



## King's Research Portal

DOI:

[10.1371/journal.pone.0235263](https://doi.org/10.1371/journal.pone.0235263)

*Document Version*

Peer reviewed version

[Link to publication record in King's Research Portal](#)

*Citation for published version (APA):*

Preston, G., Yang, L., Phillips, D., & Maier, C. (2020). Visualisation tools for dependent peptide searches to support the exploration of in vitro protein modifications. *PloS one*, 15(7), e0235263. [e0235263]. <https://doi.org/10.1371/journal.pone.0235263>

### **Citing this paper**

Please note that where the full-text provided on King's Research Portal is the Author Accepted Manuscript or Post-Print version this may differ from the final Published version. If citing, it is advised that you check and use the publisher's definitive version for pagination, volume/issue, and date of publication details. And where the final published version is provided on the Research Portal, if citing you are again advised to check the publisher's website for any subsequent corrections.

### **General rights**

Copyright and moral rights for the publications made accessible in the Research Portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognize and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the Research Portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the Research Portal

### **Take down policy**

If you believe that this document breaches copyright please contact [librarypure@kcl.ac.uk](mailto:librarypure@kcl.ac.uk) providing details, and we will remove access to the work immediately and investigate your claim.

1

2

3

4

Visualisation tools for dependent peptide searches to support the

5

exploration of *in vitro* protein modifications

6

7

8

George W. Preston<sup>1,2,#a\*</sup>, Liping Yang<sup>2</sup>, David H. Phillips<sup>1</sup>, Claudia S. Maier<sup>2\*</sup>

9

10

11

12

<sup>1</sup> MRC-PHE Centre for Environment & Health, Department of Analytical, Environmental & Forensic Sciences, School of Population Health & Environmental Sciences, Faculty of Life Sciences & Medicine, King's College London, London, United Kingdom

15

16

<sup>2</sup> Department of Chemistry, Oregon State University, Corvallis, OR, United States of America

17

18

<sup>#a</sup> Current address: Stoller Biomarker Discovery Centre, University of Manchester, Manchester, United Kingdom

19

20

21

22

\* Corresponding authors

23

E-mail: george.preston@kcl.ac.uk (GWP), claudia.maier@oregonstate.edu (CSM)

24

## 25 **Abstract**

26           Dependent peptide searching is a method for discovering covalently-modified peptides – and  
27 therefore proteins – in mass-spectrometry-based proteomics experiments. Being more permissive than  
28 standard search methods, it has the potential to discover novel modifications (e.g., post-translational  
29 modifications occurring *in vivo*, or modifications introduced *in vitro*). However, few studies have  
30 explored dependent peptide search results in an untargeted way. In the present study, we sought to  
31 evaluate dependent peptide searching as a means of characterising proteins that have been modified *in*  
32 *vitro*. We generated a model data set by analysing *N*-ethylmaleimide-treated bovine serum albumin,  
33 and performed dependent peptide searches using the popular MaxQuant software. To facilitate  
34 interpretation of the search results (hundreds of dependent peptides), we developed a series of  
35 visualisation tools (R scripts). We used the tools to assess the diversity of putative modifications in  
36 the albumin, and to pinpoint hypothesised modifications. We went on to explore the tools' generality  
37 via analyses of public data from studies of rat and human proteomes. Of 19 expected sites of  
38 modification (one in rat cofilin-1 and 18 across six different human plasma proteins), eight were  
39 found and correctly localised. Apparently, some sites went undetected because chemical enrichment  
40 had depleted necessary analytes (potential 'base' peptides). Our results demonstrate (i) the ability of  
41 the tools to provide accurate and informative visualisations, and (ii) the usefulness of dependent  
42 peptide searching for characterising *in vitro* protein modifications. Our model data are available via  
43 PRIDE/ProteomeXchange (accession number PXD013040).

44

## 45 **Introduction**

46           By the time a protein is subjected to analysis, it can have acquired one or more covalent  
47 modifications. These could include modifications of biological origin, modifications introduced

48 deliberately (e.g., to probe protein structure and function), and modifications occurring during sample  
49 preparation and storage. In bottom-up mass-spectrometry-based proteomics, where proteins are  
50 digested and analysed as peptides, prior knowledge of modifications can enable more of the acquired  
51 spectra to be identified [1]. Known or suspected modifications are specified as parameters of a  
52 database search, enabling more of the protein sequence to be mapped, and also allowing the  
53 modifications themselves to be localised and quantified. For partially characterised or unknown  
54 modifications, however, this approach is not practical: specifying a long list of variable modifications  
55 (e.g., as a way of capturing unknown modifications) would expand the database dramatically,  
56 lengthening the search duration and reducing the number of confidently identified spectra [2]. New  
57 types of search have been developed to address this problem [2-9]. They include ‘open’ database  
58 searches, which permit precursor ions with shifted masses [6]; and ‘spectral pair’ searches, in which  
59 unidentified spectra are matched to spectral libraries [5]. An example of the latter approach is  
60 dependent peptide (DP) searching [3]. In a typical case, a DP is a chromatographic feature that is not  
61 identified by a database search, but whose fragment-ion spectrum partially matches that of one of the  
62 search hits (the ‘base’ peptide). The DP is typically a modified form of the base peptide, and the two  
63 features’ masses differ. In theory, some of the DP’s product ions will be the same as the  
64 corresponding product ions of the base peptide, while others will display the mass difference ( $\Delta m$ ).  
65 Crucially,  $\Delta m$  does not need to be specified *a priori*, as it is calculated for every pair of unidentified  
66 feature and database-search hit. Identifying features in this way can take much less time than a  
67 database search [3] but does confer certain limitations: sites that are fully occupied by unknown  
68 modifications cannot be detected; and overall sequence coverage is unlikely to be extended.

69           Originally implemented as a stand-alone tool (ModifiComb [3, 10]), DP searching has  
70 recently been incorporated into the MaxQuant software [11, 12]. Within MaxQuant, the DP search  
71 can utilise hits (i.e., potential base peptides) generated by the Andromeda search engine [13]. Studies  
72 utilising MaxQuant’s DP search function have confirmed its potential to discover modifications [12,  
73 14-21]. Lassak et al. used the function to analyse a bacterial translation elongation factor, and  
74 discovered a novel type of glycosylation [14]. Mordret et al. used the function to detect single amino

75 acid substitutions, and also carried out a general test of its validity [12]. A large set of synthetic  
76 phosphopeptides and corresponding unmodified peptides was analysed, and the phosphoryl  
77 modifications were left for a DP search to find. The search identified over a thousand spectra as  
78 belonging to singly-phosphorylated peptides, and all of these were true positives [12]. Few studies,  
79 however, have explored DP search results in an untargeted way. In the present study, we sought to  
80 evaluate DP searching as a means of characterising *in vitro* modified proteins. First, we generated  
81 model data by analysing a model protein (bovine serum albumin, BSA) that had been treated with a  
82 protein-modifying reagent (*N*-ethylmaleimide, NEM). Then, we performed DP searches and  
83 attempted to rationalise the search results.

