m. Score assigned to identified peptide by search algorithm

471     n. Peptide sequence, indicating position of modification within parentheses.

472     o. Inferred name and specificity of modification indicated in (n)

473     p. Example protein containing indicated peptide sequence

474     q. modified amino acid and rounded mass of corresponding modification

475     r. Peptide sequence, indicating position of modification with numerical modification
476 representation immediately following modified residue

477     s. Deviation (Da) between observed and theoretical (unmodified) peptide masses

478

## Supplementary Table 2

479
480     **Supplementary Table 2. High-confidence TagGraph results per mass spectrum from**
481 **human proteome data set**

482 *[NB: this is a 10.2 Gb data file]*

483     a. Tissue from which mass spectrum was derived

484     b. Acquisition method for mass spectrum, using the format [separation method (SDS-PAGE
485 ("Gel") or high-pH reversed phase ("bRP"))]_[mass spectrometer (LTQ Orbitrap Velos ("Velos")
486 or Orbitrap Elite ("Elite")]

487     c. Fraction number from high-pH reversed phase concatenated fractionation (see **Methods**)

488     d. Number of the MS2 fragmentation scan

489     e Retention time (minutes) of the indicated MS2 scan

490     f. Charge state of the precursor ion that gave rise to the indicated MS2 scan

491     g. Inferred singly-charged ion mass, based on observed precursor ion's m/z ratio and
492 charge

493     h. Computed peptide's mass, based on its amino acid sequence and any additional
494 modifications

495     i. Parts-per-million mass deviation between observed and theoretical peptide masses

496     j. Probability of indicated peptide identification being correct, as computed by the
497 expectation maximization-optimized Bayesian network.

26

498      k. log-transformed, inverse probability of indicated peptide identification being correct, as
499 computed by the expectation maximization-optimized Bayesian network: -log10(1-p)

500      l. peptide sequence from the input FASTA sequence database, noting the amino acids
501 flanking the peptide, appearing outside the periods.

502      m. TagGraph-resolved peptide sequencing, noting deviations from the database sequence
503 with a "-"

504      n. Modifications assigned to TagGraph-resolved peptide: Nested series are of the format:
505 ((*'Mod1 name from Unimod if exists'*, *Mod1 delta mass from Unimod if it exists*, *Mod1 delta*
506 *mass vs. Unimod if exists*), (*Mod1 target amino acid from Unimod if exists*, *Mod1 target amino*
507 *acid location on peptide from Unimod if exists*), *indexed location of Mod1 on peptide sequence*
508 *counting from zero*),  ((*'Mod2 name from Unimod if exists'…*, *indexed location of Mod2 on*
509 *peptide sequence counting from zero*),(...)]

510      o. list of proteins from FASTA sequence database containing indicated peptide

511      p. peptide sequence derived from *de novo* sequencing

512      q. score assigned to *de novo* sequenced peptide by *de novo* algorithm

513

## Supplementary Table 3

515     Normalized spectral abundance factor (NSAF) for all proteins TagGraph measured from Kim
516 et al data set.

## Supplementary Table 4

518     All modifications found by TagGraph and their corresponding number of peptide-spectrum
519 matches and unique peptides.  Sites with at least 20 spectral counts were reported in
520 modification counts reported in text.

## Supplementary Table 5

522     Normalized spectral count (NSC) for all modified sites TagGraph measured from Kim et al
523 data set.

## Supplementary Table 6

525     Table S6: All enriched ontologies and corresponding p values for 22 noteworthy PTMs
526 (biological process and cellular compartment)

## Supplementary Table 7

528 All mono- and dimethylation sites with normalized spectral counts (NSC) and corresponding
529 protein abundances (NSAF).
530

## Supplementary Table 8

532     Estimated stoichiometry for all modified sites TagGraph measured from Kim et al data set.

533

534    **Supplementary Table 9**
535    PTMs TagGraph assigned to five representative histone isoforms, aligned with prior histone
536    PTM compilations.

537    **Supplementary Table 10**
538    Supplementary Table 10. Significant overlap between COSMIC cancer mutations and
539    TagGraph proline hydroxylations
540
541


542    **SUPPLEMENTARY NOTES**

543    **Supplementary Note 1: TagGraph**
544    **A. FM-index procedure.** The FM-Index[26] implementation used in TagGraph is a fork of an

545    existing open source implementation (https://github.com/mpetri/FM-Index). To create an FM-

546    index of the human proteome [12], we first concatenated all sequences in the FASTA-formatted

547    protein database into a flat sequence file. A separate database of protein start offsets is

548    maintained for retrieval of protein annotations. A Burrows-Wheeler Transform[27] was then

549    applied to the flat sequence file for fast substring search, the results of which were compressed

550    using an RRR Wavelet Tree [28] for efficient in-memory storage. The number of occurrences of a

551    candidate sequence pattern in an index can be computed in O(N) time, where N is the length of

552    the input pattern. The locations of the input pattern in an index can be retrieved in O(M) time,

553    where M is the number of occurrences of the input pattern in the index.


