

1 ***Mimicked synthetic ribosomal protein complex for benchmarking***  
2 ***crosslinking mass spectrometry workflows***

3 **Manuel Matzinger<sup>1,\*§</sup>, Adrian Vasiiu<sup>1,\*</sup>, Mathias Madalinski<sup>1</sup>, Fränze Müller<sup>1</sup>, Florian Stanek<sup>1</sup>,**  
4 **Karl Mechtler<sup>1,2,§</sup>**

5 <sup>1</sup>Institute of Molecular Pathology (IMP), Vienna BioCenter (VBC), Vienna, Austria

6 <sup>2</sup>Institute of Molecular Biotechnology, Austrian Academy of Sciences, Vienna BioCenter (VBC), Vienna, Austria

7 \* These authors contributed equally

8 §Correspondence to: [manuel.matzinger@imp.ac.at](mailto:manuel.matzinger@imp.ac.at), [karl.mechtler@imp.ac.at](mailto:karl.mechtler@imp.ac.at)

9 **ABSTRACT:**

---

10 The field of cross-linking mass spectrometry has matured to a frequently used tool for the investigation  
11 of protein structures as well as interactome studies up to a system wide level. The growing community  
12 generated a broad spectrum of applications, linker types, acquisition strategies and specialized data anal-  
13 ysis tools, which makes it challenging, especially for newcomers, to decide for an appropriate analysis  
14 workflow. Therefore, we here present a large and flexible synthetic peptide library as reliable instrument  
15 to benchmark crosslinkers with different reactive sites as well as acquisition techniques and data analysis  
16 algorithms. Additionally, we provide a tool, IMP-X-FDR, that calculates the real, experimentally vali-  
17 dated, FDR, compares results across search engine platforms and analyses crosslink properties in an au-  
18 tomated manner. The library was used with the reagents DSSO, DSBUS, CDI, ADH, DHSO and azide-a-  
19 DSBSO and data were analysed using the algorithms MeroX, MS Annika, XlinkX, pLink 2, MaxLynx  
20 and xiSearch. We thereby show that the correct algorithm and search setting choice is highly important to  
21 improve ID rate and FDR in combination with software and sample-complexity specific score cut-offs.  
22 When analysing DSSO data with MS Annika, we reach high identification rates of up to ~70 % of the  
23 theoretical maximum (i.e. 700 unique lysine-lysine cross-links) while maintaining a low real FDR of  
24 < 3 % at cross-link level and with high reproducibility, representatively showing that our test system de-  
25 livers valuable and statistically solid results.

26 **Keywords: Crosslinking, mass spectrometry, synthetic peptide library, FDR control**

27 Graphical abstract:

28

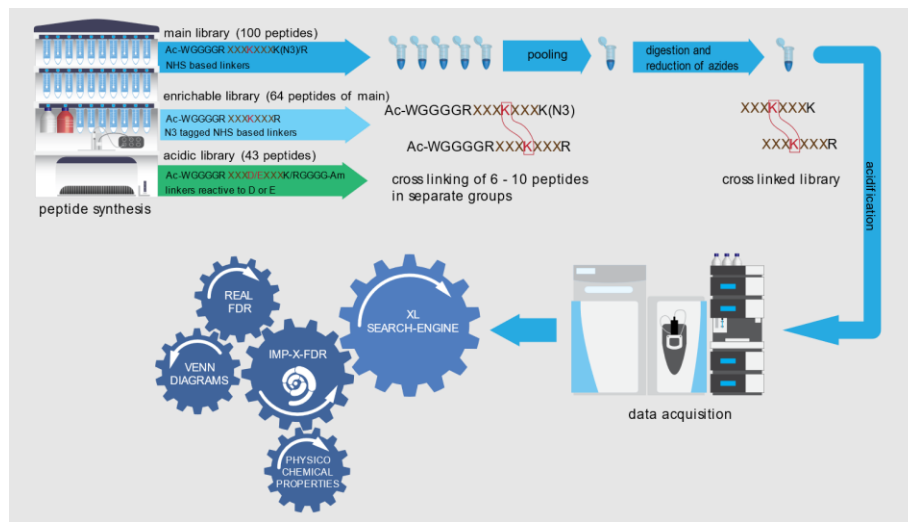
29

30

31

32

33



34

35

## INTRODUCTION

36

The field of cross-linking mass spectrometry has matured and now represents a frequently used technique for the investigation of protein structures as well as to freeze (transient) protein-protein interactions and uncover whole interactomes on a system wide level. Numerous reviews already summarized successful applications but also limitations of this technique.<sup>1-4</sup> The growing community also participated in the generation of a wide variety of cross-linker reagents bearing chemical reactivities mainly towards lysine (e.g. via N-Hydroxysuccinimide esters<sup>5,6</sup>) but also towards acidic amino acids (e.g. by amide formation<sup>7</sup> or hydrazines<sup>8</sup>), cysteine (e.g. via maleimides<sup>9,10</sup>) or even without any specificity (e.g. via diazirine groups<sup>11</sup>). With a focus on proteome wide studies and *in vivo* cross linking, MS-cleavable linkers, like DSSO<sup>12</sup> or DSBU<sup>13</sup>, are facilitating data analysis by generating characteristic doublet ions and became commonly used. Aiming to dig deeper in the interactome of complex samples, reagents bearing an affinity tag for selective enrichment of cross-linked peptides were further developed.<sup>14-16</sup> The optimization of cross-linker specific acquisition strategies<sup>17</sup> and most recently the implementation of ion-mobility<sup>18,19</sup> or FAIMS filtering<sup>20</sup> as additional separation technique further boosted the number of possible crosslink (XL) identifications.

50

The broad spectrum of applications, linker types and acquisition strategies<sup>4</sup> led to the development of lots of specialized data analysis tools<sup>21</sup> which makes it challenging, especially for newcomers, to decide for an appropriate analysis workflow.

53

Therefore, a synthetic peptide library as previously published by our group<sup>22</sup> is a valuable tool for standardization and can be used as a basis to decide for the optimal analysis tool in dependency of the used crosslinker and acquisition strategy. The previous peptide library was based on 95 synthetic peptides of the protein Cas9.

56

57

In this study we present a significantly improved and extended peptide library that now contains a total of 141 peptides from 38 different proteins of the *E. coli* ribosomal complex. This enables finding inter- and intra-protein cross-links in our results. Furthermore, the number of theoretical correct cross-link combinations is increased from 426 in the previously published version to up to 1018 in this library. In conclusion a more reliable and, if supported by the data analysis tool, separate inter/intra false discovery rate (FDR) calculation can be performed. In contrast to our previously published library system of Cas9, the peptides were now combined to 3 different libraries designed to be compatible not only with lysine but also with aspartic- and glutamic-acid reactive cross-linkers as well as for crosslinkers bearing an azide as affinity tag, respectively.

66

With the here reported peptide library, we mimic a real protein complex and a system that is appropriate to find optimal settings for real biological samples as well as to benchmark different crosslinker types and data analysis tools.

68

69 To increase the usability of that library, we additionally created a tool, IMP-X-FDR, that is capable to  
70 check the target-decoy based FDR estimation given by search engines and instead outputs the “real”,  
71 experimentally validated, FDR. Additionally, the tool can correct the number of cross-link IDs to a real  
72 FDR of 1 or 5% by applying a score-cutoff as well as to compare the results obtained from several search  
73 engines or cross-linkers in Venn Diagrams. IMP-X-FDR completes this task in an automated manner and  
74 includes an easy-to-use graphical user interface, which broadens the potential user group. IMP-X-FDR is  
75 free to use and can be downloaded from Github (<https://github.com/fstane/imp-x-fdr>).

76

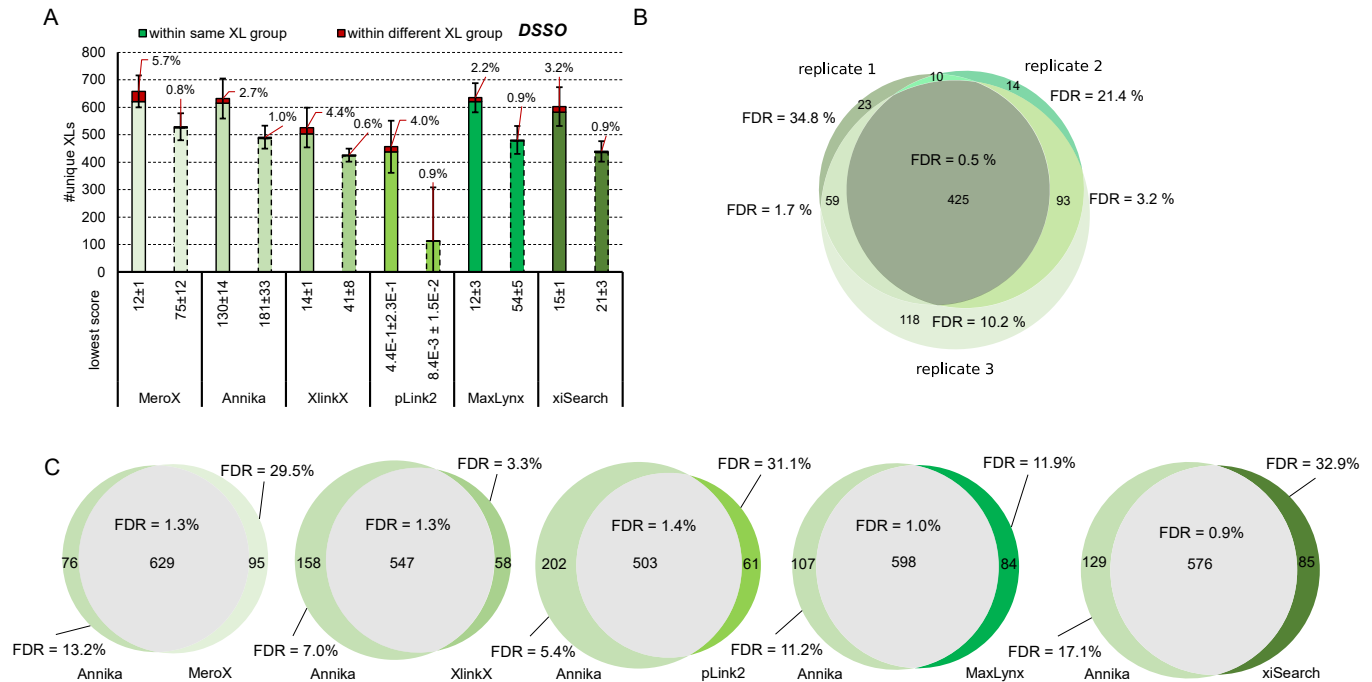
## 77 RESULTS

78 We synthesized 141 peptides based on sequences from 38 proteins of the *E. coli* ribosomal complex (Sup-  
79 plementary Table 1). They are designed to contain exactly one crosslink-able position. Peptides are  
80 grouped to 6 - 10 peptides and crosslinked groupwise. After that, all groups are pooled to obtain the  
81 crosslinked library were links between peptides of different groups or to not synthesized peptides are  
82 known false positives. The main library consists of 100 peptides containing exactly one crosslink-able  
83 lysine residue. All peptides start with the sequence WGGGGR- and their n-termini are protected by an  
84 acetate group to hinder any crosslink reaction at this position. Tryptophan thereby facilitates photometric  
85 quantitation of peptides after synthesis. C-terminal lysine residues are modified to an azide (instead of an  
86 amine) to again block the crosslink reaction. During sample processing the protected n-terminal sequence  
87 part is removed by tryptic digestion and azide modified lysine's are reduced to amines yielding ordinary  
88 tryptic peptides with a known crosslink position for MS/MS analysis. We additionally compiled a library  
89 not containing any azide protected lysine residue but instead exclusively those 64 peptides of the main  
90 library ending with arginine. This “enrichable library” is compatible with azide-based affinity enrichment  
91 as done with the reagent azide-tagged acid-cleavable disuccinimidylbissulfoxide, (DSBSO). Finally, a  
92 third library, made from 43 peptides, is designed to contain exactly one reactive aspartic-acid or glutamic-  
93 acid for use with crosslinker reagents reactive to carboxylic acids. In this “acidic library” the c-terminal  
94 peptide part is amide protected and all sequences end as GGGG after a K or R which will again release  
95 ordinary tryptic peptides after digestion.