84 Visualisation tools can greatly facilitate the interpretation of proteomic mass spectrometry  
85 data and database search results [22-25]. We identified a need for tools that visualise DP search  
86 results, and to meet this need we wrote a set of five scripts in the R language [26]. Three of the scripts  
87 are for surveying distributions of DPs (i.e., are hypothesis-generating), and the other two are for  
88 pinpointing hypothesised modifications (i.e., are hypothesis-testing). Some of the scripts can enrich  
89 DPs for modifications that are unique to a test sample. Herein we report search results and  
90 visualisations for our own data, as well as for public data from two other studies [27, 28]. The results  
91 demonstrate how a combination of DP searching and visualisation can assist in the characterisation of  
92 *in vitro* modified proteins. The approach could be useful for characterising protein targets of enzyme  
93 activities and reactive small molecules.

94

## 95 **Materials and methods**

### 96 **Preparation and analysis of modified BSA**

97 BSA (1 mg mL<sup>-1</sup>) was reacted with NEM (1 mM) in potassium phosphate buffer (100 mM) at  
98 pH 7.4. Unreacted NEM was scavenged with 1,4-dithiothreitol (DTT). The protein was purified

99 (buffer exchange), reduced (DTT), alkylated (iodoacetamide), purified again (acetone precipitation),  
100 and digested (trypsin). The peptides were analysed in duplicate (analytical replicates 1 and 2) by  
101 reversed-phase nano liquid chromatography (nanoACQUITY liquid chromatograph; Waters, Milford,  
102 Massachusetts, USA) with online data-dependent tandem mass spectrometry (Orbitrap Fusion Lumos  
103 mass spectrometer; ThermoFisher Scientific, Waltham, Massachusetts, USA). A control sample  
104 (untreated BSA) was also prepared and analysed. Further details of materials and methods can be  
105 found in S1 Text. Data for NEM-treated and untreated BSA have been deposited in  
106 PRIDE/ProteomeXchange [29] (accession number PXD013040).

## 107 **Public data**

108 Further mass spectrometry data were obtained from PRIDE/ProteomeXchange [29]. Data  
109 were selected according to the following criteria: (i) experiment involving exposure of one or more  
110 proteins to a protein-modifying reagent; (ii) data collected using standard data-dependent acquisition  
111 mass spectrometry; (iii) control data and two replicates available (not an essential criterion); (iv)  
112 results of a variable-modification database search reported in the literature. The following files/groups  
113 of files met the criteria and were included in the present study (the selection was not exhaustive):  
114 DMF\_Cofilin1A.raw from PRIDE project PXD008314 [27]; and 1362-cs774\_0\_a.raw, 1364-  
115 cs774\_0\_b.raw, 1380-cs774\_5\_a.raw and 1382-cs774\_5\_b.raw from PRIDE project PXD006663  
116 [28]. Sequence data (\*.fasta files) were obtained from the RCSB Protein Data Bank [30] (accession  
117 numbers 4F5S [31] and 1S81 [32]) and UniProt [33] (accession numbers P02042, P02647, P02766,  
118 P02768, P02787, P45592 and P68871, and the UniProtKB/Swiss-Prot human proteome, 4<sup>th</sup> June  
119 2019). Where possible, sequences were obtained without extraneous elements such as signal peptides.  
120 The data from Protein Data Bank accession number 4F5S consisted of two identical sequences, and so  
121 one of these was removed. MaxQuant's database of contaminants [11] was used either as supplied or  
122 in an edited form (see S1 Text).

123

## 124 **Database searches and dependent peptide searches**

125 All searches were done in MaxQuant (Max Planck Institute of Biochemistry, version 1.6.0.1)  
126 [11]. Database searches were done using Andromeda [13]. Individual \*.raw files were searched  
127 against databases consisting of either a protein of interest plus potential contaminants (1 + 244 or 1 +  
128 245 sequences) or the UniProtKB/Swiss-Prot human proteome plus potential contaminants (20,406 +  
129 82 sequences). The *in silico* digestion was done in ‘specific’ mode, using ‘Trypsin/P’ as the enzyme,  
130 and allowing for a maximum of two missed cleavages. The maximum peptide mass was adjusted so as  
131 to include  $\geq 95\%$  of relevant theoretical peptides (see S1 Text). The minimum peptide length was  
132 seven amino acid residues. Additional *in silico* digestions were done using PeptideMass [34] (see S1  
133 Text).

134 When the purpose of the database search was to discover potential base peptides, a minimal  
135 set of variable modifications (methionine oxidation and protein N-terminal acetylation) and an  
136 appropriate fixed modification (cysteine *S*-carbamidomethylation [28] or *S*-pyridylethylation [27])  
137 were specified. When the purpose was instead to localise a particular modification (a variable-  
138 modification search), the fixed and variable modifications were adjusted accordingly (see S1 Text).  
139 The maximum number of modifications per peptide was always five, and the ‘second peptides’  
140 function was always enabled. DP searches were appended to their respective database searches by  
141 enabling the ‘dependent peptides’ function (false discovery rate of 1%, mass bin size of 0.0065 Da).  
142 Results of DP searches were obtained from allPeptides.txt files [11] and filtered as described below.  
143 For selected DPs, we investigated whether the same chromatographic feature could also be detected  
144 by a variable-modification search (see S1 Text).

## 145 **Development of visualisation tools**

146 All tools (Scripts I-V; S1-5 Scripts) were developed in R for Windows (R Core Team, version  
147 3.4.0 or later) [26] using functions from the base packages, plus the additional function *read.fasta*  
148 from the ‘seqinR’ package (version 3.3-6 or later) [35]. Scripts I-IV were developed and tested on a

149 Dell desktop PC (Intel Core i5-7500 processor, 8 GB RAM) running Windows 10. Script V was  
150 developed and tested on a Toshiba laptop PC (AMD E1-2100 APU processor, 4 GB RAM) running  
151 Windows 8.1. Each script requires a set of search results (allPeptides.txt files), a protein sequence  
152 (\*.fasta file), and the identifier of a protein of interest (e.g., a UniProt identifier). The search results  
153 are filtered (see Table B of S1 Text) and a table of DPs is prepared. DPs are localised to segments of  
154 the protein sequence using a 'sliding window' [36]. Different scripts require different numbers of  
155 allPeptides.txt files and process the data to different extents. Scripts I, II and III are for surveying  
156 DPs'  $\Delta m$  values; they return DP localisation plots and  $\Delta m$  frequency histograms. Scripts IV and V are  
157 for pinpointing particular modifications; they require an expected  $\Delta m$  value, and they return plots of  
158 DPs' localisation probabilities [12]. Scripts II, III and IV have the ability to enrich the DPs for  
159 modifications that are unique to a test sample (see 'Results and discussion'). A set of notes explaining  
160 how the scripts work can be found in Table C of S1 Text. We will endeavour to maintain the scripts'  
161 compatibility with relevant software, and to address any limitations that come to light. Any future  
162 versions of the scripts will be made available from <https://github.com/preston-gw/>.

