# Structure-based validation can drastically under-estimate error rate in proteome-wide cross-linking mass spectrometry studies

**Kumar Yugandhar**[1,2], **Ting-Yi Wang**[1,2], **Shayne D. Wierbowski**[1,2], **Elnur Elyar Shayhidin**[1,2], **Haiyuan Yu**[1,2,*]

[1]Department of Computational Biology, Cornell University, Ithaca, NY, USA

[2]Weill Institute for Cell and Molecular Biology, Cornell University, Ithaca, NY, USA

## Abstract

Thorough quality assessment of novel interactions identified by proteome-wide cross-linking mass spectrometry (XL-MS) studies is critical. Almost all current XL-MS studies have validated cross-links against known 3D structures of representative protein complexes. Here we provide theoretical and experimental evidence demonstrating this approach can drastically underestimate error rates for proteome-wide XL-MS datasets, and propose a comprehensive set of four data-quality metrics to address this issue.

## Editorial summary:

The current standard approach for estimating error in proteome-scale crosslinking-mass spectrometry datasets has severe limitations. A proposed set of data-quality metrics provides a more accurate assessment of error rate.

---

Cross-linking mass spectrometry (XL-MS) is a powerful platform capable of unveiling protein interactions and capturing their structural dynamics[1]. The wealth of information from proteome-wide XL-MS approaches facilitate large-sale identification of protein-protein interactions[2, 3], and high-throughput three dimensional structural modeling of functional protein complexes[4-6]. With the increased throughput of these techniques, the number of false positive cross-links and incorrect interactions can quickly add up with just one large-scale XL-MS experiment, if one is not careful. Therefore, thorough quality assessment has become critically important.

It has been shown that the conventional false discovery rate (FDR) calculations for XL-MS can be susceptible to error propagation[7] (Supplementary Note 1). Currently, almost all proteome-wide XL-MS studies leverage available 3D structures of representative complexes

---

*To whom correspondence should be addressed. Tel: 607-255-0259; Fax: 607-255-5961; haiyuan.yu@cornell.edu.

Competing Financial Interests

The authors declare no competing financial interests.

for validation and quality assessment[8, 9]. Here, we demonstrate fundamental flaws in this structure-based quality assessment approach that can drastically underestimate the error rates of large-scale XL-MS datasets.

In small-scale XL-MS studies, the fraction of cross-linked residue pairs that satisfy the maximum distance a given cross-linker can span (e.g. 30Å for DSSO[10]) provides meaningful insights into protein flexibility and the quality of the cross-links detected. In proteome-wide XL-MS studies, researchers extend this concept and use representative, highly abundant complexes such as the ribosome and proteasome to estimate the quality of all cross-links reported. However, true positive and false positive cross-links in these large-scale studies are not equally likely to successfully map onto an existing 3D structure, leading to massive underestimation of false positives (Fig.1a).

To theoretically demonstrate this, let us consider a reference protein complex structure consisting of 100 subunits. Because a false positive cross-link can be detected between any two random proteins within the proteome (~20,000 proteins for human proteome-wide experiments), for a given false positive with one of its ends mapped to the reference complex, the probability that the second end also maps to this complex by random chance is $5 \times 10^{-3}$ (100/20000). It should be noted that this probability would be even lower for the often used ribosome (76 subunits: PDB ID 5T2C) or proteasome (34 subunits: PDB ID 5GJQ) complexes. However, these probabilities only hold for random peptide pairs (derived from false positive cross-links); true positive cross-links are much more likely to perfectly map to existing 3D structures. Conceptually, this is very similar to the fact that false positive cross-links are much more likely to be interprotein than intraprotein as shown by previous studies[11, 12]. We expect that almost all false positive cross-links will have only one peptide mapped to the reference complex structure. The current structural-mapping approach explicitly considers only cross-links where both peptides map to the same complex structure, and in doing so, it enriches for true positive cross-links and massively underestimates the error rates for proteome-wide XL-MS datasets. Consequentially, this validation approach may erroneously annotate artifacts as novel interactions, resulting in less reliable experimental datasets for further studies.