### 96 **Benchmarking crosslink search engines with linkers targeting lysine.**

97 To benchmark commonly used crosslink search-algorithms we applied the MS cleavable linker reagents  
98 disuccinimidyl sulfoxide (DSSO), ureido-4,4'-dibutyric acid bis(hydroxysuccinimide) ester (DSBU) and  
99 1,1'-carbonyldiimidazole (CDI) to the main library (Supplemental Table 2). As representatively shown  
100 on the data generated with DSSO the benchmarked search engines all output higher experimentally vali-  
101 dated FDRs than the estimated 1% on crosslink level (Figure 1A). For this dataset MS Annika<sup>23</sup> and  
102 MaxLynx<sup>24</sup> perform best, both by means of correct FDR estimation as well as by means of unique ID  
103 numbers. We additionally applied post-score-cutoffs to correct the experimentally validated FDR to  $\leq 1\%$ .  
104 The obtained results are in line with minimal scores recommended by the software developers (i.e. scores  
105  $>100$  are considered as good for MeroX<sup>25</sup>, 50 is default and 75 seems reliable from our data; 40 is default  
106 for XlinkX<sup>26</sup>, 41 seems reliable from our data). Although using a score-cutoff is an effective strategy to  
107 correct for acceptable FDR, our data also shows, that built in (usually target-decoy based) FDR estima-  
108 tions are not sufficient yet. Especially when using pLink 2<sup>27</sup>, we had the impression that (score-based)  
109 separation of correct and incorrect IDs does not work properly meaning that the majority of crosslink IDs  
110 is lost upon applying our FDR correction. Of note, pLink 2 was initially designed to work with non-MS-  
111 cleavable linker reagents and is not optimized for HCD data in combination with cleavable linkers, which  
112 might explain its weak performance in this dataset compared to all other tested algorithms.

113

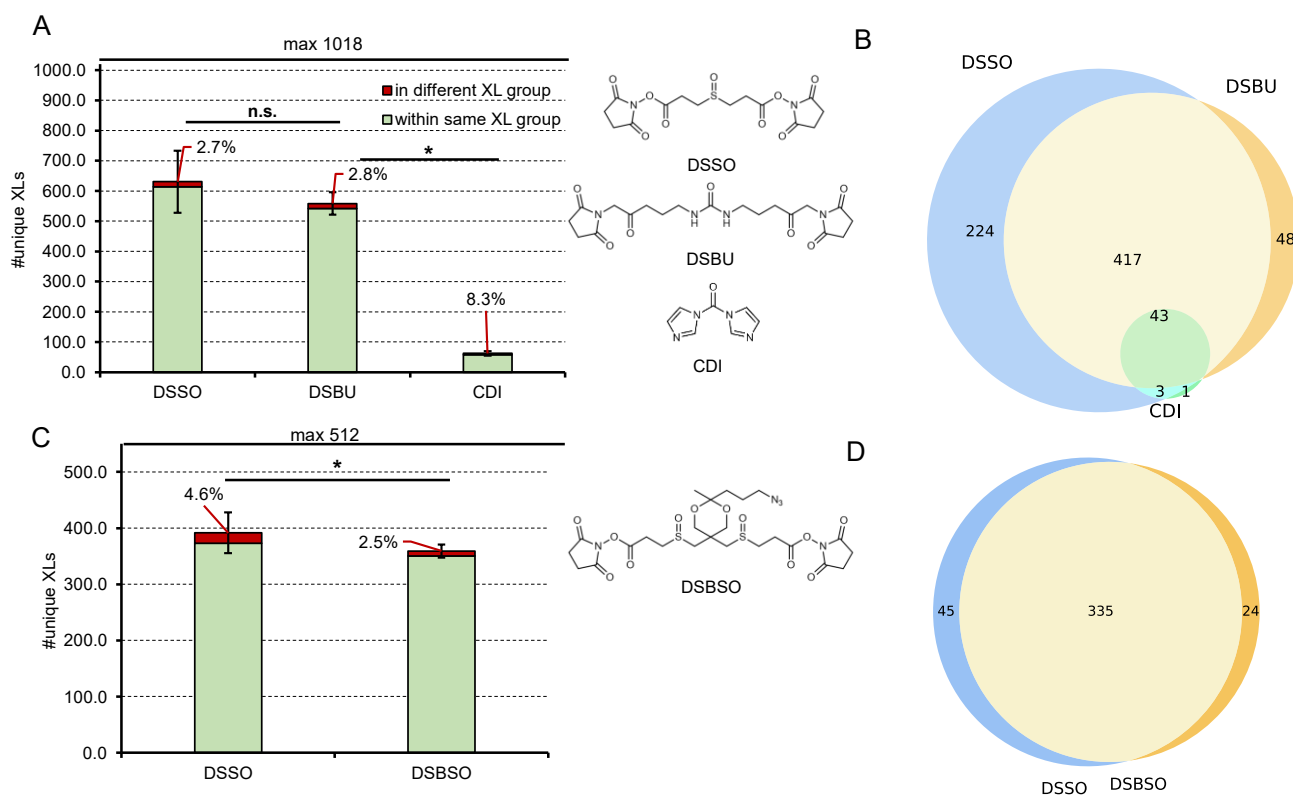


114  
115 **Figure 1: Benchmarking of data analysis tools on the example of DSSO.** (A) Average crosslink numbers using DSSO after  
116 acquisition using a stepped HCD MS2 method. Applied to the main library using the algorithms MeroX<sup>25</sup>, Annika<sup>23</sup>, XlinkX<sup>26</sup>,  
117 pLink 2<sup>27</sup>, MaxLynx<sup>24</sup> or xiSearch<sup>28,29</sup> for analysis. All results were obtained at 1% estimated FDR (solid line bars) and cor-  
118 rected by applying a post-score cutoff to reach an experimentally validated FDR ≤ 1 (dashed line bars). The experimentally  
119 validated FDR is shown as callout, error-bars indicate standard deviations, n=3 (B) Overlap of crosslinks identified in each  
120 replicate using Annika (C) Overlap of cross links identified in replicate 3 after analysis using Annika or an alternative algo-  
121 rithm as given. B & C: Experimentally validated FDRs for commonly found and exclusively found links are given.

122 Instead of using score cutoff values, the comparison of identified crosslinks with more than one search  
123 engine can significantly contribute to improve the confidence in results (Figure 1C). Using our in house-  
124 developed tool IMP-X-FDR we visualized the overlap of search results obtained from MS Annika and a  
125 second search engine and calculated the FDR in an automated manner (Examples of output Figures auto-  
126 matically created by IMP-X-FDR are shown in Supplemental Figure 5 and 6). The fraction of commonly  
127 identified unique crosslinks contains up to 629 entries (Annika + MeroX) and within this fraction the  
128 experimentally validated FDR is ≤ 1.4 % in all cases and therefore very close to the accepted 1%. On the  
129 contrary those crosslinks exclusively identified by only one search engine contain most false positives  
130 yielding to FDR rates of up to 31%. A similar effect is also observed for replicate measurements (Figure  
131 1B). Of 425 unique crosslinks commonly found in three replicates only 2 (0.5 %) were incorrect. While  
132 using crosslinks commonly found across replicate measurements seems to yield highly confident results,  
133 the accumulation of IDs from several replicate measurements to boost link numbers is prone to also ac-  
134 cumulate wrong hits and should therefore be avoided. We further investigated those two crosslinks that  
135 were incorrectly assigned in all three replicates using Annika (Figure 1B): The first one is a homeotypic  
136 link of the peptide MAKLTK that does not exist in the library (but in the database used to search the files).  
137 A peptide with the sequence MAKTIK of the same mass is however part of our library and was therefore  
138 very likely generating the wrongly annotated spectra. The second one connects two existing peptides  
139 (LSYDTEASIAKAK- VAVIKAVR) that are however within different groups.

140 In a next step we benchmarked the reagents DSSO, DSBU and CDI on the main library (Figure 2 A and  
141 B). Expectedly, the performance of DSSO and DSBU is on a similar level, since both have comparable  
142 spacer lengths of 10.1 and 12.5 Å respectively and the same reactive site. The two linkers bear different  
143 reactive groups for MS based fragmentation which might lead to the assumption that differences in spectra  
144 quality explain the slight difference in unique link numbers. Notably this effect is software specific. An-  
145 nika performs very well with DSSO and scores DSSO crosslinks better than DSBU links (average score

146 279 for all DSSO links vs 269 for all DSBU links from our main library). In contrast MeroX performs  
 147 very well with DSBU and scores those links slightly better (average score 131 for all DSSO links vs 133  
 148 for all DSBU links from our main library). In conclusion, when comparing MeroX results, DSBU (767  
 149 links on average) outperforms DSSO (658 links on average) in terms of unique crosslinks (data shown in  
 150 Supplementary Table 2).