163 The accuracy of the visualisations was tested by manually mapping selected DPs onto  
164 graphics generated by scripts (20 DPs across five graphics, all confirmed to have been mapped  
165 correctly). After this, no significant changes were made to either the main data processing code or the  
166 mapping code (changes to graphical parameters, for example, were permitted). Certain graphics were  
167 prepared in batches by iterating an appropriate script. Figures were prepared from R output files using  
168 Inkscape (Free Software Foundation, version 0.91 or later) and GNU Image Manipulation Program  
169 (GIMP Development Team, version 2.10.8) (see S1 Text). Figures such as S4 Fig and S8 Fig are  
170 representative of the graphics generated by the scripts.

## 171 **Visualisation tools' instructions for use**

- 172 1. Open R (version 3.6.0 or later)
- 173 2. Make sure that package 'seqinR' [35] is installed. Installation is achieved by entering  
174 *install.packages("seqinr")* in the R console.



- 175 3. Open the script within R (File > Open script)
- 176 4. Review the script and complete file paths as appropriate. Further instructions are included in
- 177 the header and body of the script. Explanatory notes can be found in Tables B and C of S1
- 178 Text.
- 179 5. Save a copy of the script if desired (File > Save as)
- 180 6. Run the script (Edit > Run all)
- 181 7. A graphic should appear on-screen. The times taken to visualise our model data using one of
- 182 the aforementioned PCs were as follows: Script I, 8 s (desktop); Script II, 27 s (desktop),
- 183 Script III, 62 s (desktop); Script IV, 16 s (desktop); Script V, 26 s (laptop).

## 184 **Mass calculations**

185 Expected  $\Delta m$  values were calculated in R (version 3.4.0 or later) [26] using monoisotopic

186 masses from ChemDraw (various versions, PerkinElmer) or Unimod [37]

187 (<http://www.unimod.org/masses.html>). Masses used for calculations were accurate to at least four

188 decimal places.  $\Delta m$  values mentioned in the text have been rounded to two decimal places.

## 189 **Statistical methods**

190 Pairwise comparisons of  $\Delta m$  frequency histograms were done using the *cor.test* function

191 (Spearman method) in R (version 3.4.0) [26].

192

193

194

195

## 196 **Results and discussion**

### 197 **Exploration of model data**

198 BSA was selected as a model protein because it is well-characterised and contains a number  
199 of nucleophilic (i.e., potentially modifiable) amino acid residues [38, 39]. NEM was selected as the  
200 protein-modifying reagent because it is reactive towards a variety of amino acid side chains (those of  
201 cysteine, lysine and histidine) [40]. We predicted that NEM would modify BSA's only reduced  
202 cysteine residue, Cys-34 [31], as well as one or more lysine and/or histidine residues. The primary  
203 products of the reaction of BSA with NEM were expected to be Michael adducts, in which a hydrogen  
204 atom of the protein has effectively been replaced by an *N*-ethylsuccinimidyl (NESyl) group ( $\Delta m =$   
205  $+125.05$  Da). NESyl groups attached to cysteine residues are susceptible to hydrolysis (additional  $\Delta m$   
206  $= +18.01$  Da) [41, 42], and we assumed that this would also be the case for NESyl groups attached to  
207 lysine or histidine residues. Additionally, sulfur atoms to which NESyl groups are attached may  
208 oxidise [41].

209 Large numbers of chromatographic features were detected in analyses of NEM-treated BSA  
210 ( $N \geq 34,812$ ), and also in analyses of untreated BSA ( $N \geq 44,430$ ). Five to six percent of the features  
211 ( $1958 \leq N \leq 2328$ ) were identified by MaxQuant as either potential base peptides (Andromeda search  
212 hits, 23-25% of identified features) or DPs (75-77% of identified features) (S1 Fig). Peptides of BSA  
213 (95% of identified features) were 9-10 times as numerous as predicted by *in silico* digestion and  
214 oxidation ( $N = 218$ ). The high ratio of observed to expected features implies that large numbers of  
215 modifications had occurred independently of NEM treatment (e.g., artefacts of sample preparation or  
216 modifications pre-existing in the BSA). The detection of so many 'background' modifications,  
217 although difficult to account for, is consistent with Nielsen and coauthors' estimate of 8-12 modified  
218 peptides per unmodified tryptic peptide [10]. Filters were employed to isolate the DPs, to limit the  
219 number of 'background' modifications (see Table B of S1 Text), and to limit  $\Delta m$  to  $\pm 500$  Da (for  
220 clarity of visualisation). Filtering removed 70-74% of the identified features (S1 Fig).

## 221 **Visualisation of $\Delta m$ distributions**

222           The filtered  $\Delta m$  values were visualised in two ways: firstly by mapping DPs to segments of  
223 the protein sequence, and secondly using a frequency histogram [3, 25]. These modes of visualisation,  
224 both achieved using Script I, revealed a diversity of putative modifications in NEM-treated BSA (Fig  
225 1A, B). Visualisations of putative modifications to porcine trypsin, also detected in analyses of NEM-  
226 treated BSA, demonstrated the flexibility of Script I (S2 Fig).

227 **Fig 1. Localisation plots and mass-shift ( $\Delta m$ ) frequency histograms for dependent peptides**  
228 **(DPs) of *N*-ethylmaleimide-treated bovine serum albumin.** In the localisation plots (left-hand  
229 panels), the protein sequence is represented as a dashed line that becomes solid in regions for which  
230 peptides were observed. X-axis values refer to positions in the protein sequence (position 1 = N-  
231 terminal amino acid residue). Each DP is represented as a rectangle whose height is proportional to  
232  $\Delta m$ , and whose grey border is partially transparent. The  $\Delta m$  values are summarised in frequency  
233 histograms (right-hand panels). The DPs were unenriched (A, B), enriched using Script II (C, D) or  
234 enriched using Script III (E, F).

235            $\Delta m$  frequency histograms were used to investigate whether modifications in NEM-treated and  
236 untreated BSA were the same. Histograms for the two samples were similar (nearly as similar, or  
237 more similar, than histograms for analytical replicates; see S3 Fig). This suggested that the samples  
238 had many modifications in common. In order to selectively visualise the NEM-derived modifications,  
239 we investigated ways of enriching DPs. We started with a method (Script II) that subtracts the DPs  
240 observed in an analysis of untreated protein from those observed in an analysis of treated protein (Fig  
241 1C, D). For this purpose, DPs were regarded as simple combinations of sequence and  $\Delta m$  (in  
242 principle, retention time could also be used, but this was not attempted here). We then took the  
243 enrichment idea a stage further (Script III) by looking for DPs that were ‘constantly conjoined’ with  
244 NEM treatment (i.e., observed in both analyses of NEM-treated BSA, but in neither analysis of  
245 untreated BSA). Enrichment was quantified as an increase in the percentage of DPs having either of  
246 two putative NEM-derived groups: intact NESyl ( $\Delta m \pm \text{tolerance} = +125.05 \pm 0.01$  Da) or hydrolysed

247 NESyl ( $\Delta m \pm \text{tolerance} = +143.06 \pm 0.01 \text{ Da}$ ). Both modifications were observed for NEM-treated  
248 BSA (8.7% of the DPs from analysis 1, Fig 1B) and neither was observed for untreated BSA. Script II  
249 effected 2.3-fold enrichment of DPs from analysis 1 of NEM-treated BSA (Fig 1D). Script III effected  
250 6.0-fold enrichment of these same DPs, but its greater stringency led to the exclusion of eleven  
251 relevant DPs (Fig 1F). The results suggest that our scripts should be able to enrich DPs even when the  
252 modifications of interest are unknown. In other words, the scripts might be able to discover novel  
253 modifications and attribute them to a given reagent or condition.