554    **B. TagGraph algorithm.** The TagGraph algorithm takes as input a set of high-resolution

555    MS/MS spectra, a corresponding set of *de novo* sequence interpretations, and an FM-Index

556    constructed from a protein FASTA file (**Supplementary Note 1a**). The algorithm then generates

557    and ranks a list of candidate peptide-spectrum matches for each input MS/MS spectrum with

558    respect to the indexed protein sequences.


559    TagGraph first computes the maximum matching substring between an input *de novo* sequence

560    and the FM-index, then retrieves all candidate protein sequences from the index which contain

561    this substring. For each candidate, TagGraph first computes all amino acid dimers which match

562    between the *de novo* sequence and protein sequence. Contiguous dimers are then merged into

563    longer sequence substring "tags," which are then input as nodes into a directed acyclic graph.

564    Edges are drawn between any two nodes to represent regions in which the *de novo* sequence

565    and matching candidate peptide sequence disagree. This disagreement could be due to a *de*

566    *novo* sequencing error or the presence of a sequence variant, post-isolation modification, post-

567 translational modification, or previously uncharacterized mass shift, in the peptide which gave

568 rise to the source spectrum. Each edge is annotated with its possible interpretations and

569 weighted based on a heuristic scoring scheme designed to weight more likely explanations

570 more highly (**Supplementary Fig. 1**). The top scoring peptide matches for each candidate are

571 retrieved from the graph as the top scoring paths between a set of start and end nodes [29].

572 These start and end nodes represent potential start and end sites for a peptide interpretation in

573 the candidate matching protein. The scored evidence supporting a peptide-spectrum

574 interpretation from its start to its end nodes is referred to as its **path score**.

575 Once identified, each candidate peptide match is scored against the observed MS/MS spectrum

576 using a probabilistic model to derive its **spectrum score**. This model scores all fragment ions

577 that support the peptide identification according to the relative likelihood of measuring the

578 observed fragment ion intensity versus random chance (Poisson) (**Supplementary Fig. 1**). The

579 probabilistic fragmentation model described above was trained on a library of 20,000 high

580 confidence in-house generated HCD MS/MS spectra, collected from peptides unrelated to the

581 current study.

582 To improve TagGraph's ability to discriminate true PTMs over confounding isobaric

583 interpretations (e.g., two phosphorylation events on neighboring serines versus the less likely

584 explanation of pyrophosphorylation on a single serine, **Supplementary Fig. 1**), the algorithm

585 creates a second set of candidate peptide-spectrum matches to explore alternate, isobaric

586 modifications learned from the input dataset: 1) All combinations of mass shift and

587 corresponding amino acid are tallied from the initial set of candidate peptide spectrum matches,

588 to create a mass shift – amino acid frequency matrix. 2) Each candidate peptide's mass shift is

589 evaluated with respect to the frequency matrix generated in (1). 3) If the peptide's mass shift is

590 equal to an existing modification corresponding with the same peptide assigned to a different

591 MS/MS spectrum, is more prevalent in the entire dataset than the existing modification, and has

592 a valid modifiable amino acid on the unmodified peptide sequence, then a peptide with this

593 alternate modification is added as a match candidate to the corresponding spectrum. 4) If the

594 mass of a candidate modification can be explained by the sum of two modifications represented

595 in the list learned in (1), an additional peptide carrying these two PTMs will be added as a match

596 candidate to the corresponding spectrum. This will occur, however, only if the expected number

597 of peptides with this combination of modifications in the dataset is greater than one, and if both

598 modifications have valid sites on the unmodified candidate peptide sequence. 5) All primary and

599 secondary match candidates are assigned a path score and spectrum score as described

29

600    above. 6) These are combined and ranked with the existing peptide-spectrum matches from the

601    first round of candidate generation (**Supplementary Fig. 1**).