151

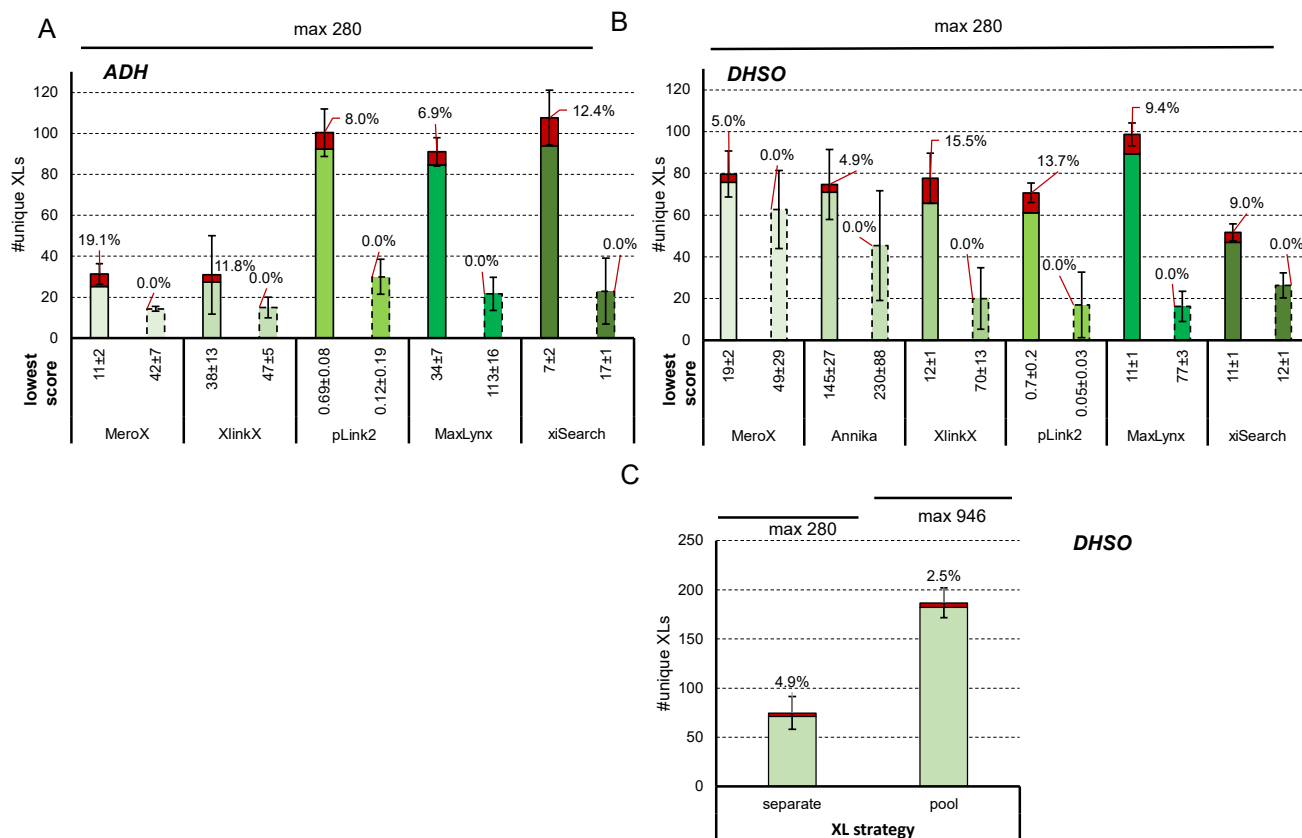
152 **Figure 2: Benchmarking the linker reagents DSSO, DSBU, CDI and DSBSO.** (A&C) Average number of unique crosslink  
 153 IDs after acquisition using a stepped HCD MS2 strategy and maximal theoretical number of true link combinations when  
 154 applying the indicated crosslinkers to the main library (A) or the enrichable library (C) and data analysis using Annika at 1%  
 155 estimated FDR. The experimentally validated FDR is shown as callout, error-bars indicate standard deviations, n=3, unpaired  
 156 Student's t-test, two tailed,  $\alpha = 0.05$ , \*  $P < 0.05$ , n.s. = not significant. (B&D) Overlap of identified links from one representa-  
 157 tive replicate of A (in B) or C (in D) respectively.

158 The zero length crosslinker CDI yielded in  $\leq 80$  unique crosslinks identified with all tested algorithms.  
 159 This low number might be reasoned by no “real” interaction sites within the peptide library that relies on  
 160 crosslink connections formed between freely moving peptides in solution. Therefore, the likelihood of  
 161 two peptides being connected by a crosslinker with a very short spacer is lowered compared to those  
 162 linkers with a 10 – 12 Å spacer. A full list of unique link numbers and experimentally validated FDRs  
 163 using all tested algorithms can be found in Supplemental Table 2.

164 Next, we compared detectability of DSSO vs DSBSO using the enrichable peptide library (Figure 2 C and  
 165 D). In this artificial system any potential steric hindrance of the azide tag of DSBSO can be neglected,  
 166 hence we assume that differences in observed crosslinks are reasoned mainly by the ionizability of the  
 167 resulting connected peptides. As illustrated in the Venn diagram in Figure 2 D, the overlap of identified  
 168 crosslinks is indeed very high and could not be distinguished to an overlap of replicate measurements  
 169 from the same linker (compare to Figure 1 B). Furthermore, DSSO only slightly, though significantly,  
 170 outperforms DSBSO by means of crosslink numbers indicating a slightly increased reactivity or ioniza-  
 171 bility.

172 **Benchmarking crosslink search engines with carboxylic acid reactive linkers**

173 Next, we investigated two reagents targeting acidic amino acids: The non-cleavable adipic acid dihydra-  
 174 zide (ADH) and the cleavable dihydrazide sulfoxide (DHSO) (Figure 3). These linkers were applied to a  
 175 smaller peptide library with a reduced number of only 280 theoretically possible crosslinks formed, how-  
 176 ever, less than 40 % of this number was identified in all cases. This indicates a lowered reaction efficiency  
 177 compared to the more established NHS ester-based linkers, where more than 60 % of the theoretical  
 178 crosslink number was reached (Figures 1 and 2).



179

180 **Figure 3: Benchmarking linker reagents reactive to acidic amino acids.** (A&B) Average number of unique  
 181 link IDs and maximal theoretical number of true crosslink combinations after acquisition using a stepped HCD MS2 strategy  
 182 when using ADH (A) or DHSO (B) to crosslink the acidic library. Data analysis was performed using the indicated algorithm  
 183 at 1% estimated FDR (solid line bars) and corrected by applying a post-score cutoff to reach an experimentally validated  
 184 FDR≤1 (dashed line bars). The experimentally validated FDR is shown as callout, error-bars indicate standard deviations,  
 185 n=3. (C) As B but when crosslinking the library either in separate groups or adding the linker to a pool of all peptides to boost  
 186 resulting ID numbers. Data analysis using Annika at 1 % estimated FDR.

187 For the non-cleavable ADH linker, pLink 2 and MaxLynX seem to perform significantly better, both by  
 188 means of reliability and ID numbers, than their competitors. However, calculated experimentally validated  
 189 FDR values seem extraordinarily high for both reagents and every software tested on the acidic library.  
 190 A proper FDR calculation might be hindered by the relatively small number of crosslinks available in this  
 191 system.

192 Of note, 4-methylmorpholinium chloride (DMTMM), that was used as coupling reagent for ADH and  
 193 DHSO, could form zero-length connections between amines and carboxylic acids. However, only two  
 194 synthetic peptides of the acidic library contain any lysine residue except for those that are terminal after  
 195 tryptic cleavage. We investigated the presence of DMTMM crosslinked (undigested) peptides and found  
 196 no evidence for such a side reaction. The low number of crosslink identifications in the acidic library  
 197 might be reasoned by a slow reaction kinetics and the fact that two steps (activation of carboxylic acids  
 198 by DMTMM followed by nucleophilic attack of the hydrazine group) are required instead of only one as  
 199 is the case for NHS based reagents. To boost the number of crosslink IDs we further tested DHSO on a

200 pool of all 41 peptides of the acidic library. This increases the number to possible crosslink combinations  
201 from 280 to 946 and therefore close to the value of the main library. The number of identified crosslinks  
202 maintained low at ~20% of the theoretical maximum when using Annika (Figure 3 C).

203 To better understand the reaction chemistry of these hydrazine-based linkers, we analyzed the results  
204 obtained for DHSO using our in house developed tool IMP-X-FDR to investigate the distribution of amino  
205 acids in detected crosslink-sequence-matches (CSMs) (Supplementary Figure 4 A and B). We thereby  
206 compared the average frequency of specific amino acids in proximity to the crosslinked aspartic- or glu-  
207 tamic-acid in identified CSMs to the theoretically expected distribution. The theoretical distribution was  
208 calculated from all, in silico generated, crosslinks that can exist within the acidic library (either cross-  
209 linked in separate groups or within one pooled group) under the assumption that every crosslink combi-  
210 nation led to exactly one CSM. By that, missing or predominant combinations can be visualized. For both  
211 DHSO based datasets (pooled and separate, as shown in Figure 3 C) similar dependencies popped up:  
212 Histidine, isoleucine, phenylalanine, tryptophan, tyrosine, and glutamine seem to reproducibly hinder the  
213 formation or identification of a crosslink from the acidic library. The frequency of amino acids within  
214 identified linear peptide sequence matches of the (non-crosslinked) acidic library was compared to the  
215 theoretical amino acid distribution under the assumption of equimolar peptide quantities (Supplemental  
216 Figure 4 C) in an additional experiment. Thereby the same MS method as for crosslink samples was used,  
217 meaning that exclusively ions with a charge  $\geq 3$  were selected for fragmentation. With this we bias the  
218 method to detect longer and higher charged peptides while not recording the majority of linear peptides.  
219 This alters the expected amino-acid distribution as peptides with amino acids carrying a positive charge  
220 are preferentially detected, enabling a fair comparison to the amino acid-distribution seen in our cross-  
221 linked samples. We indeed found fewer peptides containing isoleucine than expected, but clearly more  
222 peptides than expected containing histidine. All other amino acids that seemed to have a negative impact  
223 on crosslink formation were found in relative frequencies as expected. Except for isoleucine this data  
224 strengthens the hypothesis that those amino acids negatively influence the reactivity of DHSO. Especially  
225 the basic histidine might cross-react with the activated carboxylic acid to form an intrapeptidal link, there-  
226 fore impeding the reaction to DHSO.

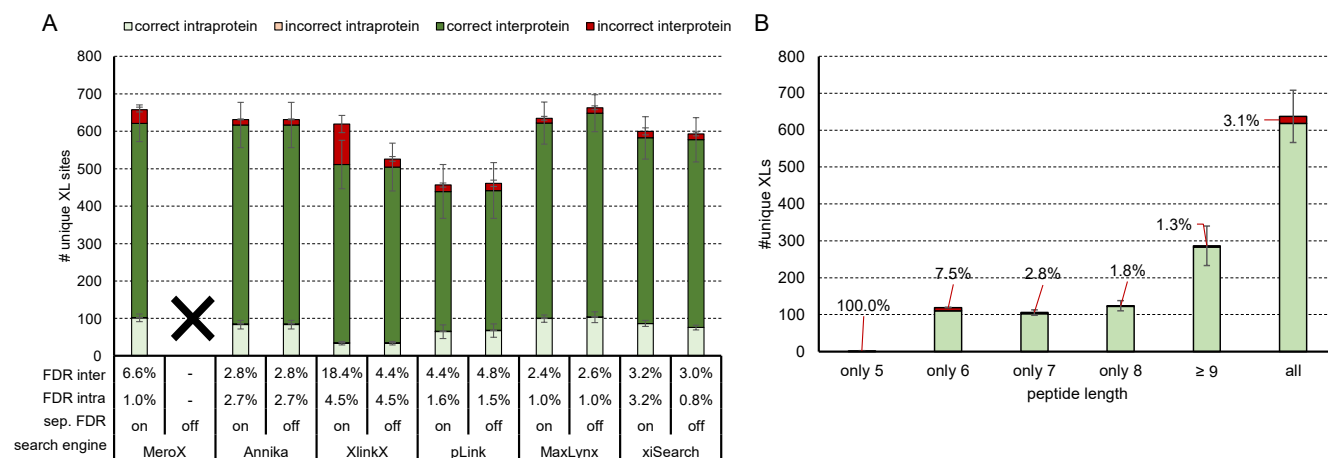
## 227 **Testing the influence of separate FDR calculation and minimal peptide length**