## 254 **Visualisation of expected modifications**

255 As well as surveying the diversity of modifications, we were also interested in visualising the  
256 distributions of specific hypothesised modifications. For this purpose, we developed a method that maps  
257 ‘constantly conjoined’ DPs to the protein sequence and highlights their probable sites of modification  
258 (Script IV). A sliding window is used as before, but in Script IV its role is to direct the entry of values  
259 into matrices. Each DP is ‘etched’ (as a line of ones) into a blank matrix, and the localisation  
260 probabilities for that DP’s modification are deposited in a corresponding zero matrix (localisation  
261 probability is a computed measure of the likelihood of a modification occurring at a given site [11, 12]).  
262 The two matrices are subsequently converted to images (R function *image* [26]) and merged. Fig 2 is a  
263 formatted version of one of the resulting graphics, showing the distribution of putative hydrolysed  
264 NESyl groups. S4 Fig is an example of an unformatted graphic, showing the distribution of putative  
265 intact NESyl groups. Note how some of the DPs in Fig 2 occur as pairs of putative diastereoisomers  
266 [43] (i.e., modified peptides with identical  $m/z$  values and fragmentation patterns but different retention  
267 times).

268 **Fig 2. Localisation plot for putative hydrolysed *N*-ethylsuccinimidyl groups in *N*-ethylmaleimide-**  
269 **treated bovine serum albumin ( $\Delta m \pm \text{tolerance} = +143.06 \pm 0.01 \text{ Da}$ ).** The protein sequence is  
270 represented as a dashed line that becomes solid in regions for which peptides were observed. X-axis  
271 values refer to positions in the protein sequence (position 1 = N-terminal amino acid residue).  
272 Dependent peptides are represented as coloured strips with shading to indicate the localisation

273 probability (darker = more probable). Any site with a non-zero probability is annotated. One dependent  
274 peptide (amino acid residues 336-347, probably modified at Arg-336 or His-337) does not appear  
275 because the relevant matrices were full.

276 Using the hydrolysed NESyl group as an example, we examined whether the putative  
277 modifications had been localised to plausible sites. In 23 of 34 DPs with putative hydrolysed NESyl  
278 groups (68%), the highest localisation probability had been assigned to a cysteine, histidine or lysine  
279 residue (S1 Table). In cases where the same probability had been assigned to multiple sites, we used a  
280 system of prioritisation to narrow down the possibilities (see S1 Text for details). Six of the 23  
281 plausible localisations were confirmed by a variable-modification search (S1 Table). In all six cases,  
282 the modification could be localised to a histidine or lysine residue. Modifications to cysteine residues,  
283 which represent a special case, are discussed below. A seventh, less plausible localisation (to Asp-13)  
284 was also confirmed by the variable-modification search (S1 Table).

285 In order to discover the DPs with modified cysteine residues, we had to account for the fact  
286 that the corresponding residues in the base peptides would also be modified (*S*-carbamidomethylated).  
287 Subtracting the  $\Delta m$  for carbamidomethylation gave new values for NESyl (+68.03 Da) and  
288 hydrolysed NESyl (+86.04 Da), both of which we recognised from the  $\Delta m$  frequency histograms (Fig  
289 1). Surprisingly, neither modification was localised to Cys-34 (S5 and S6 Figs), and no modified  
290 cysteine residues were confirmed by the variable-modification search. Seeking to understand the  
291 apparent absence of modifications to Cys-34, we turned to a group of unexplained DPs ( $\Delta m = -25.03$   
292 Da; Fig 1F), which we speculated might contain oxidised cysteine residues (cysteinesulfinic acid).  
293 The modification was found to have been localised to Cys-34 in some DPs, but none that were  
294 'constantly conjoined' with NEM treatment (S7 Fig). The ambiguous results for Cys-34 are possibly a  
295 consequence of modifications at this site having decomposed prior to or during analysis. We did see  
296 some evidence of modification to cysteine residues other than Cys-34 (S1 Table, S5-7 Figs), and this  
297 was unexpected because these residues are normally disulfide-bonded to other cysteine residues [31].  
298 It is possible that modifications to cysteine residues other than Cys-34 occurred when DTT was added

299 to scavenge unreacted NEM. It is also possible that some of the other modifications observed in the  
300 study occurred following this addition of DTT.

301 The above results highlight the fact that observed  $\Delta m$  values do not always correspond to real  
302 chemical transformations, and cannot always be interpreted directly. Direct interpretation is permitted  
303 if the DP is singly-modified and the base peptide is truly unmodified, and likewise if the DP and base  
304 peptide contain modifications that ‘cancel out’ (e.g., cysteine *S*-carbamidomethylation). However, if  
305 the base peptide contains modifications not found in the DP, or if the base peptide and DP contain  
306 different modifications at the same site, then interpretation will be less straightforward. Problems of  
307 this nature could be avoided by omitting modifications from the database search, but this would of  
308 course restrict the overall number of identifications.

## 309 **Validation of tools using public data**

310 We explored the scripts’ generality by applying them to analyses of public data. We  
311 hypothesised that a combination of DP searching and visualisation would reveal the same adducts as  
312 other authors had found by variable-modification database searching. We expected to observe these  
313 adducts via targeted visualisation (e.g., using Script IV). We also predicted that they would be evident  
314 from an untargeted survey (e.g., using Script III).