To demonstrate our theory experimentally, we obtained a subset of 122 raw files from our recent proteome-wide human K562 XL-MS study[2]. Next, in order to generate three sets of cross-links with drastically different qualities, we ran the XlinkX search engine (Proteome Discoverer 2.2) using three criteria of increasing stringency ('10% FDR', '1% FDR' and '1% FDR with  XlinkX score  50'; Methods). As shown in Fig.1b, at 10% FDR, a set of 35,561 interprotein cross-links were identified (we intentionally chose 10% FDR to obtain a low-quality set of cross-links with many false positives), 1% FDR yielded 16,591 interprotein cross-links, whereas '1% FDR with  XlinkX score  50' yielded 985 interprotein cross-links. We mapped the interprotein cross-link residue pairs from these three sets separately onto the 3D structure of the human proteasome following the conventional methodology. We then calculated the percentage of mapped residue pairs that satisfied DSSO's theoretical constraint (  30Å). We observed that there was no significant difference (all *P*>0.85) among the three sets in terms of their percentage of residue pairs satisfying the distance constraint (Fig.1c), even though the overall qualities of these three sets are

drastically different by design. Additionally, we utilized our recently published search engine, MaXLinker[2], to repeat the analysis and observed similar results, confirming that these findings are software-independent (Extended Data Fig.1a,1b). We further re-analyzed raw files from two other publicly available studies representing different organisms (*E. coli*[13] and Mouse[10]) and cellular compartments (mitochondria[10]) (Fig.2a,b and Extended Data Fig.2a,b). These experimental results confirm that the current structure-mapping approach fails to capture the underlying error rate and indicate an urgent need for reliable metrics to estimate the quality of proteome-wide cross-linking datasets.

To address the pitfalls of the current validation approach, we propose the following comprehensive set of four measurements:

*(i) Fraction of structure-corroborating identifications (FSI)*: The current structure-based validation approach considers only those cross-links where both peptides mapped to the reference structure. Here, we propose FSI as an improved structure-based metric that uses the number of ***all*** interprotein cross-links with at least one peptide mapped to the reference structure, not just those with both peptides mapped, as the denominator (Methods).

*(ii) Fraction of mis-identifications (FMI)*: Including the proteome of an unrelated organism in the search database as an internal negative control can be an efficient way to independently assess the underlying error rate of the cross-link search algorithm[2, 14, 15] (Methods).

*(iii) Fraction of interprotein cross-links from known interactions (FKI)*: Using prior knowledge of experimentally-detected protein interactions to calculate the FKI provides a comparative quality estimate (Methods).

*(iv) Fraction of validated novel interactions using orthogonal experimental assays*: It is essential to validate a representative set of novel interactions identified in proteome-wide XL-MS studies using an orthogonal experimental assay (e.g., Y2H, PCA), to ensure data quality and reproducibility (Methods). Furthermore, using a Bayesian framework[16, 17] (Supplementary Note 2) and leveraging the validation rates among a positive reference set (PRS) of well-known interactions and a negative reference set (RRS) of random pairs, we can calculate the absolute precision of the novel interactions detected in a XL-MS study.

We next applied our proposed metrics on our human proteome-wide XL-MS results, and demonstrated how each of them efficiently captures the differences in data quality among the three filtered sets (Fig.1d-g). Fig.1d shows that our improved structure-based metric FSI, differentiates the three sets with statistical significance, which could not be achieved by the conventional structure-based approach (Fig.1c). The results are consistent with our prior theoretical expectation that applying more stringent quality filters would remove predominantly (likely false positive) cross-links with only one peptide mapped to the structure, and thereby result in higher FSI values (Fig.1b,d). Furthermore, Fig.1e reveals the exact same trend: FMI is significantly lower for the '1% FDR with XlinkX score 50' set compared to the other two sets. Moreover, as shown in Fig. 1f, FKI exhibits great agreement with the expected data quality of different data sets (at '1% FDR with XlinkX score 50', FKI is 55.5%; but at '10% FDR', FKI is merely 4.4%; $P \ll 1 \times 10^{-20}$).