228 Apart from MeroX, all tested search engines allow their users to decide on performing a separate inter-/  
229 intra-crosslink FDR calculation. MeroX calculates FDR of intra- and inter-protein crosslinks as well as  
230 dead-end-links in separate groups by default. A separated target-decoy based FDR calculation is consid-  
231 ered useful as the group of interprotein (heteromeric) connections is much larger compared to the theo-  
232 retical intra-protein crosslinks that can form. This might lead to an underestimated error for the group of  
233 heteromeric crosslinks if the FDR is estimated on the total set of CSMs. Lenz et al. showed that by calcu-  
234 lating the target-decoy based FDR separately, the final FDR of their DSSO dataset was lowered from  
235 36 % to 15 %.<sup>30</sup> This is in line with findings from others that found most wrong identifications in the  
236 group of interprotein connections especially when using large databases<sup>31,32</sup>. They estimated the error rate  
237 to be in the range of 20 – 25 % false positives within a dataset of 2 % overall FDR.<sup>31</sup> In contrast to our  
238 previously published peptide library<sup>22</sup> consisting of peptides from only one protein, the main library of  
239 this study contains 842 theoretical inter-protein crosslinks, 100 intra-protein crosslinks and 100 homeo-  
240 typic crosslinks (link between peptides of the same sequence). The distribution of inter and intra links  
241 nicely represent the theoretical distribution of a real protein mix sample (i.e. *E.coli* ribosome). In conclu-  
242 sion we were wondering if the FDR calculation in separate groups does also influence our results using  
243 the peptide library. In line with our expectations, all tools suffered from a higher error rate within the  
244 group of inter crosslinks (Figure 4 A). Interestingly, xiSearch does not show any difference in inter links  
245 but a lower FDR for intra links when selecting "ignore groups" (=separate FDR set to off) in xiFDR.  
246 Using pLink 2 or MaxLynx the number of crosslinks but also FDR slightly increases when disabling  
247 separate FDR calculation. XlinkX predominantly adds wrong crosslinks to its result file upon enabling  
248 separate FDR calculation. Enabling or disabling this option does however not influence the result when



249 using Annika. In contrast to our expectations, the separate FDR calculation did not significantly improve  
 250 overall FDR or ID numbers independently of the search tool used. This might still be reasoned by the  
 251 nature of our artificial library system that was searched against a database of 171 ribosomal proteins.  
 252 Hence, peptides of 133 proteins contained in the database are not existent in the sample, leading to a  
 253 disproportional large number of theoretical vs existing inter-protein crosslinks. Furthermore, the actual  
 254 number of identified interprotein connections was higher than those of intra-protein links. This corre-  
 255 sponds to the expected theoretical distribution but differs from actual real proteome-wide searches where  
 256 intra-protein links are more abundant. Aiming to further investigate a more complex system we spiked  
 257 the peptide library into a non-crosslinked background of tryptic HEK peptides (1:5 mass ratio) and ana-  
 258 lyzed the resulting data again with or without a separate FDR calculation set in each search algorithm.  
 259 This however led to a very similar result with little to no effect on the final crosslink IDs. Only with  
 260 XlinkX we now identify 348 instead of 259 correct interprotein links while maintaining the error rate  
 261 (Supplementary Figure 1)

262 Next, we tested the influence of the peptide length of the shorter peptide within a linked pair on result  
 263 quality. Figure 4 B clearly illustrates that shorter peptides are more prone for wrongly annotated spectra.  
 264 This fits our expectations as (too) short peptides will generate fewer fragments and therefore yield in less  
 265 confident identifications. In a large database the chance of sequences from different protein overlapping  
 266 by chance is furthermore increased with decreased peptide length potentially leading to ambiguous iden-  
 267 tifications. Based on our data a minimal peptide length of 6 or even 7 seems beneficial, although >100  
 268 unique crosslinks are lost when excluding results containing peptides of 6 amino acids length. Of note,  
 269 our library contains no peptide that has a sequence length of only 5, which is why that group contains  
 270 exclusively wrong hits.



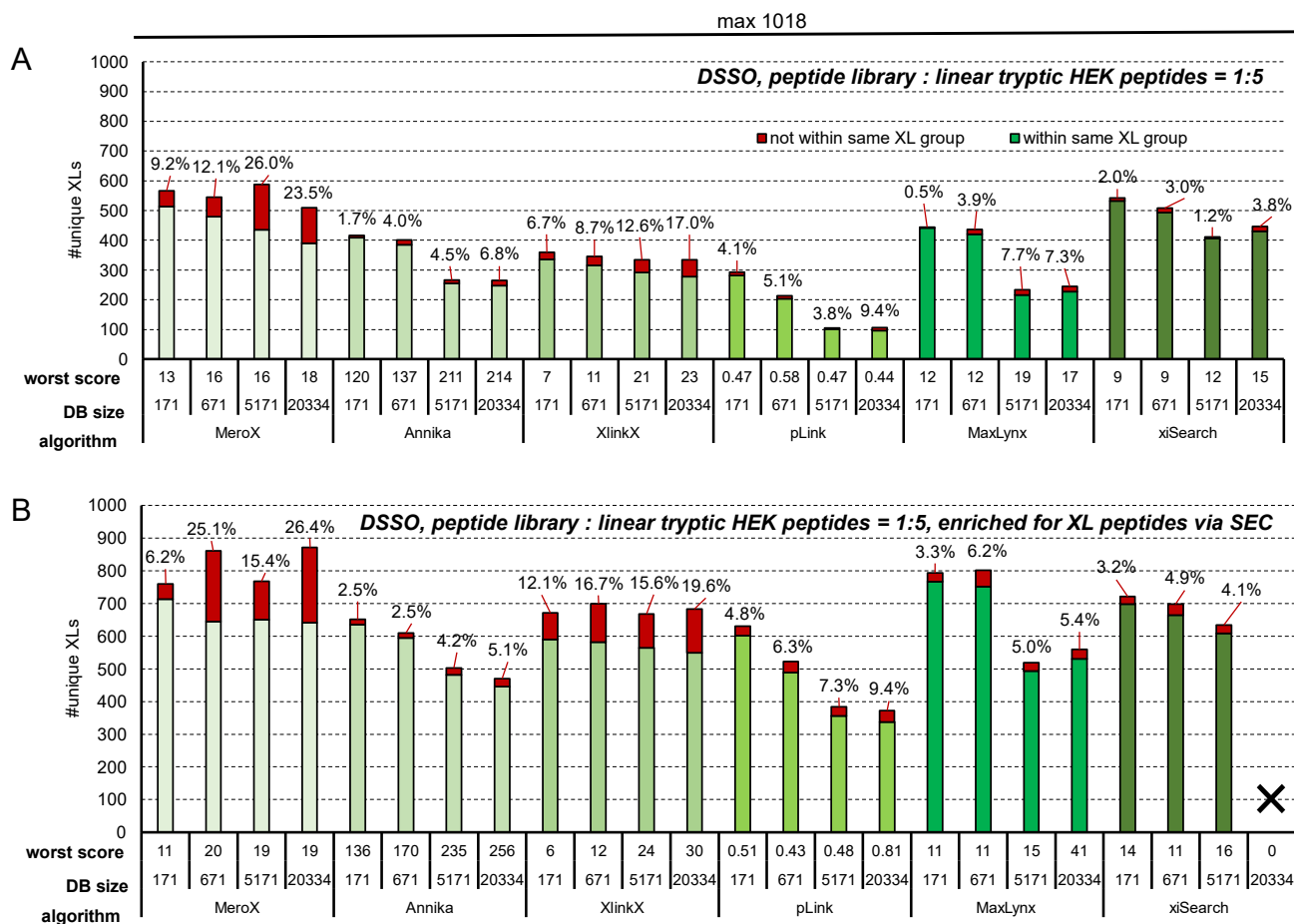
271

272 **Figure 4: Effect of separate inter/intra FDR calculation and minimal peptide length on FDR.** (A) Average number of unique  
 273 crosslinks from the DSSO crosslinked main library after acquisition using a stepped HCD MS2 strategy with separate FDR  
 274 calculation for inter- and intra-crosslinked peptides set on or off. Although synthetic peptides were used for crosslinking their  
 275 sequences are based on ribosomal protein sequences. “Intraprotein” are defined as homomeric links and “interprotein” are  
 276 heteromeric links based on the proteins the synthetic peptides correspond to. Error bars indicate standard deviations, experi-  
 277 mentally validated FDR is shown as callout, n=3 (B) Average number of crosslinks from the DSSO linked main library iden-  
 278 tified with Annika at 1 % estimated FDR when filtering for crosslinked peptides of the given length (meaning the length of the  
 279 shorter peptide within the crosslinked sequence).

## 280 Influence of increased sample complexity and crosslink enrichment.

281 To mimic more realistic conditions – where non-crosslinked linear peptides are way more abundant – we  
 282 spiked the main library into a tryptic digest of linear HEK peptides at a mass ratio of 1:5. The resulting  
 283 mixture was analyzed by means of LC-MS/MS and crosslink searches were performed against databases  
 284 of different sizes, starting with the ribosomal database (171 proteins) that was also used for all other  
 285 searches and ending with proteome wide searches (Figure 5 A). Annika, pLink 2, MaxLynx and xiSearch

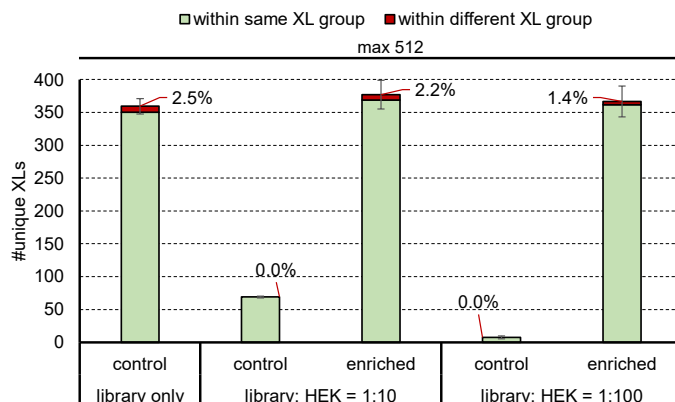
286 maintain FDR at levels below 10 % but loose up to 50 % or more of their identifications upon increasing  
287 the database size to a proteome wide search. Notably, xiSearch maintained high numbers of unique links  
288 for all database sizes at quite low FDR rates. XiSearch therefore clearly outperforms all tested search  
289 engines by means of ID numbers and FDR in case of proteome-wide searches. On the downside it con-  
290 sumes high computational power, which is why analysis of more than one raw-file at once did not work  
291 out for proteome-wide searches on our computer (IntelXeon CPU@ 2.6 GHz, 128 GB RAM). In con-  
292 trast, MeroX and XlinkX maintain their identification numbers at a high level at cost of data reliability,  
293 leading to very high FDR values of up to 26 %. The database size furthermore influences the minimal  
294 score (maximal for pLink 2) accepted as more decoy hits can be found (e.g. Annika increases its accepted  
295 minimal score from 120 to 214). As described by Weisbrod and coworkers<sup>16</sup>, the number of redundant  
296 sequences within the database increases with increasing size leading to ambiguous crosslink IDs. This  
297 cannot be visualized with our peptide library system as a correctly annotated crosslinked peptide is still  
298 correct in case its sequence is contained several times in the database. Within our CSMs at 1% target-  
299 decoy FDR we did not obtain any protein ID for any database size that was ambiguous. We representa-  
300 tively checked on this issue in Annika without any FDR filter and found a maximum of 6.3 % of all CSMs  
301 contained at least one redundant sequence. However, in real samples that need biological interpretation  
302 such redundant sequences impede proper annotation of CSMs to the respective protein-protein interaction.  
303 When enriching the spiked library by size exclusion chromatography (SEC) we were able to (re-)boost  
304 identifications to the level seen without spiking (Figure 5 B) and obtaining the same trends with regards  
305 to FDR. We additionally applied post-score cutoffs to the results using the largest database based on the  
306 scores that yielded in 1 % experimentally validated FDR in our initial non-spiked measurements (shown  
307 in Figure 1A) and that are more stringent than those cutoff values recommended by the authors of the  
308 respective search engines. This improves the experimentally validated FDR, that is however still ranging  
309 from 2.1 % for pLink 2 to 10.5 % for MeroX in the spiked samples (Supplementary Figure 2 A and B).  
310 The database size dependent effects we observed within our spiked samples are furthermore reproducible  
311 when analyzing the non-spiked library with the same set of databases, as representatively analyzed with  
312 Annika and shown in Supplementary Figure 2 C. Our results suggest that the choice of a properly sized  
313 database is of high importance for the reliability of the results as well as that post-score cutoffs to mini-  
314 mize effects of improper FDR estimation need to be empirically determined in dependence of used soft-  
315 ware and complexity of the sample.