315 First, we analysed data from a study by Piroli et al. [27]. In this work, cultured cells (rat  
316 astrocytes) were exposed to the drug dimethyl fumarate. The authors collected proteins from the  
317 exposed cells, resolved them on gels, and then subjected individual protein bands to reduction,  
318 alkylation (4-vinylpyridine), digestion (trypsin) and analysis. In one of the protein bands, the authors  
319 detected a monomethyl fumarate adduct of cofilin-1 and localised the modification to a cysteine  
320 residue (Cys-139). Using the methods developed for the analysis of NEM-treated BSA, we performed  
321 an independent analysis of Piroli and coauthors’ cofilin-1 data. Script I was used to survey  $\Delta m$  values,  
322 and an additional script was developed for mapping the expected  $\Delta m$  (Script V, an analogue of Script  
323 IV that works with single allPeptides.txt files). Script I did not highlight the expected nominal  $\Delta m$

324 (+25 Da; S8 Fig), which is perhaps understandable given this script's inability to enrich DPs. The  
325 overall sparsity of S8 Fig could reflect a real lack of modifications in the rat cofilin-1, or alternatively  
326 it could reflect qualities of the sample and/or data. Script V revealed that the expected  $\Delta m$  ( $+24.97 \pm$   
327  $0.01$  Da) was present in one DP, and that it had been localised to the correct cysteine residue (S9 Fig).  
328 The  $\Delta m$  itself is also evidence of correct localisation, since this is the difference in mass with respect  
329 to a pyridylethylated base peptide. The discrepancy in the site numbers (138 versus 139) arose  
330 probably because we used the sequence of mature rat cofilin-1 (no N-terminal methionine residue)  
331 whereas Piroli et al. used the sequence of the full-length protein (we used mature sequences where  
332 possible to ensure that N-terminal peptides could be found by Andromeda).

333 Further data were from a study by Salomón et al. [28]. In this work, the reactivity of human  
334 blood proteins towards a metabolite, 3-hydroxy-2,5-hexanedione, was explored using an alkyne-  
335 functionalised probe ('alk-3-HHD'). The authors prepared plasma from probe-treated blood and  
336 collected the plasma proteins. The proteins were reduced, alkylated (iodoacetamide) and digested  
337 (trypsin), and the resulting peptides were chemically enriched for alk-3-HHD adducts. The authors  
338 detected two different types of modification ('HTO' and 'HDMP', both specific to lysine residues)  
339 among six polypeptides (apolipoprotein A-I, haemoglobin  $\beta$ - and  $\delta$ -chains, serotransferrin, serum  
340 albumin, and transthyretin). In total, the authors detected 18 unique sites of modification among 22  
341 analytes. Again, we applied methods that had been developed for the BSA adducts. Salomón and  
342 coauthors' dataset included replicates and a control, permitting the use of Scripts III and IV. To  
343 maximise contrast, we used data for the highest concentration of alk-3-HHD. Script III revealed  
344 multiple 'constantly conjoined' DPs of each of the six aforementioned polypeptides, but did not  
345 highlight any of the expected  $\Delta m$  values (Fig 3A, B; S10-14 Figs). Six DPs mapped to both of the  
346 haemoglobin chains, and therefore could not be localised unambiguously. Script IV was used to map  
347 HTO- and HDMP-type modifications (with or without sulfonation [28]) to the sequences of the  
348 polypeptides. In total, 14 DPs with putative alk-3-HHD-derived modifications were detected (11  
349 unique combinations of sequence and modification) (S2 Table). In each of the DPs, a lysine residue  
350 was either the site with the highest localisation probability, or was one of multiple such sites. Most of

351 the DPs (93%) were of either haemoglobin  $\beta$ -chain (Fig 3C) or serum albumin (S15 Fig). Of the 22  
352 analytes reported by Salomón et al., we detected six (27%) as DPs. Of the 18 expected sites of  
353 modification, we observed seven (39%). It is perhaps unsurprising that some of the expected DPs  
354 were not detected, since the chemical enrichment performed by Salomón et al. had the potential to  
355 remove their corresponding base peptides (an effect alluded to by Tyanova et al. in their protocol  
356 [11]). Indeed, for seven of the 22 expected analytes (32%), the absence of a required base peptide was  
357 sufficient to explain the absence of the DP.

358 **Fig 3. Visualisation of dependent peptides of human haemoglobin  $\beta$ -chain using public data**  
359 **from the study by Salomón et al. [28].** (A) Dependent peptide localisation plot. (B) Mass-shift ( $\Delta m$ )  
360 frequency histogram. (C) Localisation plot for putative sulfonated ‘HDMP’-type modifications ( $\Delta m \pm$   
361 tolerance =  $+390.10 \pm 0.01$  Da) [28]. Six of 106 dependent peptides also mapped to the sequence of  
362 haemoglobin  $\delta$ -chain (S11 Fig). Base peptide VLGA $\beta$ SDGLAHL $\beta$ DNLKGT $\beta$ FATLSELHCDK went  
363 undetected in analyses of untreated proteins and therefore does not appear (see Table C of S1 Text).

## 364 **Scope for extending the present study**

365 There is scope beyond the present study for developing and integrating the visualisation tools.  
366 As R scripts, they are highly amenable to modification, and could be adapted for more specialised  
367 purposes. The filters and the  $\Delta m$  tolerance could be made more stringent or permissive as required. The  
368 plots could be customised by changing the colour scheme or narrowing the limits of the  $\Delta m$  axes.  
369 Another idea would be to rotate the histogram so as to align its  $\Delta m$  axis with that of the DP  
370 localisation plot.

371 One area in which there is significant scope for development is annotation. Currently, the  
372 scripts highlight mass shifts but do not attempt to identify them. Some identifications are already  
373 made by MaxQuant, and these could be easily transferred onto the plots. Further identifications could  
374 be made via public protein-modification databases such as Unimod [37] and RESID [44]. These  
375 databases contain calculated  $\Delta m$  values via which observed  $\Delta m$  values could be linked to



376 modifications' identities. Identities could be added to the plots or visualised separately (e.g., a word-  
377 cloud of modifications' names). Another idea would be to highlight particular features of the protein  
378 sequence, such as nucleophilic amino acid residues. This could be done using lines, symbols or text.

379         Currently, each script visualises results for a single protein (Scripts I-III) or combination of  
380 protein and  $\Delta m$  (Scripts IV and V). One way of extending the approach would be to iterate scripts so  
381 that they cycle through lists of proteins and/or  $\Delta m$  values (in fact, we prepared certain groups of  
382 figures in this way). In theory, this could be done in a 'data-dependent' fashion by extracting the lists  
383 directly from allPeptides.txt. If this were attempted, each graphic (e.g., \*.svg file) would have to be  
384 stamped with the protein identifier and/or  $\Delta m$  value.

385         Scripts II, III and IV are able to enrich DPs for modifications that are unique to a test sample.  
386 The modes of enrichment employed by the scripts are simple but should work well for many *in vitro*  
387 modifications (especially modifications not found *in vivo*). The visualisation of *in vivo* modifications  
388 would be an impactful next step, but one that might require a more quantitative approach: it would be  
389 helpful to visualise differences in abundance in addition to the simple difference between presence  
390 and absence.

391         Finally, there is scope for combining the scripts in an R package [45]. This would promote  
392 their usability beyond the present study.