602    ## Supplementary Note 2. Target-decoy error estimation is poorly suited to
603    **unrestricted search results**

604    The target-decoy validation methodology is readily applicable to search results generated by

605    conventional database search engines. Peptides identified by unrestricted search engines like

606    TagGraph pose several challenges that were not anticipated in our initial description and of this

607    error estimation tool, and which violate the major assumptions we proposed[30–32]. For target-

608    decoy to accurately estimate false discovery rates, the set of decoys must be chosen such that

609    incorrect identifications have an equal chance of matching either the target or decoy databases.

610    For conventional database search, choosing a decoy database composed of the reversed

611    counterparts of sequences in the target database largely satisfies this criterion. However, even

612    relatively simple searches permitting just one variable modification (e.g., phosphorylation),

613    secondary validation methods (e.g., Ascore [33] are needed to measure the modification's site

614    localization accuracy. This issue arises because the assumptions underpinning target-decoy are

615    violated: the likelihood of a correct (target) peptide bearing an incorrectly-localized modification

616    matching a given MS/MS spectrum is far greater than an incorrect, decoy peptide. As we will

617    demonstrate below, this problem is greatly compounded when considering hundreds of

618    modifications simultaneously (as Peaks PTM does), and becomes exponentially worse still

619    when allowing arbitrary mass modifications and no protease specificity. In all cases, errors

620    pertaining to the identity and localization of the annotated modifications and the location of the

621    peptide sequence in the proteome cannot be accurately estimated using target-decoy with

622    reversed sequence decoys.

623    **Case 1: Improper modification annotation and localization**

624    This case is the most common source of error in unrestricted database search, and is not

625    addressed by target-decoy based validation. When considering a peptide sequence with a

626    potential modification, an unrestricted search algorithm must score all possible localizations of

627    that modification on the peptide sequence.  The best, or highest scoring localization is often

628    returned, although other localizations could be considered as well. Every residue in the peptide

629    as well as the peptide's N- and C-termini serve as potential modification sites.

630    Algorithms which consider large numbers of known residue-specific modifications but not

631    undefined mass shifts, such as Peaks PTM and Byonic (not in wildcard mode) are still

632    confronted by extremely large numbers of possible localizations. For instance, according to

633    Unimod, methylation (+14.01565 Da) can be localized to the peptide N-terminus or C-terminus,

634    as well as the amino acids C,H,K,N,Q,R,D,E,S, and T. Furthermore, many modifications in

635    Unimod have similar, near-isobaric masses, increasing the set of potential localizations beyond

636    those suggested by the modification itself.

637    Another potential source of FDR estimation error lies in determining the number and types of

638    modifications assigned to a peptide. When considering a peptide with a known mass shift

639    corresponding to a limited set of potential modifications, an algorithm must decide whether the

640    mass shift corresponds to a single modification, two modifications in combination, or three or

641    more. Even when considering only well-defined modifications, the combined mass of two

642    modifications often equals the mass of third modification – two phosphorylations equal the mass

643    of one pyrophosphorylation; two methylations equal the mass of one dimethylation, etc. There

644    are also a considerable number of cases in which a single parent mass could correspond to

645    various combinations of different modifications. For instance, the mass 86.00 Daltons could

646    correspond to two carbamylations or one acetylation and one carboxylation. In the case of

647    arbitrary mass modifications, a given mass shift could correspond with an effectively infinite

648    range of modifications and modification combinations. The enormity of the set of possible

649    modifications also limits the effectiveness of the open search method [16], necessitating user-

650    supervised follow-up analysis to identity the modification corresponding with a measured mass

651    shift, and to localize it on the returned peptide sequence.

652    Due to the abundance of candidates considered by the search engine, errors in modification

653    annotation are common in all unrestricted search methods. These errors can arise from low-

654    quality mass spectra, resulting from noise, incomplete fragmentation, co-isolated precursor ions,

655    or other confounding features.  They can also arise from the search engine's configuration, such

656    as internal scoring functions that inappropriately give extra weight modifications believed to be

657    more likely *a priori,* or because the algorithm did not consider the true modification annotation

658    as a candidate. In general, peptide-spectrum matches with incorrect modification annotations

659    share many of the same b- and y-ions as the correct modification annotations. Thus, it is far

660    more likely for a peptide to match an incorrect modification annotation than a reversed decoy

661    sequence. A few examples of such errors produced by previously published unrestricted search

662    algorithms are provided (**Supplementary Figure 7**).