316

317 *Figure 5: Performance benchmarking in a mimicked complex environment and upon increased database size. The DSSO*  
 318 *linked main library was mixed with linear tryptic HEK peptides (1:5 w/w). Bars indicate the number of unique crosslinks after*  
 319 *acquisition using a stepped HCD MS2 strategy and identified using the indicated algorithm at 1% estimated FDR when using*  
 320 *databases containing exclusively 171 E. coli ribosomal proteins, or additional 500, 5000 or 20163 human proteins. (A) direct*  
 321 *measurement (B) measurement after enrichment for crosslinked peptides by size exclusion chromatography. Of note, analysis*  
 322 *of the 5 SEC fractions did reproducibly not work with our largest 20334 protein database and xiSearch, as the software*  
 323 *crashes. This data is therefore missing.*

324 To check for the performance of affinity-enrichment using the azide tagged linker DSBSO we also spiked  
 325 the enrichable library, containing no azide-protected lysines, with linear tryptic HEK peptides in mass  
 326 ratios of 1:10 or 1:100 (Figure 6). Enrichment was performed by clicking crosslinked peptides to beads  
 327 functionalized with dibenzocyclooctyne (DBCO) as previously described.<sup>33</sup>



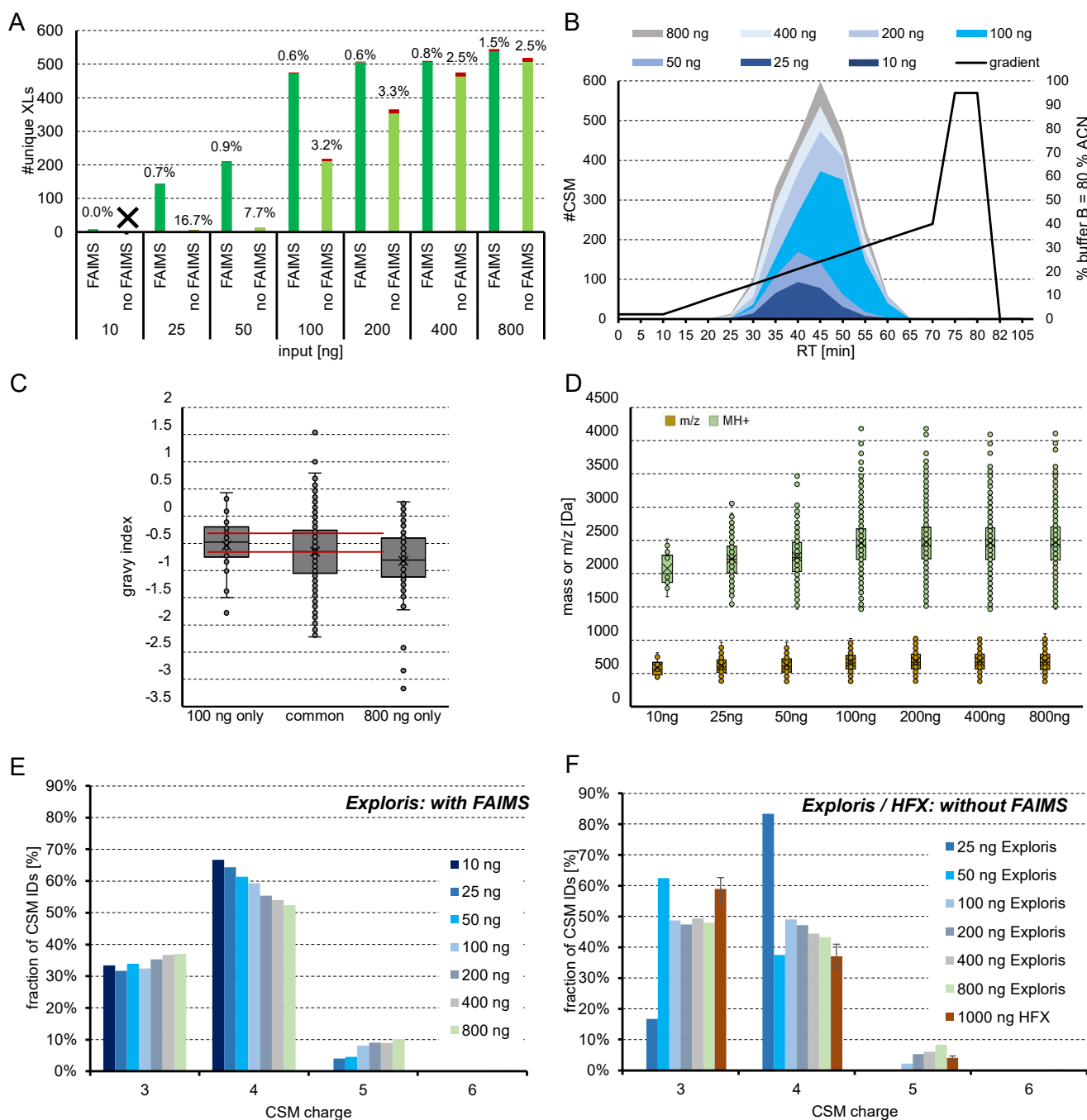
328

329 *Figure 6: Affinity enrichment of DSBSO crosslinked synthetic peptides from a complex environment. Average number of*  
330 *unique crosslinks after acquisition using a stepped HCD MS2 strategy identified in the DSBSO crosslinked enrichable li-*  
331 *brary with or without spiking using linear tryptic HEK peptides as indicated. Controls were directly used for measurement;*  
332 *enriched samples were subjected to affinity enrichment.*

333 The total amount of peptides subjected to MS analysis was kept constant at 1  $\mu$ g for all injections as this  
334 seemed maximal for our LC-MS setting. This means that the 1:100 control sample contains 10 ng cross-  
335 linked peptides. For enrichment 20  $\mu$ g crosslinked peptides were spiked with 200  $\mu$ g or 2 mg HEK pep-  
336 tides resulting in 1.3  $\mu$ g total peptides on average, and independent of the spike ratio, in the enriched  
337 fraction. Although quite some input material was lost during enrichment, the theoretical input can be  
338 freely upscaled to compensate. The resulting enriched samples were of high purity, enabling the injection  
339 of close to 1  $\mu$ g cross-linked material even in samples with high amounts of background (instead of only  
340 10 ng, as in the 1:100 spiked control) and therefore maintaining constantly high crosslink numbers and  
341 low FDR values independently of the sample complexity prior to enrichment.

### 342 Influence of additional FAIMS separation on resulting crosslink identification numbers and prop- 343 erties

344 High-field asymmetric-waveform ion-mobility spectrometry (FAIMS) adds another separation dimension  
345 and therefore decreases spectrum complexity and reduces noise. Both effects were reported beneficial for  
346 the identification of crosslinked peptides.<sup>20</sup> We probed the effect of FAIMS on the here presented syn-  
347 thetic peptide system using DSSO (Supplementary Figure 3). In line with the results of Schnirch and  
348 coworkers<sup>20</sup>, we observe a maximum number of unique crosslinks when using compensation voltages  
349 (CV) in the range of -50 to -60V. Furthermore, we observed very high reproducibility when using FAIMS  
350 and a trend to lower FDR values upon lowered CV. The combination of 3 CVs within one run boosted  
351 our overall identification number to 700 which is a 10 % increase compared to our measurements without  
352 FAIMS (numbers from analysis using Annika at 1 % estimated FDR).



353

354 **Figure 7: Variation of input amount and physicochemical properties of crosslinks.** Lowered input amounts, as indicated, of  
 355 the DSSO cross-linked main library were measured on an Orbitrap Exploris 480 using a stepped HCD MS2 strategy. Data  
 356 was analyzed using Annika at 1 % FDR. (A) Unique crosslinks with or without FAIMS attached and experimentally validated  
 357 FDR is given above bars. (B) Distribution of spectral matches over retention time and used gradient with FAIMS. (C) Distri-  
 358 bution of hydrophobicity of CSMs identified exclusively in the 100 or 800 ng sample or in both samples (measured with  
 359 FAIMS). (D) Distribution of identified cross-linked peptide masses of CSMs (M+H) and m/z in dependence of used input  
 360 amount. (measured with FAIMS) (E) Charge distribution obtained with FAIMS using CVs -50, -55 and -60 for acquisition and  
 361 stepwise changing the input amount from 10 – 800 ng or (F) without FAIMS when acquiring on HFX or Exploris instrument.

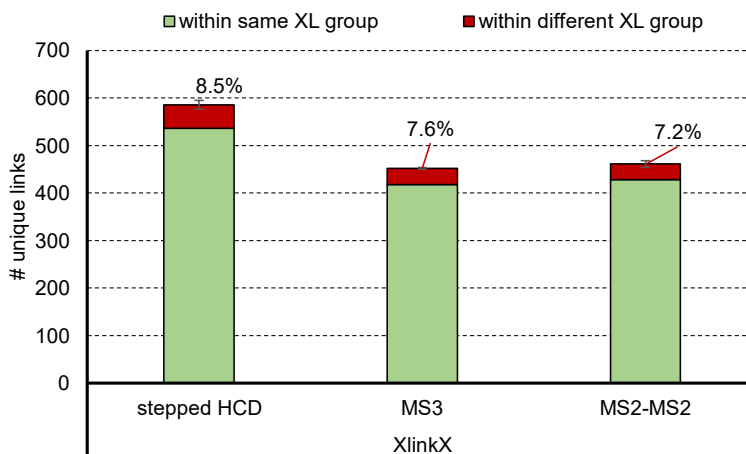
362 Next, we investigated if lower injection amounts are sufficient for identification of crosslinks thanks to  
 363 the improved signal to noise ratios when using FAIMS.<sup>34,35</sup> Without FAIMS, on the HFX 1000 ng were  
 364 injected for all samples, which was stepwise lowered down to 10 ng on the Exploris with and without  
 365 FAIMS. Judged from ID numbers, 100 ng, only 1/8 of the maximal amount used, was sufficient to still  
 366 identify close to 500 unique crosslinks with FAIMS which is ~90% of the identified links using 800 ng  
 367 input (Figure 7 A). Further lowered amounts lead to a drastic decrease of crosslinks. In a direct comparison

368 of data acquired with/without the FAIMS device attached, we see a clear advantage of FAIMS especially  
369 for lowered injection amounts based on the number of identified crosslinks. Upon injection of higher  
370 peptide input the effect seems to diminish but the number of CSMs found per unique crosslink is still  
371 increased by using FAIMS (Supplemental Figure 8A).