## 393 **Summary**

394         We have developed and tested a set of analytical tools with which to interpret the results of a  
395 DP search. The tools visualise putative modifications ( $\Delta m$  values) for a protein of interest (either an  
396 isolated protein or a component of a proteome). Some of the tools are able to enrich DPs for  
397 modifications that are unique to a test sample. We envisaged that an untargeted survey of DPs (using  
398 Scripts I-III) might generate hypotheses that could then be tested via targeted visualisation (using  
399 Scripts IV and V). This approach helped us to achieve our aim of rationalising DP search results for  
400 NEM-treated BSA. Expected modifications were found, and the majority of these were localised to

401 chemically plausible sites. In formal tests involving public data, a number of expected modifications  
402 were detected and correctly localised (although here the methods for surveying  $\Delta m$  values proved less  
403 helpful than they had for the BSA study). Analyses of cysteine-specific modifications led us to  
404 consider the effect of fixed and variable modifications on  $\Delta m$ ; and analyses of data for chemically  
405 enriched peptides led us to consider the potential of chemical enrichment to limit DP searches. We  
406 conclude (i) that the tools can summarise DP search results accurately and informatively, and (ii) that  
407 DP searching can be useful for characterising *in vitro* modified proteins.

408

## 409 **References**

- 410 1. Cottrell JS. Protein identification using MS/MS data. *J Proteomics*. 2011;74(10):1842-51. doi:  
411 10.1016/j.jprot.2011.05.014.
- 412 2. Ahrné E, Müller M, Lisacek F. Unrestricted identification of modified proteins using MS/MS.  
413 *Proteomics*. 2010;10(4):671-86. doi: 10.1002/pmic.200900502.
- 414 3. Savitski MM, Nielsen ML, Zubarev RA. ModifiComb, a New Proteomic Tool for Mapping  
415 Substoichiometric Post-translational Modifications, Finding Novel Types of Modifications, and  
416 Fingerprinting Complex Protein Mixtures. *Mol Cell Proteomics*. 2006;5(5):935-48. doi:  
417 10.1074/mcp.T500034-MCP200.
- 418 4. Ahrné E, Nikitin F, Lisacek F, Müller M. QuickMod: A Tool for Open Modification  
419 Spectrum Library Searches. *J Proteome Res*. 2011;10(7):2913-21. doi: 10.1021/pr200152g.
- 420 5. Na S, Paek E. Software eyes for protein post-translational modifications. *Mass Spectrom Rev*.  
421 2015;34(2):133-47. doi: 10.1002/mas.21425.
- 422 6. Chick JM, Kolippakkam D, Nusinow DP, Zhai B, Rad R, Huttlin EL, et al. A mass-tolerant  
423 database search identifies a large proportion of unassigned spectra in shotgun proteomics as modified  
424 peptides. *Nat Biotechnol*. 2015;33:743-9. doi: 10.1038/nbt.3267.

- 425 7. Bagwan N, Bonzon-Kulichenko E, Calvo E, Lechuga-Vieco AV, Michalakopoulos S,  
426 Trevisan-Herraz M, et al. Comprehensive Quantification of the Modified Proteome Reveals Oxidative  
427 Heart Damage in Mitochondrial Heteroplasmy. *Cell Rep.* 2018;23(12):3685-97. doi:  
428 10.1016/j.celrep.2018.05.080.
- 429 8. Devabhaktuni A, Lin S, Zhang L, Swaminathan K, Gonzalez CG, Olsson N, et al. TagGraph  
430 reveals vast protein modification landscapes from large tandem mass spectrometry datasets. *Nat*  
431 *Biotechnol.* 2019;37(4):469-79. doi: 10.1038/s41587-019-0067-5.
- 432 9. Na S, Kim J, Paek E. MODplus: Robust and Unrestrictive Identification of Post-Translational  
433 Modifications Using Mass Spectrometry. *Anal Chem.* 2019;91(17):11324-33. doi:  
434 10.1021/acs.analchem.9b02445.
- 435 10. Nielsen ML, Savitski MM, Zubarev RA. Extent of Modifications in Human Proteome  
436 Samples and Their Effect on Dynamic Range of Analysis in Shotgun Proteomics. *Mol Cell*  
437 *Proteomics.* 2006;5(12):2384-91. doi: 10.1074/mcp.M600248-MCP200.
- 438 11. Tyanova S, Temu T, Cox J. The MaxQuant computational platform for mass spectrometry-  
439 based shotgun proteomics. *Nat Protoc.* 2016;11:2301-19. doi: 10.1038/nprot.2016.136.
- 440 12. Mordret E, Dahan O, Asraf O, Rak R, Yehonadav A, Barnabas GD, et al. Systematic  
441 Detection of Amino Acid Substitutions in Proteomes Reveals Mechanistic Basis of Ribosome Errors  
442 and Selection for Translation Fidelity. *Mol Cell.* 2019;75(3):427-41. doi:  
443 10.1016/j.molcel.2019.06.041.
- 444 13. Cox J, Neuhauser N, Michalski A, Scheltema RA, Olsen JV, Mann M. Andromeda: A Peptide  
445 Search Engine Integrated into the MaxQuant Environment. *J Proteome Res.* 2011;10(4):1794-805.  
446 doi: 10.1021/pr101065j.
- 447 14. Lassak J, Keilhauer EC, Fürst M, Wuichet K, Gödeke J, Starosta AL, et al. Arginine-  
448 rhamnosylation as new strategy to activate translation elongation factor P. *Nat Chem Biol.*  
449 2015;11:266-70. doi: 10.1038/nchembio.1751.
- 450 15. Soufi B, Krug K, Harst A, Macek B. Characterization of the E. coli proteome and its  
451 modifications during growth and ethanol stress. *Front Microbiol.* 2015;6:103. doi:  
452 10.3389/fmicb.2015.00103.

- 453 16. Bogdanow B, Zauber H, Selbach M. Systematic Errors in Peptide and Protein Identification  
454 and Quantification by Modified Peptides. *Mol Cell Proteomics*. 2016;15(8):2791-801. doi:  
455 10.1074/mcp.M115.055103.
- 456 17. Cvetesic N, Semanjski M, Soufi B, Krug K, Gruic-Sovulj I, Macek B. Proteome-wide  
457 measurement of non-canonical bacterial mistranslation by quantitative mass spectrometry of protein  
458 modifications. *Sci Rep*. 2016;6:28631. doi: 10.1038/srep28631.
- 459 18. Rykær M, Svensson B, Davies MJ, Hägglund P. Unrestricted Mass Spectrometric Data  
460 Analysis for Identification, Localization, and Quantification of Oxidative Protein Modifications. *J*  
461 *Proteome Res*. 2017;16(11):3978-88. doi: 10.1021/acs.jproteome.7b00330.
- 462 19. van Mierlo G, Wester RA, Marks H. Quantitative subcellular proteomics using SILAC  
463 reveals enhanced metabolic buffering in the pluripotent ground state. *Stem Cell Res*. 2018;33:135-45.  
464 doi: 10.1016/j.scr.2018.09.017.
- 465 20. Tiwari MK, Hägglund PM, Møller IM, Davies MJ, Bjerrum MJ. Copper ion / H<sub>2</sub>O<sub>2</sub> oxidation  
466 of Cu/Zn-Superoxide dismutase: Implications for enzymatic activity and antioxidant action. *Redox*  
467 *Biol*. 2019;26:101262. doi: 10.1016/j.redox.2019.101262.
- 468 21. Sinitcyn P, Rudolph JD, Cox J. Computational Methods for Understanding Mass  
469 Spectrometry-Based Shotgun Proteomics Data. *Annu Rev Biomed Data Sci*. 2018;1(1):207-34. doi:  
470 10.1146/annurev-biodatasci-080917-013516.
- 471 22. Martín-Campos T, Mylonas R, Masselot A, Waridel P, Petricevic T, Xenarios I, et al. MsViz:  
472 A Graphical Software Tool for In-Depth Manual Validation and Quantitation of Post-translational  
473 Modifications. *J Proteome Res*. 2017;16(8):3092-101. doi: 10.1021/acs.jproteome.7b00194.
- 474 23. Avtonomov DM, Kong A, Nesvizhskii AI. DeltaMass: Automated Detection and  
475 Visualization of Mass Shifts in Proteomic Open-Search Results. *J Proteome Res*. 2019;18(2):715-20.  
476 doi: 10.1021/acs.jproteome.8b00728.
- 477 24. Tyanova S, Temu T, Carlson A, Sinitcyn P, Mann M, Cox J. Visualization of LC-MS/MS  
478 proteomics data in MaxQuant. *Proteomics*. 2015;15(8):1453-6. doi: 10.1002/pmic.201400449.