663    Rare but biologically important PTMs can often occur at rates orders of magnitude lower than

664    the most common post-isolation modification in a dataset. Without an accurate method to

665    estimate FDRs of modification annotations, systematic algorithm-dependent errors are subject

666    to being mistaken for these kinds of rare peptide-spectrum matches, limiting the utility of

667    unrestricted modification searches for finding them. TagGraph's hierarchical Bayesian model

668    addresses this issue by making no assumptions about types of modifications present in the

669    dataset *a priori* (or about any other attribute of the dataset). Thus, the validation model learns

670    the distribution of modifications present, and weights its confidence in a particular modification

671    annotation against other attributes of the peptide-spectrum match, such as the evidence

672    provided in the spectrum, localization of the peptide sequence, and other attributes.

673    **Case 2: Incorrect base peptide sequence**

674    The large set of possible modifications an unrestricted search algorithm must consider also

675    makes it possible for such algorithms to choose an incorrect base peptide sequence. This

676    incorrect peptide sequence could originate from a protein other than the correct sequence (i.e.,

677    a sequence from a homologous protein containing a single amino acid polymorphism), or could

678    be a slight deviation from the correct peptide sequence (i.e., with several amino acids added or

679    truncated from the N- or C- terminus). In all cases, these incorrect base peptide sequence

680    predictions can be reconciled with the MS/MS spectrum's precursor mass through erroneous

681    modification annotations. As above, these incorrect base peptide matches are more likely than

682    decoy sequence matches to a spectrum since they often share many b- and y-ions with the

683    correct peptide-spectrum match. A few examples of such errors are provided (**Supplementary**

684    **Figure 7**).

685    The rate at which these types of errors occur is determined by the degree of sequence self-

686    similarity present in the sequence database and the number of modifications considered. As the

687    number of modifications considered increases, more base peptide sequences can be

688    considered as candidate matches for a peptide spectrum match. Similarly, as sequence self-

689    similarity increases (i.e., due to homology), more sequences similar to the true peptide

690    sequence can be considered as candidates: mass shifts assigned to one or more modifications

691    can make them consistent with the spectrum's precursor mass, even with very high degrees of

692    measurement mass accuracy.

693    **Exacerbation of Case I and Case II errors when allowing arbitrary protease specificity**

694    We have shown with specific examples and with global analysis that previously described

695    unrestricted search algorithms are prone to the above errors cases.  Furthermore, the rates of

696    these errors greatly exceed what would be expected using target-decoy FDR estimation. This

697    error underestimation problem could be worse for TagGraph, since it is the first to perform truly

698    unconstrained search, with respect to protease specificity, and arbitrary types and numbers of

699    modifications, and sequence variants. This larger search space raises the likelihood of matching

700    an incorrect peptide sequence to an input spectrum by chance.  We explore several of these

701    possibilities below.

702    Peptides that are inconsistent with protease specificity can result from either in-source

703    fragmentation or endogenous protease activity.  Though often excluded from reported datasets,

704    TagGraph was designed to readily detect them in either modified and unmodified forms

705    (**Supplementary Note 1**). Spectra corresponding to such peptides can return high-scoring

706    incorrect identifications with a protease-constrained unrestricted search (i.e., returning the

707    nearest protease-specific peptide with a modification annotation -- Open search is prone to this

708    error). As such results must then be manually separated from true, modification-bearing

709    peptides, an algorithm such as TagGraph which can correctly annotate these peptide-spectrum

710    matches from the outset is highly desirable.

711    However, opening up the search space in this way poses substantial challenges in validation in

712    addition to those presented above. Namely, the set of candidates which must be considered by

713    the algorithm increases exponentially. For instance, consider an MS/MS spectrum for which the

714    correct identification is the following peptide-spectrum match:

715    …KAAREE(+43)TDFCEEDSIEKSS…

716    Where the underlined sequence indicates the returned peptide-spectrum match, placed in the

717    context of its surrounding protein sequence in the database. The +43 indicates an N-terminal

718    carbamylation. An algorithm which allows arbitrary mass modifications and no protease

719    specificity must consider an effectively infinite number of candidates. However, unlike the case

720    in which protease specificity is constrained *a priori*, the algorithm must also consider multiple

721    base peptide sequences. For example, here are several incorrect annotations which could be

722    candidate matches for the spectrum yielding the true match above:

723    …KAAREE(-86)TDFCEEDSIEKSS…

724    …KAARE(-86)ETDFCEEDSIEKSS…

725    …KAAREET(+172)DFCEEDSIEKSS…

726    …KAAREET(+43)DFCE(+129)EDSIEKSS…

727    …KAAREE(+43)TDFCEEDSIE(+128)KSS…

728    …KAARE(-86)ETDFCEEDSIE(+128)KSS…

729

730    These annotations have similar predicted fragmentations to the correct peptide-spectrum match.

731    As described in Case I and II above, this similarity makes them far more likely than a match to a

732    decoy sequence in the case of a false match, rendering target-decoy unsuitable. Algorithms

733    which constrain the set of modifications and protease specificity face the possibility of Case I

734    and II errors on a subset of spectra, and commit them involuntarily by excluding the correct

735    peptide-spectrum match from the set of candidates on another subset of spectra. By allowing

736    arbitrary mass modifications and no protease specificity, TagGraph must be able to discriminate

737    correct peptide spectrum matches from Case I and Case II errors for *every* spectrum.