372 Since our data with FAIMS outperforms those without we focused further investigations on FAIMS data:  
373 When comparing retention times of crosslinked peptides IDs, we observed a slight shift towards higher  
374 retention times with higher peptide amounts (Figure 7 B). This effect is not reasoned by an overloaded  
375 column as the retention time of individual CSMs did not change (Supplemental Figure 8B), but rather by  
376 identifying additional CSMs. When looking into the physicochemical properties of the additionally iden-  
377 tified CSMs, a shift in hydrophobicity (Figure 7 C) as well as an increase of molecular weights and  $m/z$   
378 (Figure 7 D) can be observed, which might explain shifted retention times upon increasing input amounts.  
379 Furthermore, the relative charge distribution (Figure 7 E) depends on the input amount. The fraction of  
380 high charged  $z = 5$  ions is increased, while the fraction of  $z = 4$  charged ions is lowered with higher input.  
381 This observation is in line with the seen dependency of molecular weight to input amount as larger pep-  
382 tides are more likely highly charged in an acidic environment. When comparing this charge distribution  
383 to those observed without FAIMS (Figure 7 F) a clear shift from more dominant  $z = 3$  charged ions  
384 without FAIMS to predominant  $z = 4$  charged ions with FAIMS can be observed, which is beneficial for  
385 the detection of predominantly higher charged crosslinked ions over linear peptides. Of note this effect  
386 seems to be dependent on the used instrument type as well, since the relative fraction of +3 charged CSMs  
387 is biggest in our data from the HFX. Those results without FAIMS obtained from 25 – 50 ng input contain  
388 only a total of 24 and 6 CSMs respectively, which is why the relative results seem not to fit to those results  
389 with higher ID numbers. With only 10 ng of input no CSMs were identified at 1% FDR (Supplemental  
390 Figure 8A).

### 391 **Benchmarking acquisition strategies.**

392 We finally investigated the FDR of crosslinks using different MSn acquisition strategies on an Eclipse  
393 Tribrid mass spectrometer (Thermo Fisher Scientific). Data analysis was performed using XlinkX at 1 %  
394 FDR on CSM and residue pair level. MS3<sup>36</sup> methods are reported as more reliable for crosslink identi-  
395 fication when using cleavable crosslinker. Therefore, MS3 (acquired as described in<sup>36</sup>) was compared to an  
396 MS2-MS2 method (acquired as described in<sup>22</sup>) and our standard stepped HCD acquisition method for the  
397 main library crosslinked using DSSO (Figure 8). In line with previous results from our group<sup>17,22</sup>, stepped  
398 HCD outperformed MS3 and MS2-MS2 results in terms of unique crosslink numbers. We however ob-  
399 served a slight increase in FDR (8.5% vs 7.6 and 7.2% for MS3 and MS2-MS2 respectively). Surprisingly,  
400 the experimentally validated FDR for MS3 is higher with lower unique residue pair numbers than for the  
401 MS2-MS2 method which contrasts with the literature. Although MS3 and MS2-MS2 methods are thought  
402 to give advantage for crosslink identification, stepped HCD performed better in our hands.



403

404 *Figure 8: Comparison of experimentally validated FDR of MSn methods. DSSO cross-linked main library were measured*  
405 *on an Eclipse Tribrid Mass Spectrometer using a stepped HCD MS2, MS3 and CID MS2- ETD MS2 acquisition strategy. Data*  
406 *was analyzed using XlinkX at 1 % FDR. Red bars indicate false crosslinks whereas green bars represent true unique crosslinks.*  
407 *Error bars indicate standard deviations, experimentally validated FDR is shown as callout, n=3.*

408

## 409 DISCUSSION

410 The here presented peptide libraries represent a new and highly flexible standard to benchmark crosslinker  
411 reagents of different chemical reactivity as well as for comparison of search engines or acquisition strat-  
412 egies. Thanks to peptide sequences originating from 38 different proteins, the library represents a realistic  
413 digest from the *E. coli* ribosomal protein complex, allowing for in-depth analysis of search-engine specific  
414 FDR calculation. Our in house developed tool IMP-X-FDR comes with an easy-to-use user interface and  
415 allows FDR calculation, comparison of crosslink results across software platforms and investigation of  
416 crosslink-properties in an automated manner also for non-bioinformaticians.

417 Our results suggests that additional and empirical score cutoffs are a valuable instrument to correct the  
418 actual FDR. The height of this cutoff value is not only software specific but also dependent on database  
419 size and sample complexity. Our library helps to find such specific score-cut-off values but also showed  
420 that built-in target-decoy based FDR estimation overestimates correct results in case of all tested search  
421 engines. Of note, not all tested tools allow to perform their target-decoy based FDR estimation on unique  
422 crosslink or even protein-protein interaction level. To give an example, MeroX and pLink 2 do their FDR  
423 estimation on spectrum level, leading to false hits propagating during grouping to the final unique cross-  
424 links and higher than expected final experimentally validated FDR<sup>28</sup>. We believe, machine learning ap-  
425 proaches as well as the inclusion of additional parameters like retention time or ion mobility to (re-)score  
426 identifications might improve crosslink numbers and validity of FDR in future algorithms.

427 The most challenging part of data analysis seems to find a good compromise in the trade between high  
428 identification numbers and low FDR. In our data, MS Annika and MaxLynx seem to find this optimum  
429 best for cleavable crosslinkers, while pLink 2 seems to perform very well for the tested non-cleavable  
430 linker. The performance of xiSearch was stable for non-cleavable and cleavable crosslinkers as well as  
431 for database sizes up to 5000 proteins and is therefore an allrounder within the crosslink search engines.  
432 We however faced troubles in software stability for proteome-wide searches. We further observed that  
433 some search engines fit better to specific linker reagents than others leading to an additional linker de-  
434 pendent performance difference that is not caused by the crosslinker chemistry itself but by technical  
435 reasons as an altered spectrum complexity for MS cleavable vs non-cleavable reagents. This is even re-  
436 flected in issues to properly define specific new linker reagents. To give an example, sulfoxy based linkers  
437 as DSSO or DSBSO were reported to generate characteristic doublet peaks of two different delta masses  
438 upon MS fragmentation thanks to water elimination<sup>12,37</sup>. Upon the tested algorithms, MaxLynx and pLink

439 2 do not allow a definition of more than two fragments and their results might be further improved upon  
440 implementation. Another example is Thermos Proteome Discoverer 2.5, where the definition of a fragment  
441 mass 0 is impossible but needed for the zero-length linker CDI. As a workaround a very low mass  $\geq 1E^{-5}$   
442 must be defined. This affects search engines as XlinkX or Annika when running as a node within Proteome  
443 Discoverer. In line with these observations, it seems that the developers of search engines focus on  
444 specific linker types for optimization of their algorithms and this yields in boosted results and better score-  
445 based separation of target vs decoy hits for linkers of the exact same chemistry.

446 Many studies aim to minimize error rate and maximize confidence in crosslink results with alternative  
447 approaches: The Rappsilber laboratory investigated this issue by separately crosslinking fractions after  
448 size exclusion chromatography (SEC) and accepting only those protein-protein interactions as confident  
449 that are between proteins of the same SEC fraction.<sup>30</sup> Another common way to validate crosslink data is  
450 by comparison to 3D structures of representative protein complexes in the dataset. Yugandhar et al.<sup>38</sup>  
451 showed that this approach can lead to a significant underestimation of the actual error rate by implement-  
452 ing additional quality criteria, including the validation of interactions by orthogonal techniques, by known  
453 interactions or by adding the proteome of unrelated organisms to the search and checking for misidentifi-  
454 cations. In line with our results their results further show, that applying minimal score-cutoffs can drast-  
455 ically reduce error rate and might therefore be highly beneficial to obtain interpretable and confident re-  
456 sults. The Bruce lab estimates the error rate for large scale studies by determining the theoretical maximal  
457 number of inter- and intra-protein crosslinks based on available 3D structures. They demonstrate that  
458 those inter-protein crosslink fractions greater than the theoretical maximum value are most likely occur-  
459 ring from false positive IDs.<sup>39</sup>

460 Complementing these studies, a synthetic library system serves as ground truth model to experimentally  
461 validate observed FDR. We believe that a gold standard in the field of crosslinking MS must be estab-  
462 lished in the future for robust data analysis. Further software updates or novel algorithms will improve  
463 the reliability of the results and increase the coverage of crosslink identification.

464 Our data will therefore provide valuable input to benchmark new or updated search engines. The freely  
465 available IMP-X-FDR can be easily adopted for automated FDR calculation with any novel crosslink  
466 search engine thanks to the open-source code. Furthermore, improvements in crosslinker reagents, MS  
467 instrumentation or chromatography can be validated using the physical library where the exact number of  
468 theoretically reachable crosslinks is well defined.

## 469 **METHODS**

### 470 Peptide synthesis

471 Solid phase peptide synthesis was done using Fmoc chemistry on a SYRO with Tip Synthesis Module  
472 (MultiSynTech GmbH). Each coupling step was performed as double coupling using HATU/DIEA for  
473 carboxylic acid activation. Lysine residues at the C-terminus bore an azide group instead of an amine to  
474 hamper any cross-linking at this position. N-termini were designed as acetyl protected WGGGGR se-  
475 quence tag and C-termini were designed as amide protected RGGGG sequence tag (for peptides to be  
476 used with linkers reactive to acids, see Supplementary Table 1 for all sequences). For this Fmoc-L-  
477 Arg(Pbf)-TCP (# PC-01-0126), Fmoc-Rink-Amide-(aminomethyl) (#PC-01-0501) or Fmoc-L-Lys(N3)-  
478 TCP (custom synthesized) resins were used respectively (all: INTAVIS Peptide Services GmbH & Co.  
479 KG). Purification was performed using a C18 kinetex column (5  $\mu$ m) and a 30 min gradient. All peptides  
480 were analyzed using a 4800 MALDI TOF/TOF (Applied Biosystems) for quality control purposes. Ly-  
481 ophilized peptides were solubilized in water and their concentration was estimated by measuring their  
482 absorption via a nanodrop (DeNovix DS-11 FX+) at 280 nm and calculating the sequence specific extinc-  
483 tion coefficient using the ProtParam tool<sup>40</sup>. Peptide solutions were dried under reduced pressure, resolu-  
484 bilized in 50 mM HEPES pH 7.5 at a concentration of 5 mM and mixed to groups for cross-linking (Sup-  
485plementary Table 1).



## 486 Sample preparation

487 For lysine reactive cross-linker reagents (DSSO, DSBSO, DSBU, CDI) 9.3 mM cross-linker reagent stock  
488 solutions were freshly prepared in dry DMSO. 0.5  $\mu$ L of stock solution was added to 1  $\mu$ L of each peptide  
489 group in separate vials. Additional stock solution was added 4x every 30 min adding up to a total of 2.5  $\mu$ L  
490 cross-linker stock solution. The resulting 3.5  $\mu$ L reaction mix were quenched using 31.5  $\mu$ L 100 mM  
491 ammoniumbicarbonate (ABC) buffer for 30 min and pooled to a single tube. The resulting mix was di-  
492 gested by addition of 5 ng trypsin/group over night at 37 °C. Azide protection groups were finally reduced  
493 to the respective amines by incubation to 50 mM (final concentration) tris(2-carboxyethyl)phosphine  
494 (TECEP) for 30 min at room temperature. Reduced peptides were pooled to a single vial, aliquoted and  
495 stored at -70 °C upon further usage.