- 479 25. Tyanova S, Temu T, Sinitcyn P, Carlson A, Hein MY, Geiger T, et al. The Perseus  
480 computational platform for comprehensive analysis of (prote)omics data. *Nat Methods*.  
481 2016;13(9):731-40. doi: 10.1038/nmeth.3901.
- 482 26. R Core Team. A language and environment for statistical computing. Vienna, Austria: R  
483 Foundation for Statistical Computing; 2017.
- 484 27. Piroli GG, Manuel AM, Patel T, Walla MD, Shi L, Lanci SA, et al. Identification of Novel  
485 Protein Targets of Dimethyl Fumarate Modification in Neurons and Astrocytes Reveals Actions  
486 Independent of Nrf2 Stabilization. *Mol Cell Proteomics*. 2019;18(3):504-19. doi:  
487 10.1074/mcp.RA118.000922.
- 488 28. Salomón T, Sibbersen C, Hansen J, Britz D, Svart MV, Voss TS, et al. Ketone Body  
489 Acetoacetate Buffers Methylglyoxal via a Non-enzymatic Conversion during Diabetic and Dietary  
490 Ketosis. *Cell Chem Biol*. 2017;24(8):935-43. doi: 10.1016/j.chembiol.2017.07.012.
- 491 29. Perez-Riverol Y, Csordas A, Bai J, Bernal-Llinares M, Hewapathirana S, Kundu DJ, et al.  
492 The PRIDE database and related tools and resources in 2019: improving support for quantification  
493 data. *Nucleic Acids Res*. 2019;47(D1):D442-50. doi: 10.1093/nar/gky1106.
- 494 30. Berman HM, Westbrook J, Feng Z, Gilliland G, Bhat TN, Weissig H, et al. The Protein Data  
495 Bank. *Nucleic Acids Res*. 2000;28(1):235-42. doi: 10.1093/nar/28.1.235.
- 496 31. Bujacz A. Structures of bovine, equine and leporine serum albumin. *Acta Crystallogr D*.  
497 2012;68(10):1278-89. doi: 10.1107/S0907444912027047.
- 498 32. Transue TR, Krahn JM, Gabel SA, DeRose EF, London RE. X-ray and NMR characterization  
499 of covalent complexes of trypsin, borate, and alcohols. *Biochemistry*. 2004;43(10):2829-39. doi:  
500 10.1021/bi035782y.
- 501 33. The UniProt Consortium. UniProt: a worldwide hub of protein knowledge. *Nucleic Acids*  
502 *Res*. 2018;47(D1):D506-15. doi: 10.1093/nar/gky1049.
- 503 34. Wilkins MR, Lindskog I, Gasteiger E, Bairoch A, Sanchez J-C, Hochstrasser DF, et al.  
504 Detailed peptide characterization using PEPTIDEMASS – a World-Wide-Web-accessible tool.  
505 *Electrophoresis*. 1997;18(3-4):403-8. doi: 10.1002/elps.1150180314.

- 506 35. Charif D, Lobry JR. SeqinR 1.0-2: A Contributed Package to the R Project for Statistical  
507 Computing Devoted to Biological Sequences Retrieval. In: Bastolla U, Porto M, Roman E,  
508 Vendruscolo M, editors. Structural Approaches to Sequence Evolution: Molecules, Networks,  
509 Populations. New York: Springer Verlag; 2007. pp. 207-32.
- 510 36. States DJ, Boguski MS. Dot Matrix Methods. In: Gribskov M, Devereux J, editors. Sequence  
511 Analysis Primer. New York: Oxford University Press; 1992. pp. 92-124.
- 512 37. Creasy DM, Cottrell JS. Unimod: Protein modifications for mass spectrometry. *Proteomics*.  
513 2004;4(6):1534-6. doi: 10.1002/pmic.200300744.
- 514 38. Diez MJF, Osuga DT, Feeney RE. The sulfhydryls of avian ovalbumins, bovine  $\beta$ -  
515 lactoglobulin, and bovine serum albumin. *Arch Biochem Biophys*. 1964;107(3):449-58. doi:  
516 10.1016/0003-9861(64)90301-7.
- 517 39. Alaiz M, Giron J. Modification of Histidine Residues in Bovine Serum Albumin by Reaction  
518 with (E)-2-Octenal. *J Agric Food Chem*. 1994;42(10):2094-8. doi: 10.1021/jf00046a005.
- 519 40. Brewer CF, Riehm JP. Evidence for possible nonspecific reactions between N-  
520 ethylmaleimide and proteins. *Anal Biochem*. 1967;18(2):248-55. doi: 10.1016/0003-2697(67)90007-  
521 3.
- 522 41. Boyatzis AE, Bringans SD, Piggott MJ, Duong MN, Lipscombe RJ, Arthur PG. Limiting the  
523 Hydrolysis and Oxidation of Maleimide-Peptide Adducts Improves Detection of Protein Thiol  
524 Oxidation. *J Proteome Res*. 2017;16(5):2004-15. doi: 10.1021/acs.jproteome.6b01060.
- 525 42. Fontaine SD, Reid R, Robinson L, Ashley GW, Santi DV. Long-Term Stabilization of  
526 Maleimide-Thiol Conjugates. *Bioconjugate Chem*. 2015;26(1):145-52. doi: 10.1021/bc5005262.
- 527 43. Kuninori T, Nishiyama J. Some Properties of Diastereomers formed in the Reactions of N-  
528 Ethylmaleimide with Biological Thiols. *Agric Biol Chem*. 1985;49(8):2453-4. doi:  
529 10.1080/00021369.1985.10867100.
- 530 44. Garavelli JS. The RESID Database of Protein Modifications as a resource and annotation tool.  
531 *Proteomics*. 2004;4(6):1527-33. doi: 10.1002/pmic.200300777.