738

739    We address this problem by using an expert-designed hierarchical Bayesian model

740    (**Supplementary Note 3**) which can learn the specific attributes of correct and incorrect peptide-

741    spectrum matches from each dataset individually, including the attributes of correct and

742    incorrect modification annotations. This avoids problems present in the scoring functions of

743    some previously published unrestricted algorithms, which make *a priori* assumptions on the

744    likelihood of observing certain modifications in the dataset that are often untrue. Unlike target-

745    decoy, this model is able to assign higher or lower confidence to a peptide-spectrum match

746    based on its modification annotation, and consider this attribute in the context of the peptide

747    localization, evidence in the spectrum, etc. The model is both flexible and extensible, enabling

748    further refinement based on discriminatory criteria discovered to be useful in the future.

749    **Supplementary Note 3. Estimating peptide identification error through a**
750    **hierarchical Bayes probabilistic model, optimized by Expectation Maximization**
751    **A. Overview.** We developed a peptide identification error model that overcomes the limitations

752    conventional target-decoy searching has for assessing modification-bearing peptide

753    identifications. Our hierarchical Bayes model[7] creates hypothetical structures for the probability

754    distributions corresponding with observing a set of data features D given that the peptide

755    interpretation is correct P(D|+) or incorrect P(D|-). D is defined as the set of peptide and

756    fragmentation spectrum attributes that are represented by the hierarchical Bayes model

757    (**Supplementary Fig. 3a**). These attributes were empirically chosen based on their utility in

758    discriminating correct peptide-spectrum matches from incorrect ones. Considering both of the

34

759 above probability distributions, we calculate P(+|D), the probability that any given peptide-

760 spectrum interpretation is correct, using Bayes Theorem.

761 **B. Model attributes.** The following describe the general and modification-specific attributes

762 used in our hierarchical Bayes model, as represented in **Supplementary Fig. 3a**.

763 **General peptide attributes**

764 *Spectrum Score*: The spectrum score of a candidate peptide-spectrum match, as defined in (*iii*),

765 above.

766 *Peptide Length:* The length of an entire candidate peptide-spectrum match. This attribute is

767 modeled jointly with the spectrum score due to the observation that the spectrum score is

768 negatively correlated with peptide length (**Supplementary Fig. 3b**).

769 *Cleavage Specificity:* If a protease was used to digest the sample, this attribute encodes

770 whether a peptide-spectrum match has full protease specificity, a nonspecific n-terminus, a

771 nonspecific c-terminus, or both nonspecific termini. This attribute is ignored if no protease was

772 used.

773 *Missed Cleavages:* If a protease was used to digest the sample, this attribute records how many

774 missed protease cleavage sites are present in a peptide-spectrum match.

775 *Matching Tag Length:* The length of the maximal matching substring between the *de novo*

776 peptide interpretation and the candidate peptide-spectrum match.

777 *Sibling Peptides*: Number of other unique peptides found from the same protein that produced

778 the current peptide-spectrum match.

779 *Mass Error*: The difference in Daltons between the mass of a peptide-spectrum match and the

780 predicted mass based on the peptide sequence and the set of its annotated modifications.

781 **Modification-derived Attributes**

782 *Modified?:* Categorical attribute denoting whether or not the current peptide-spectrum match

783 has modifications with respect to the underlying database sequence.

784 *Maximum Modification Mass:* The mass (Da) of the largest modification assigned to a candidate

785 peptide-spectrum match.

786 *Number Unique Occurrences:* The number of unique peptides observed in the dataset with the

787 same set of modifications as the candidate peptide-spectrum match.

35

788    *Modification Type:* Categorical labeling of annotated modification as either an amino acid

789    substitution, defined modification (i.e., present in Unimod but not an amino acid substitution),

790    insertion/deletion, or undefined mass shift.

791    *Number of Modifications:* The number of modifications annotated on the peptide-spectrum

792    match.