Finally, we performed a thorough orthogonal experimental validation of randomly chosen novel interactions from the three sets using a protein complementation assay (PCA)[18, 19]. The fraction of PCA-positive novel interactions from '1% FDR with XlinkX Score 50' set (the highest quality set) is distinctively higher compared to the other two sets and indistinguishable from that of PRS (*P*=0.17; Fig.1g). Notably, the fraction of PCA-positive interactions for '1% FDR' and '10% FDR' are indistinguishable from that of RRS. Furthermore, using the Bayesian framework[16, 17] (Supplementary Note 2), we calculated the absolute precision of the novel interactions detected in our human XL-MS study (Extended Data Fig.3). Especially since the true FDR at the protein pair level can be significantly higher than the estimated FDR at the peptide pair level[7, 15], absolute precision will be critically important for confirming the quality of novel protein-protein interactions identified in a large-scale cross-linking study. Finally, we confirmed the usefulness and robustness of the three computational metrics (namely FSI, FMI and FKI) on the re-analyzed *E. coli* (Fig.2c-e) and Mouse mitochondrial (Extended Data Fig.2c-e) XL-MS datasets, and using the additional search engine MaXLinker[2] (Extended Data Fig.1c-e).

Taken together, our four metrics constitute a comprehensive framework to facilitate both relative comparison across different datasets and absolute estimation of error rates. Moreover, because these metrics stem from different principles, they provide complementary insights to various aspects of the data quality. FMI provides an orthogonal estimation of FDR and serves as an absolute measure of error rate. In fact, other methods[11, 14, 20, 21] have been reported to provide complementary error estimates for XL-MS studies, and show good agreement with FMI in terms of the relative data quality across different datasets (Supplementary Note 3). Since FSI typically leverages thoroughly studied complexes, in theory, it should provide an absolute estimate of quality. Nonetheless, we do note that it may only provide relative comparison especially in cases where limited or incomplete 3D reference structures are available (Fig.2c and Extended Data Fig.2c). FKI and 'fraction of validated novel interactions using orthogonal experimental assays' specifically addresses the quality of detected interactions inferred from interprotein cross-links. Because a large fraction of true protein interactions is yet to be discovered, FKI only provides relative estimates of quality among comparable datasets. Finally, even if high-throughput orthogonal assays are not available, we recommend that low-throughput validation assays (such as co-immunoprecipitation[22]) be performed on a meaningful subset of the interactions identified (Supplementary Note 4).

In conclusion, we theoretically and experimentally illustrated the limitation of the current structure-based validation approach for evaluating proteome-wide XL-MS results. Furthermore, we proposed a comprehensive set of four metrics, and demonstrated their ability to distinguish datasets with varying qualities. Moreover, we acknowledge that this drastic under-estimation of the error rate by the conventional structure-based method is unlikely to pose a serious issue for XL-MS studies focused on specific proteins and individual complexes as long as the cross-link search is performed against only proteins that are included in the experiment. Importantly, this issue is highly relevant for the increasingly popular proteome-wide XL-MS experiments[8, 9] and cross-linking immunoprecipitation MS (xIP-MS) studies[23]. Going forward, a comprehensive and accurate quality assessment

framework such as the one proposed in this work needs to be adapted to aid in the advancement of XL-MS technologies.

# Methods

## Data processing

Cross-links were identified using XlinkX software (Proteome discoverer 2.2). PD templates for different XlinkX search methodologies were obtained from Rosa Viner (Thermo fisher Scientific). The raw files for *E. coli* XL-MS dataset (MS2-MS3) were obtained through e-mail request to Dr. Fan Liu. In addition to filtering XLs at '10% FDR' and '1% FDR', we further filtered the '1% FDR' set using ' XlinkX score' cut off 50. XlinkX score is a CSM-level scoring parameter in XlinkX software that indicates confidence in identifying a peptide pair over the next best competing peptide pair for a given precursor mass (higher score implies better quality). In addition to the three sets, we also filtered cross-links at '20% FDR' and carried out the structure-based mapping analyses to verify that the trend observed in Fig. 1c, Fig. 2b and Extended Data Fig. 2b holds at a this much higher FDR threshold (Extended Data Fig.4). During generation of 20% FDR set using MaXLinker software, the FDR was estimated at the CSM level.