532 45. Gentleman RC, Carey VJ, Bates DM, Bolstad B, Dettling M, Dudoit S, et al. Bioconductor:  
533 open software development for computational biology and bioinformatics. *Genome Biol.*  
534 2004;5(10):R80. doi: 10.1186/gb-2004-5-10-r80.

535

536

537

## 538 **Supporting information**

539 **S1 Fig. Numbers of features identified in analyses of NEM-treated and untreated BSA.** After  
540 filtering, the numbers of dependent peptides (DPs) were all similar (no two counts differed by more  
541 than 5%). ‘Conjoined’ DPs were those detected in analysis 1 of NEM-treated BSA and not detected in  
542 analysis 1 of untreated BSA. ‘Constant’ DPs were those detected in both analyses of NEM-treated  
543 BSA. ‘Constantly conjoined’ DPs were those detected in both analyses of NEM-treated BSA and not  
544 detected in either analysis of untreated BSA.

545 **S2 Fig. Localisation plot and mass-shift frequency histogram for dependent peptides of porcine**  
546 **trypsin.** MaxQuant identified three putative deamidations and a putative methylation. Two of the  
547 deamidations were localised to asparagine residues.

548 **S3 Fig. Similarity of mass-shift frequency histograms.** Three pairwise comparisons are shown:  
549 treated/treated, untreated/untreated and treated/untreated. The treated/untreated pair shown is the least  
550 similar of four possible combinations.  $\rho$  = Spearman correlation coefficient.

551 **S4 Fig. Probability localisation plot for putative intact NESyl groups in NEM-treated BSA ( $\Delta m$**   
552  **$\pm$  tolerance = +125.05  $\pm$  0.01 Da).** X-axis values refer to positions in the protein sequence.

553 **S5 Fig. Probability localisation plot for additional intact NESyl groups in NEM-treated BSA**  
554 **( $\Delta m \pm$  tolerance = +68.03  $\pm$  0.01 Da).** X-axis values refer to positions in the protein sequence.

555 **S6 Fig. Probability localisation plot for additional putative hydrolysed NESyl groups in NEM-**  
556 **treated BSA ( $\Delta m \pm \text{tolerance} = +86.04 \pm 0.01 \text{ Da}$ ). X-axis values refer to positions in the protein**  
557 **sequence.**

558 **S7 Fig. Probability localisation plot for putative oxidations in NEM-treated BSA ( $\Delta m \pm$**   
559 **tolerance =  $-25.03 \pm 0.01 \text{ Da}$ ). X-axis values refer to positions in the protein sequence.**

560 **S8 Fig. Localisation plot and mass-shift frequency histogram for dependent peptides of cofilin-1**  
561 **from dimethyl-fumarate-treated rat astrocytes [27]. The nominal  $\Delta m$  of +105 Da is consistent with**  
562 **pyridylethylation of non-cysteine residues.**

563 **S9 Fig. Probability localisation plot for 1-carboxy-2-methylcarboxyethyl groups in cofilin-1**  
564 **from dimethyl-fumarate-treated rat astrocytes [27] ( $\Delta m \pm \text{tolerance} = +24.97 \pm 0.01 \text{ Da}$ ). X-axis**  
565 **values refer to positions in the protein sequence.**

566 **S10 Fig. Localisation plot and mass-shift frequency histogram for dependent peptides of**  
567 **apolipoprotein A-I from alk-3-HHD-treated human blood [28].**

568 **S11 Fig. Localisation plot and mass-shift frequency histogram for dependent peptides of**  
569 **haemoglobin  $\delta$ -chain from alk-3-HHD-treated human blood [28]. All six dependent peptides also**  
570 **mapped to the sequence of haemoglobin  $\beta$ -chain (Fig 3).**

571 **S12 Fig. Localisation plot and mass-shift frequency histogram for dependent peptides of**  
572 **serotransferrin from alk-3-HHD-treated human blood [28].**

573 **S13 Fig. Localisation plot and mass-shift frequency histogram for dependent peptides of serum**  
574 **albumin from alk-3-HHD-treated human blood [28].**

575 **S14 Fig. Localisation plot and mass-shift frequency histogram for dependent peptides of**  
576 **transthyretin from alk-3-HHD-treated human blood [28].**

577 **S15 Fig. Probability localisation plot for putative sulfonated HDMP-type modifications in serum**  
578 **albumin from alk-3-HHD-treated human blood [28] ( $\Delta m \pm \text{tolerance} = +390.10 \pm 0.01 \text{ Da}$ ). X-**  
579 **axis values refer to positions in the protein sequence.**



580 **S1 Script. Script I.** An R script that filters dependent peptides and generates a dependent-peptide  
581 localisation plot and a mass-shift frequency histogram.

582 **S2 Script. Script II.** An R script that filters dependent peptides, enriches them on the basis of  
583 ‘conjunction’ and generates a dependent-peptide localisation plot and a mass-shift frequency  
584 histogram.

585 **S3 Script. Script III.** An R script that filters dependent peptides, enriches them on the basis of  
586 ‘constant conjunction’ and generates a dependent-peptide localisation plot and a mass-shift frequency  
587 histogram.

588 **S4 Script. Script IV.** An R script that isolates dependent peptides with a specified mass shift,  
589 enriches them on the basis of ‘constant conjunction’ and generates a probability localisation plot.

590 **S5 Script. Script V.** An R script that isolates dependent peptides with a specified mass shift and  
591 generates a probability localisation plot.

592 **S1 Table. Dependent peptides and matching variable-modification search results.** Localisation  
593 probabilities for hydrolysed NESyl groups are given in parentheses after the respective amino acid  
594 symbols. Potential sites of modification are underlined, with the most plausible sites in boldface (see  
595 S1 Text). Dependent peptides were allowed to have either of two mass shifts:  $+86.04 \pm 0.01$  Da or  
596  $+143.06 \pm 0.01$  Da (see ‘Results and discussion’).

597 **S2 Table. Dependent peptides with putative alk-3-HHD-derived modifications.** Dependent  
598 peptides containing HDMP-/HTO-type modifications were identified using Script IV. All  
599 modifications were detected as sulfonyl derivatives, and all could be localised to lysine residues. DPs  
600 were matched to Salomón and coauthors’ search results [28] (‘+’ = match) by sequence, site of  
601 modification and modification type (‘Expected analyte’), or by protein site only (‘Expected site’).  
602 Two site numbers are given: the number used for the matching (first number); and the equivalent  
603 number for the mature protein (second number, in parentheses).

604 **S1 Text. Supplementary methods.** Chemicals; Preparation of BSA adducts; Sample preparation for  
605 mass spectrometry; Nano liquid chromatography and mass spectrometry; Enumeration of tryptic  
606 peptides; Calculation of maximum peptide mass; Contaminant databases; Comparison of dependent-  
607 peptide and variable-modification search results; Figure preparation; References; Table A (Gradient  
608 elution timetable); Table B (Criteria used to filter dependent-peptide search results); Table C  
609 (Explanatory notes to accompany scripts).