496 For aspartic acid and glutamic acid reactive cross-linker reagents (DHSO, ADH) 300 mM cross-linker  
497 reagent and 1.2 M (4-(4,6-dimethoxy-1,3,5-triazin-2-yl)-4-methyl-morpholinium chloride) (DMTMM)  
498 stock solutions were prepared in 25 mM HEPES pH 7.5. 0.25  $\mu$ L of cross-linker and DMTMM stock-  
499 solution were added 5x every 30 min to 1  $\mu$ L of each peptide group. The reaction was quenched by adding  
500 trifluoroacetic acid (TFA) to a final concentration of 4 % (w/v) for 20 min followed by re-neutralization  
501 by addition of 50  $\mu$ L 1M Tris pH 7.5 buffer. Peptides were pooled and digested as described above.

## 502 Enrichment strategies

503 To mimic complex mixtures, cross-linked and digested peptide pools were mixed with a 5-100x excess  
504 (by mass) of tryptic HEK peptides. The resulting spiked samples were enriched either by size exclusion  
505 chromatography (SEC) or via affinity enrichment.

506 For SEC, ~10  $\mu$ g of cross-linked peptide-library + tryptic HEK peptides were fractionated on a TSKgel  
507 SuperSW2000 column (300 mm  $\times$  4.5 mm  $\times$  4  $\mu$ m, Tosoh Bioscience), which was operated at 200  $\mu$ L/min  
508 in 30 % ACN, 0.1 % TFA. Fractions were collected every minute, ACN was removed under reduced  
509 pressure to obtain a concentrated sample for LC-MS/MS

510 DSBSO cross-linked peptides (+ linear HEK peptides in varying mass ratios) were affinity enriched using  
511 dibenzylcyclooctyne (DBCO) immobilized to beads as described elsewhere<sup>33</sup>.

512 Tryptic HEK peptides were generated as follows: HEK cells were lysed in 10 M urea in 100 mM Tris by  
513 ultrasonication. The cleared lysate was reduced at a final concentration of 10 mM dithiothreitol in the  
514 presence of benzonase for 1 h at 37 °C. This was followed by alkylation at a final concentration of 20 mM  
515 iodoacetamide for 30 min at room temperature in the dark. Digestion was performed using LysC (1:200  
516 w/w) for 2 h at 37 °C in 6 M urea followed by addition of trypsin (1:200 w/w) for 16 h t 37 °C in 2.5 M  
517 urea.

## 518 Chromatographic separation and mass spectrometry

519 Samples were separated using a Dionex UltiMate 3000 HPLC RSLC nano-system coupled to an Q Exac-  
520 tive™ HF-X Orbitrap mass spectrometer or to an Orbitrap Exploris™ 480 mass spectrometer equipped  
521 with a FAIMS pro interface (all: Thermo Fisher Scientific). Samples were loaded onto a trap column  
522 (Thermo Fisher Scientific, PepMap C18, 5 mm  $\times$  300  $\mu$ m ID, 5  $\mu$ m particles, 100 Å pore size) at a flow  
523 rate of 25  $\mu$ L min<sup>-1</sup> using 0.1 % TFA as mobile phase. After 10 min, the trap column was switched in line  
524 with the analytical column (Thermo Fisher Scientific, PepMap C18, 500 mm  $\times$  75  $\mu$ m ID, 2  $\mu$ m, 100 Å).  
525 Peptides were eluted using a flow rate of 230 nl min<sup>-1</sup>, with the following gradient: 0 -10 min 2 % buffer  
526 B, followed by an increasing concentration of buffer B up to 40 % until min 130. This is followed by a 5  
527 min gradient from reaching 95 % B, washing for 5 min with 95% B, followed by re-equilibration of the  
528 column in buffer A at 30°C (buffer B: 80 % ACN, 19.92 % H<sub>2</sub>O and 0.08 % TFA, buffer A: 99.9% H<sub>2</sub>O,  
529 0.1% TFA).

530 The mass spectrometer was operated in a data-dependent mode, using a full scan (m/z range 375-1500,  
531 nominal resolution of 120.000, target value 1E6). MS/MS spectra were acquired by stepped HCD using  
532 an NCE (normalized collision energy) of 27±6 for sulfoxy group linkers (DSSO, DSBSO, DHSO), 30±3  
533 for urea-based linkers (DSBU, CDI) and 28±4 for non-cleavable linkers (ADH). An isolation width of  
534 1.0 m/z, a resolution of 30.000 and a target value of 5E4 (on HF-X) and 1E5 (on Exploris) was set. Pre-  
535 cursor ions selected for fragmentation (± 10 ppm, including exclusively charge states 3-8) were put on a  
536 dynamic exclusion list for 30 s. Measurements including FAIMS were performed on the Orbitrap Exploris  
537 under alteration of used compensation voltages as given for each result.

538 MS3 and MS2-MS2 acquisitions were performed on a Orbitrap Eclipse Tribid (Thermo) using the same  
539 HPLC setting as described above. Acquisition strategies were designed as described in Wheat *et al.*<sup>36</sup> and  
540 Beveridge *et al.*<sup>22</sup> respectively.

#### 541 Data Analysis and post processing

542 Data analysis was performed against a custom shotgun database containing 171 *E. coli* ribosomal proteins  
543 at 1 % FDR level. For analyses using Annika or XlinkX, Thermo raw files were loaded to Thermo Pro-  
544 teome Discoverer 2.5 that and both search engines were used as node within that software. MaxLynx was  
545 used as part of MaxQuant v 2.0.2.0 by direct usage of Thermo raw files as well. For MeroX, raw files  
546 were converted to mzML and for pLink 2 and xiSearch files were converted to mgf using MSConvertGUI  
547 v3.0.21084. The result files are available for download in the PRIDE repository<sup>41</sup> using the identifier  
548 PXD029252 (User: reviewer\_pxd029252@ebi.ac.uk; Password: sihLJE67). The software specific set-  
549 tings are furthermore summarized in Supplementary Table 3.

550 Post processing was done using the graphical user interface of our in house developed tool IMP-X-FDR  
551 (Supplemental Figure 5 E). It enables to calculate the experimentally validated FDR and therefore validate  
552 the target-decoy based FDR estimated by the crosslink search engine according to the following formulae:

$$553 \quad FDR_{\text{experimentally validated}} = \frac{\#target \text{ XLs across peptides not within same XL group}}{\#target \text{ XLs total}}$$

554 When calculation FDR on CSM level, unique residue pairs (XLs) are replaced by CSM IDs in the above  
555 formulae. Some search engines allow the export of target-decoy filtered XL lists, but not all of them. To  
556 ensure functionality with all search engines and enable the direct usage of the search engine result file as  
557 input for IMP-X-FDR, our tool automatically filters away IDs marked as decoy and exclusively selects  
558 inter- and intra-protein crosslinks (but excludes dead-end links or linear peptides).

559 FDR validation is done based on crosslinks only allowed as correct in case they are formed within the  
560 same crosslink group (see Supplementary Table 1 for allocation of peptides to groups). We call this func-  
561 tionality “FDR recalculation” and adopted the code for each crosslink search engine, due to differences  
562 in their output format. For a correct FDR recalculation, a support file containing all group-allocated pep-  
563 tides of all used (sub) peptide libraries is provided with the software. The tool outputs a csv file containing  
564 a list of all XLs within the same or different group as well as informative graphs showing the number of  
565 IDs and the score vs experimentally validated FDR or number of crosslinks (Supplemental Figure 5 A-  
566 C). The functionality “Venn diagrams” of IMP-X-FDR was used to visualize the overlap of replicates of  
567 searches from different search algorithms (example output shown in Supplemental Figure 5 D). This func-  
568 tionality uses the output of “FDR recalculation” as input, which ensures a uniform format and compares  
569 peptide sequences, their originating protein, and the position of the peptide in that protein.

570 The third function of IMP-X-FDR is to investigate physicochemical properties of crosslinks. To do so the  
571 freely available tools from Biopython<sup>42</sup>, specifically from the Bio package, Bio.SeqUtils subpackage and  
572 Bio.SeqUtils.ProtParam module, were used. Crosslinked peptides were represented in a linearized form  
573 to ensure compatibility with the used packages originally designed for linear peptides. IMP-X-FDR out-  
574 puts a csv file containing calculated crosslink properties, which includes the isoelectrical point, fraction  
575 of aromatic amino acids, molecular mass, gravy value and amino acid distribution. The obtained data is  
576 automatically compared to the respective properties of all (in silico generated) theoretically formed cross-  
577 links within the library. Thereby we assume the identification of exactly one CSM for each theoretical  
578 crosslink. The unnormalized output graphics are constructed on the crosslink level and histograms con-  
579 structed on CSM level are normalized to a total area of 1. Finally IMP-X-FDR investigates amino acid  
580 motives using the module seqlogo 5.29.8<sup>43</sup> to create position probability matrices. Thereby the closest  
581 three neighboring amino acids of the linker's binding site are investigated for frequent or rare amino acids  
582 and can be compared to the (theoretically expected) crosslinks within the library. Representative output  
583 graphs are illustrated in Supplementary Figure 6. A user's manual, containing a detailed explanation of  
584 each output file and used functions is delivered with IMP-X-FDR. The code is freely available  
585 (<https://github.com/fstaneke/imp-x-fdr>) and can be used on command line basis or via a graphical user  
586 interface.

## 587 **ASSOCIATED CONTENT**

### 588 **Supporting Information**

- 589 Supplementary Table 1: List of all synthesized peptides and their annotation to groups for crosslinking
- 590 Supplementary Table 2: List of all crosslink IDs at 1% estimated FDR from crosslinked library samples  
591 measured without FAIMS and experimentally validated FDRs from main, enrichable and acidic library.
- 592 Supplementary Table 3: Search settings used from MeroX, Annika, XlinkX, pLink 2, MaxLynx and  
593 xiSearch analyses.
- 594 Figure S 1: Effect of separate inter/intra FDR calculation in a complex environment.
- 595 Figure S 2: Performance benchmarking in a mimicked complex environment searching against the full  
596 proteome.
- 597 Figure S 3: DSSO linked peptides on a FAIMS equipped device
- 598 Figure S 4: Influence of specific amino acids in proximity to the crosslink-site influencing the formation  
599 of a crosslink
- 600 Figure S 5: Exemplary output figures of IMP-X-FDR.
- 601 Figure S 6: Exemplary output figures of the physicochemical cross-link properties functionality of IMP-  
602 X-FDR
- 603 Figure S 7: Effect of site localization searches (KSTY) for FDR estimation.
- 604 Figure S8: Impact of FAIMS and input amount to identified CSMs

## 605 AUTHOR INFORMATION

### 606 Author Contributions

607 The study was designed by MM and KM. Experiments were performed by ADV. IMP-X-FDR was created  
608 by ADV and wrapped into a user interface by FS. M. Madalinski performed peptide synthesis. Experi-  
609 ments were performed and the manuscript was written by MM. FM added data and figures for results of  
610 xiSearch and MS3, MS2-MS2 acquisition strategies and helped with the revised version of the manuscript.  
611 All authors have given approval to the final version of the manuscript.