793    *Number of Context Variants:* The number of unique peptides present in the dataset with the

794    same base peptide sequence as the candidate peptide-spectrum match. A unique peptide is

795    defined by its primary amino acid sequence, all modifications it may contain, and the

796    corresponding modification positions along the sequence.

797    *Number of Single Modifications Found On Same Context:* If a given peptide-spectrum match

798    contains multiple modifications, this attribute records the number of distinct singly-modified

799    forms of the same peptide found in the entire dataset.

800    **C. Learning the distributions $P(D|+)$ and $P(D|-)$.** As the set of correct and incorrect peptide-

801    spectrum matches is not known *a priori*, we use an expectation maximization algorithm to learn

802    the distributions $P(D|+)$ and $P(D|-)$ from observed data features D. Learning these distributions

803    is accomplished through two steps: a re-ranking step and a convergence step.

804    Several model attributes in the hierarchical Bayes model rely on quantities calculated from the

805    set of all peptide-spectrum matches in a dataset. These attributes can significantly affect the

806    probabilities $P(+|D)$ of candidate peptide-spectrum matches for each spectrum. Thus, the

807    estimate of $P(+|D)$ within a given iteration of the expectation-maximization process can change

808    both the attributes D for all peptide-spectrum matches and the optimal match for a particular

809    spectrum relative to the previous iteration. We expect the rankings produced by the probabilities

810    $P(+|D)$ to be superior to the initial rankings TagGraph produces using the sum of the spectrum

811    and path scores alone. Given that the estimate of $P(+|D)$ changes with each iteration, we

812    further expect more sensitive discrimination between correct and incorrect identification by

813    allowing the rankings of candidates for each spectrum to shift as the estimate of $P(+|D)$

814    increases in accuracy. This intuition informs the basis of the re-ranking step of model training. In

815    this step, the estimates of $P(D|+)$ and $P(D|-)$ are calculated using the top scoring candidates for

816    each spectrum per expectation-maximization iteration. The candidates for each spectrum are

817    then re-ranked according to their corresponding probabilities $P(+|D)$, and the estimates of

818    $P(D|+)$ and $P(D|-)$ in the next expectation-maximization iteration are calculated based on this

819    new set of top-ranked candidates. After 20 rounds of re-ranking, the top ranked candidates are

820    fixed in position and the algorithm is allowed to converge on stable estimates of P(D|+) and

821    P(D|-) during the convergence step. Model variance is calculated as the Euclidean distance

822    between the vectors P(+|D) for all spectra between the current and previous iteration.

823    Convergence of the model to near zero variance typically occurs within 100 iterations.

824    Within each iteration, the distributions P(D|-) and P(D|+) are learned as follows: the spectrum

825    score and peptide length attributes are fitted as a multivariate Gaussian using a maximum

826    likelihood estimator, weighted by the probability estimates derived from the previous expectation

827    maximization iteration. The remaining model attributes are discretized into bins: the probability

828    of observing each bin B for a given attribute A for the (+) distribution in the current expectation-

829    maximization iteration is calculated using the estimates of P(+|D) from the previous iteration

830    according to the familiar formula[34]:

831
$$P(B|+) = \frac{\sum_{\{i|A_i \in B\}} P(+|D)_i}{\sum P(+|D)_i} \quad \text{(Equation 1)}$$

832    The formula for calculating attribute probabilities for the (-) distribution can be analogously

833    generated using estimates of P(-|D). Before the first iteration, the EM algorithm is supplied with

834    an initial guess for the parameters of the multivariate Gaussian describing the peptide length

835    and spectrum score and for the distribution of the matching tag length attribute (**Supplementary**

836    **Fig. 3**). Initial guesses for the parameters of both the correct (+) and incorrect (-) distributions

837    are supplied. These guesses are used to populate initial estimates for the probabilities P(+|D)

838    for each peptide-spectrum match in the dataset. These probabilities are then iteratively

839    improved using the expectation maximization algorithm as described above. The spectrum-level

840    probabilities P(-|D) can be readily converted to a global false discovery rate for a given set of

841    spectra S using the following formula:

842
$$FDR_S = \frac{\sum_{\{j|j \in S\}} P(-|D)_j}{N_S} \times 100 \quad \text{(Equation 2)}$$

843    Where $N_S$ is the number of spectra in S. The large number of free parameters used to generate

844    this model could be susceptible to overtraining with datasets of small size. However, we found

845    that the algorithm converges onto accurate probability estimates for datasets of the size typically

846    produced in modern proteomics experiments (**Supplementary Fig. 4, Supplementary Note**

847    **3D**).

848 **D. Evaluating EM model stability.** We probed the robustness of the expectation-maximization

849 based learning approach in two different ways, both using the cell line dataset described in

850 **Figure 1**. First, the initial guesses used to seed the model training were randomly varied. The

851 EM algorithm was run as described in section (iv), and the estimates of P(+|D) of all peptide-

852 spectrum matches in the dataset were recorded for each iteration following the random initial

853 guess. The root mean square deviation (RMSD) of each probability estimate was computed and

854 averaged over the entire probability vector to derive the Mean RMSD over all initial guesses.

855 This deviation asymptotically approaches zero with increased iterations, demonstrating that the

856 final probability estimates are independent of the initial guess used (**Supplementary Fig. 4a**).

857 Second, to assess whether or not the EM model was susceptible to over-fitting, we employed

858 five-fold cross validation. First, the dataset was randomly split into testing and training subsets,

859 consisting of 10% and 90% of the peptide-spectrum matches, respectively. Slices consisting of

860 80% of the training dataset were randomly chosen and used to train the parameters of the

861 P(D|+) and P(D|-) distributions. Estimates of the probabilities P(+|D) for the peptide-spectrum

862 matches in the testing dataset were computed at each iteration using the models trained from

863 each slice. The final probability estimates derived from the testing set were found to be

864 independent of the slice used for training, demonstrating that the model learned general

865 features of the data and did not over-fit to specific attributes of the randomly chosen subsets

866 (**Supplementary Fig. 4b**).

867 **REFERENCES**

868 1. Fok, J. Y., Ekmekcioglu, S. & Mehta, K. Implications of tissue transglutaminase

869 expression in malignant melanoma. *Mol. Cancer Ther.* **5,** 1493–503 (2006).

870 2. Yang, F., Shen, Y., Camp, D. G. & Smith, R. D. High-pH reversed-phase chromatography

871 with fraction concatenation for 2D proteomic analysis. *Expert Rev. Proteomics* **9,** 129–34

872 (2012).

873 3. Vizcaíno, J. A. *et al.* A guide to the Proteomics Identifications Database proteomics data

874 repository. *Proteomics* **9,** 4276–83 (2009).

875 4. Kim, M.-S. *et al.* A draft map of the human proteome. *Nature* **509,** 575–81 (2014).

876 5. Vizcaíno, J. A. *et al.* A guide to the Proteomics Identifications Database proteomics data

877 repository. *Proteomics* **9,** 4276–83 (2009).

878  6.  Ma, B. *et al.* PEAKS: powerful software for peptide de novo sequencing by tandem mass
879      spectrometry. *Rapid Commun. Mass Spectrom.* **17,** 2337–42 (2003).

880  7.  Frank, A. & Pevzner, P. PepNovo: de novo peptide sequencing via probabilistic network
881      modeling. *Anal. Chem.* **77,** 964–73 (2005).

882  8.  Chi, H. *et al.* pNovo: de novo peptide sequencing and identification using HCD spectra. *J.*
883      *Proteome Res.* **9,** 2713–24 (2010).

884  9.  Devabhaktuni, A. & Elias, J. E. Application of de novo sequencing to large-scale complex
885      proteomics datasets. *Journal of Proteome Research* (2016). Available at:
886      http://pubs.acs.org/doi/pdfplus/10.1021/acs.jproteome.5b00861. (Accessed: 19th January
887      2016)

888  10. Dancík, V., Addona, T. A., Clauser, K. R., Vath, J. E. & Pevzner, P. A. De novo peptide
889      sequencing via tandem mass spectrometry. *J. Comput. Biol.* **6,** 327–42

890  11. Eng, J. K., McCormack, A. L. & Yates, J. R. An approach to correlate tandem mass
891      spectral data of peptides with amino acid sequences in a protein database. *J. Am. Soc.*
892      *Mass Spectrom.* **5,** 976–989 (1994).

893  12. O'Donovan, C. *et al.* High-quality protein knowledge resource: SWISS-PROT and
894      TrEMBL. *Brief. Bioinform.* **3,** 275–84 (2002).

895  13. Han, X., He, L., Xin, L., Shan, B. & Ma, B. PeaksPTM: Mass spectrometry-based
896      identification of peptides with unspecified modifications. *J. Proteome Res.* **10,** 2930–2936
897      (2011).

898  14. Bern, M., Kil, Y. J. & Becker, C. Byonic: advanced peptide and protein identification
899      software. *Curr. Protoc. Bioinformatics* **Chapter 13,** Unit13.20 (2012).