612

## 613 ACKNOWLEDGMENT

614 This work supported by the EPIC-XS, Project Number 823839, funded by the Horizon 2020 Program of  
615 the European Union, by the project LS20-079 of the Vienna Science and Technology Fund and the by the  
616 ERA-CAPS I 3686 and P35045-B project of the Austrian Science Fund. We thank the IMP for general  
617 funding and access to infrastructure and especially the technicians of the protein chemistry facility for  
618 continuous laboratory support. Our gratitude further goes to Dr. Elisabeth Roitinger for fruitful discus-  
619 sions and for her valuable inputs to the manuscript.

## 620 REFERENCES

- 621 1. Iacobucci, C., Götze, M. & Sinz, A. Cross-linking/mass spectrometry to get a closer view on protein  
622 interaction networks. *Curr. Opin. Biotechnol.* **63**, 48–53 (2020).
- 623 2. Piersimoni, L. & Sinz, A. Cross-linking/mass spectrometry at the crossroads. *Anal. Bioanal. Chem.*  
624 (2020) doi:10.1007/s00216-020-02700-x.
- 625 3. Belsom, A. & Rappsilber, J. Anatomy of a crosslinker. *Curr. Opin. Chem. Biol.* **60**, 39–46 (2021).
- 626 4. Matzinger, M. & Mechtler, K. Cleavable Cross-Linkers and Mass Spectrometry for the Ultimate Task  
627 of Profiling Protein-Protein Interaction Networks in Vivo. *J. Proteome Res.* **20**, 78–93 (2021).
- 628 5. Pilch, P. F. & Czech, M. P. Interaction of cross-linking agents with the insulin effector system of  
629 isolated fat cells. Covalent linkage of 125I-insulin to a plasma membrane receptor protein of 140,000  
630 daltons. *J. Biol. Chem.* **254**, 3375–3381 (1979).
- 631 6. Staros, J. V. N-hydroxysulfosuccinimide active esters: bis(N-hydroxysulfosuccinimide) esters of two  
632 dicarboxylic acids are hydrophilic, membrane-impermeant, protein cross-linkers. *Biochemistry* **21**,  
633 3950–3955 (1982).
- 634 7. D’Este, M., Eglin, D. & Alini, M. A systematic analysis of DMTMM vs EDC/NHS for ligation of  
635 amines to hyaluronan in water. *Carbohydr. Polym.* **108**, 239–246 (2014).
- 636 8. Leitner, A. *et al.* Chemical cross-linking/mass spectrometry targeting acidic residues in proteins and  
637 protein complexes. *Proc. Natl. Acad. Sci.* **111**, 9455–9460 (2014).
- 638 9. Smyth, D., Blumenfeld, O. & Konigsberg, W. Reactions of N-ethylmaleimide with peptides and  
639 amino acids. *Biochem. J.* **91**, 589–595 (1964).
- 640 10. Partis, M. D., Griffiths, D. G., Roberts, G. C. & Beechey, R. B. Cross-linking of protein by  $\omega$ -malei-  
641 mido alkanoylN-hydroxysuccinimido esters. *J. Protein Chem.* **2**, 263–277 (1983).
- 642 11. Gomes, A. F. & Gozzo, F. C. Chemical cross-linking with a diazirine photoactivatable cross-linker  
643 investigated by MALDI- and ESI-MS/MS. *J. Mass Spectrom.* **45**, 892–899 (2010).
- 644 12. Kao, A. *et al.* Development of a novel cross-linking strategy for fast and accurate identification of  
645 cross-linked peptides of protein complexes. *Mol. Cell. Proteomics MCP* **10**, M110.002212 (2011).

- 646 13. Müller, M. Q., Dreiocker, F., Ihling, C. H., Schäfer, M. & Sinz, A. Cleavable cross-linker for protein  
647 structure analysis: reliable identification of cross-linking products by tandem MS. *Anal. Chem.* **82**,  
648 6958–6968 (2010).
- 649 14. Burke, A. M. *et al.* Synthesis of two new enrichable and MS-cleavable cross-linkers to define protein–  
650 protein interactions by mass spectrometry. *Org Biomol Chem* **13**, 5030–5037 (2015).
- 651 15. Steigenberger, B., Pieters, R. J., Heck, A. J. R. & Scheltema, R. A. PhoX: An IMAC-Enrichable  
652 Cross-Linking Reagent. *ACS Cent. Sci.* **5**, 1514–1522 (2019).
- 653 16. Chavez, J. D., Weisbrod, C. R., Zheng, C., Eng, J. K. & Bruce, J. E. Protein Interactions, Post-trans-  
654 lational Modifications and Topologies in Human Cells. *Mol. Cell. Proteomics* **12**, 1451–1467 (2013).
- 655 17. Stieger, C. E., Doppler, P. & Mechtler, K. Optimized Fragmentation Improves the Identification of  
656 Peptides Cross-Linked by MS-Cleavable Reagents. *J. Proteome Res.* **18**, 1363–1370 (2019).
- 657 18. Ihling, C. H., Piersimoni, L., Kipping, M. & Sinz, A. Cross-linking/Mass Spectrometry Combined  
658 with Ion Mobility on a timsTOF Pro Instrument for Structural Proteomics. (2021)  
659 doi:10.1101/2021.03.26.437136.
- 660 19. Steigenberger, B. *et al.* Benefits of Collisional Cross Section Assisted Precursor Selection (caps-  
661 PASEF) for Cross-linking Mass Spectrometry. *Mol. Cell. Proteomics MCP* (2020)  
662 doi:10.1074/mcp.RA120.002094.
- 663 20. Schnirch, L. *et al.* Expanding the depth and sensitivity of cross-link identification by differential ion  
664 mobility using FAIMS. *Anal. Chem.* (2020) doi:10.1021/acs.analchem.0c01273.
- 665 21. Yılmaz, Ş. *et al.* Cross-linked peptide identification: A computational forest of algorithms. *Mass Spec-*  
666 *trom. Rev.* **37**, 738–749 (2018).
- 667 22. Beveridge, R., Stadlmann, J., Penninger, J. M. & Mechtler, K. A synthetic peptide library for bench-  
668 marking crosslinking-mass spectrometry search engines for proteins and protein complexes. *Nat.*  
669 *Commun.* **11**, 742 (2020).
- 670 23. Pirklbauer, G. J. *et al.* MS Annika: A New Cross-Linking Search Engine. *J. Proteome Res.* (2021)  
671 doi:10.1021/acs.jproteome.0c01000.
- 672 24. Yılmaz, Ş., Busch, F., Nagaraj, N. & Cox, J. *Accurate and automated high-coverage identification of*  
673 *chemically cross-linked peptides with MaxLynx.* 2021.08.26.457759 [https://www.biorxiv.org/con-](https://www.biorxiv.org/content/10.1101/2021.08.26.457759v1)  
674 [tent/10.1101/2021.08.26.457759v1](https://www.biorxiv.org/content/10.1101/2021.08.26.457759v1) (2021) doi:10.1101/2021.08.26.457759.
- 675 25. Iacobucci, C. *et al.* A cross-linking/mass spectrometry workflow based on MS-cleavable cross-linkers  
676 and the MeroX software for studying protein structures and protein–protein interactions. *Nat. Protoc.*  
677 **13**, 2864 (2018).
- 678 26. Liu, F., Lössl, P., Scheltema, R., Viner, R. & Heck, A. J. R. Optimized fragmentation schemes and  
679 data analysis strategies for proteome-wide cross-link identification. *Nat. Commun.* **8**, 15473 (2017).
- 680 27. Chen, Z.-L. *et al.* A high-speed search engine pLink 2 with systematic evaluation for proteome-scale  
681 identification of cross-linked peptides. *Nat. Commun.* **10**, 3404 (2019).
- 682 28. Fischer, L. & Rappsilber, J. Quirks of Error Estimation in Cross-Linking/Mass Spectrometry. *Anal.*  
683 *Chem.* **89**, 3829–3833 (2017).
- 684 29. Mendes, M. L. *et al.* An integrated workflow for crosslinking mass spectrometry. *Mol. Syst. Biol.* **15**,  
685 e8994 (2019).
- 686 30. Lenz, S. *et al.* Reliable identification of protein-protein interactions by crosslinking mass spectrome-  
687 try. *Nat. Commun.* **12**, 3564 (2021).

- 688 31. de Jong, L., Roseboom, W. & Kramer, G. A composite filter for low FDR of protein-protein interac-  
689 tions detected by in vivo cross-linking. *J. Proteomics* **230**, 103987 (2021).
- 690 32. de Jong, L. *et al.* In-Culture Cross-Linking of Bacterial Cells Reveals Large-Scale Dynamic Protein-  
691 Protein Interactions at the Peptide Level. *J. Proteome Res.* **16**, 2457–2471 (2017).
- 692 33. Matzinger, M., Kandioller, W., Doppler, P., Heiss, E. H. & Mechtler, K. Fast and Highly Efficient  
693 Affinity Enrichment of Azide-A-DSBSO Cross-Linked Peptides. *J. Proteome Res.* **19**, 2071–2079  
694 (2020).
- 695 34. Swearingen, K. E. & Moritz, R. L. High-field asymmetric waveform ion mobility spectrometry for  
696 mass spectrometry-based proteomics. *Expert Rev. Proteomics* **9**, 505–517 (2012).
- 697 35. Gerbasi, R. *et al.* Deeper Protein Identification by Using FAIMS in Top-down Proteomics. (2021)  
698 doi:10.26434/chemrxiv.13653578.v1.
- 699 36. Wheat, A. *et al.* Protein interaction landscapes revealed by advanced in vivo cross-linking–mass spec-  
700 trometry. *Proc. Natl. Acad. Sci.* **118**, (2021).
- 701 37. Kaake, R. M. *et al.* A New *in Vivo* Cross-linking Mass Spectrometry Platform to Define Protein-  
702 Protein Interactions in Living Cells. *Mol. Cell. Proteomics* **13**, 3533–3543 (2014).
- 703 38. Yugandhar, K., Wang, T.-Y., Wierbowski, S. D., Shayhidin, E. E. & Yu, H. Structure-based valida-  
704 tion can drastically underestimate error rate in proteome-wide cross-linking mass spectrometry stud-  
705 ies. *Nat. Methods* **17**, 985–988 (2020).
- 706 39. Keller, A., Chavez, J. D., Felt, K. C. & Bruce, J. E. Prediction of an Upper Limit for the Fraction of  
707 Interprotein Cross-Links in Large-Scale In Vivo Cross-Linking Studies. *J. Proteome Res.* **18**, 3077–  
708 3085 (2019).
- 709 40. Gasteiger, E. *et al.* Protein Identification and Analysis Tools on the ExPASy Server. in *The Prote-*  
710 *omics Protocols Handbook* (ed. Walker, J. M.) 571–607 (Humana Press, 2005). doi:10.1385/1-59259-  
711 890-0:571.
- 712 41. Perez-Riverol, Y. *et al.* The PRIDE database and related tools and resources in 2019: improving sup-  
713 port for quantification data. *Nucleic Acids Res.* **47**, D442–D450 (2019).
- 714 42. Cock, P. J. A. *et al.* Biopython: freely available Python tools for computational molecular biology  
715 and bioinformatics. *Bioinformatics* **25**, 1422–1423 (2009).
- 716 43. Crooks, G. E., Hon, G., Chandonia, J.-M. & Brenner, S. E. WebLogo: a sequence logo generator.  
717 *Genome Res.* **14**, 1188–1190 (2004).
- 718