900  15. Na, S., Bandeira, N. & Paek, E. Fast multi-blind modification search through tandem
901      mass spectrometry. *Mol. Cell. Proteomics* **11,** M111.010199 (2012).

902  16. Chick, J. M. *et al.* A mass-tolerant database search identifies a large proportion of
903      unassigned spectra in shotgun proteomics as modified peptides. *Nat. Biotechnol.* **33,**
904      743–9 (2015).

905  17. Deutsch, E. mzML: a single, unifying data format for mass spectrometer output.
906      *Proteomics* **8,** 2776–7 (2008).

907   18.   Nesvizhskii, A. I. A survey of computational methods and error rate estimation
908         procedures for peptide and protein identification in shotgun proteomics. *J. Proteomics* **73,**
909         2092–123 (2010).

910   19.   Huttlin, E. L. *et al.* A tissue-specific atlas of mouse protein phosphorylation and
911         expression. *Cell* **143,** 1174–1189 (2010).

912   20.   Zhang, Y., Wen, Z., Washburn, M. P. & Florens, L. Refinements to label free proteome
913         quantitation: how to deal with peptides shared by multiple proteins. *Anal. Chem.* **82,**
914         2272–81 (2010).

915   21.   Choi, H., Fermin, D. & Nesvizhskii, A. I. Significance analysis of spectral count data in
916         label-free shotgun proteomics. *Mol. Cell. Proteomics* **7,** 2373–85 (2008).

917   22.   Huang, D. W., Sherman, B. T. & Lempicki, R. A. Systematic and integrative analysis of
918         large gene lists using DAVID bioinformatics resources. *Nat. Protoc.* **4,** 44–57 (2009).

919   23.   Forbes, S. A. *et al.* COSMIC: Exploring the world's knowledge of somatic mutations in
920         human cancer. *Nucleic Acids Res.* **43,** D805–D811 (2015).

921   24.   Uhlén, M. *et al.* A human protein atlas for normal and cancer tissues based on antibody
922         proteomics. *Mol. Cell. Proteomics* **4,** 1920–32 (2005).

923   25.   Schwartz, D. & Gygi, S. P. An iterative statistical approach to the identification of protein
924         phosphorylation motifs from large-scale data sets. *Nat. Biotechnol.* **23,** 1391–8 (2005).

925   26.   Ferragina, P. & Manzini, G. Opportunistic data structures with applications. in
926         *Proceedings 41st Annual Symposium on Foundations of Computer Science* 390–398
927         (IEEE Comput. Soc, 2000). doi:10.1109/SFCS.2000.892127

928   27.   Li, H. & Durbin, R. Fast and accurate long-read alignment with Burrows-Wheeler
929         transform. *Bioinformatics* **26,** 589–95 (2010).

930   28.   *String Processing and Information Retrieval*. **5280,** (Springer Berlin Heidelberg, 2009).

931   29.   Lu, B. & Chen, T. A suboptimal algorithm for de novo peptide sequencing via tandem
932         mass spectrometry. *J. Comput. Biol.* **10,** 1–12 (2003).

933   30.   Peng, J., Elias, J. E., Thoreen, C. C., Licklider, L. J. & Gygi, S. P. Evaluation of
934         Multidimensional Chromatography Coupled with Tandem Mass Spectrometry
935         (LC/LC−MS/MS) for Large-Scale Protein Analysis: The Yeast Proteome. *J. Proteome*

936        *Res.* **2,** 43–50 (2003).

937    31.    Moore, R. E., Young, M. K. & Lee, T. D. Qscore: an algorithm for evaluating SEQUEST
938           database search results. *J. Am. Soc. Mass Spectrom.* **13,** 378–86 (2002).

939    32.    Elias, J. E. & Gygi, S. P. Target-decoy search strategy for increased confidence in large-
940           scale protein identifications by mass spectrometry. *Nat. Methods* **4,** 207–14 (2007).

941    33.    Beausoleil, S. A., Villén, J., Gerber, S. A., Rush, J. & Gygi, S. P. A probability-based
942           approach for high-throughput protein phosphorylation analysis and site localization. *Nat.*
943           *Biotechnol.* **24,** 1285–92 (2006).

944    34.    Nesvizhskii, A. I., Keller, A., Kolker, E. & Aebersold, R. A Statistical Model for Identifying
945           Proteins by Tandem Mass Spectrometry. *Anal. Chem.* **75,** 4646–4658 (2003).

946