Target protein sequences were downloaded from Uniprot database[24] (with filter 'reviewed'): [(i) *Escherichia coli*: 5268 sequences; downloaded on 28th October 2017, (ii) *Saccharomyces cerevisiae:* 7904 sequences; downloaded on 28th September 2017 ('reviewed: yes'), (iii) Human (*Homo sapiens*): 42202 sequences (20206 canonical; 21996 isoforms); downloaded on 23rd June 2017), and (iii) Mouse (*Mus musculus*): 17019 sequences; downloaded on 8th July 2019]. More specifically, the human database consists of 21,996 isoform sequences in addition to the 20,206 canonical sequences. The mouse database consists of the canonical sequences for 17019 proteins. The *E. coli* database contains of 5268 sequences in total, consisting of 4436 sequences from K12 strain (4436; most common) and the remaining 832 sequences from other less common strains. Similarly, for *Saccharomyces cerevisiae* the fasta database consists of 6721 sequences from the common strain 'ATCC 204508' and the remaining sequences comes from other lesser common strains such as 'YJM789', 'RM11-1a' and 'JAY291'. We utilized the full list of protein entries (did not rely the protein grouping) to classify each cross-link as 'interprotein' or 'intraprotein', to avoid any inconsistencies that might occur due to potential protein grouping artifacts. When performing searches for Fig.1e and Extended Data Fig.5a, XlinkX crashed multiple times given the huge number of raw files (122 files) and the enormous search space (*H. sapiens* + *S. cerevisiae* ). Hence, we ran the searches on a smaller set of raw files (25 files) to generate Fig.1e and Extended Data Fig.5a.

## Mapping of cross-links to existing PDB structures

Cross-links from our human K562 proteome-wide XL-MS dataset were mapped to the 3D structure of human 26S proteasome (PDB id: 5GJQ) utilizing residue level mappings between Uniprot and PDB entries obtained from SIFTS[25] database. In cases where multiple positions within the PDB structure were valid, the mapping with the shortest distance was prioritized. For the re-analyzed mouse mitochondrial XL-MS dataset[10], the cross-links were

mapped to homologous complexes (PDB IDs 1EUC, 1T9G, 5LNK, 1ZOY, 1NTM, 1V54) as shown previously[10]. In brief, the protein sequences for all proteins involved in detected cross-links were aligned against a reference database containing PDB sequences of interest using BLAST[26]. All BLAST matches with significant E-value and percent identity greater than 70% were retained. Exact positions for each cross-link were mapped against homologous PDB structures using a pairwise alignment, and cross-links were only considered successfully mapped if the cross-linked lysine was conserved in the structure. In cases where multiple positions within the PDB structure were valid, the mapping with the shortest distance was prioritized. Any cross-links where the exact position of the cross-linked lysine was not structurally resolved in a homologous PDB structure, were considered partially mapped. Because SIFT residue-level mapping for most of the representative structures (PDB IDs 2VRH, 1DKG, 1PCQ, 3JCD, 4PC1, and 2LRX) was unavailable for the *E. coli* dataset[13], we utilized the above mentioned homology based approach and the closest homologous complexes (PDB IDs 5MY1, 5ADY, 5ME0, 2RDO, 2VRH, 4JK2, 4YLN, 4YLO, 4XO2, 4YFH, 4YF0).

### Fraction of structure-corroborating identifications (FSI)

FSI can be calculated using the following equation:

$$
\begin{aligned}
&\text{FSI}(\%) \\
&= \frac{\text{Number interprotein XLs with in the maximum Euclidean distance constraint of the linker}}{\text{Number interprotein XLs with at least one of the two linked residues mapped to structure}} \\
&\quad X \\
&\quad 100
\end{aligned}
\tag{1}
$$

In this work, we used 30Å as the maximum distance constraint for DSSO.

### Fraction of mis-identifications (FMI)

FMI is the fraction of cross-link identifications from a false search space (from an unrelated organism) among all the identified cross-links. It can be calculated using the following equation:

$$
\text{FMI}(\%) = \frac{Number\ of\ mis-identifications}{Total\ number\ of\ identifications}\ X\ 100
\tag{2}
$$

In the current work, all the raw files were re-analyzed against a sequence database containing all the sequences from the target organism's proteome and all the sequences from *S. cerevisiae* proteome. Then the fraction of mis-identifications i.e., cross-links with at least one of the two linked residues unambiguously mapped to proteins from *S. cerevisiae* is calculated (if any cross-link has peptide shared between homologous proteins from the target organism and *S. cerevisiae*, it was considered a true identification). Importantly, when choosing an unrelated organism, it is important to make sure that there is no potential experimental contamination with proteins from that organism. It should be noted that another decoy database (reverse sequences of proteomes from both organisms) is generated for the FDR calculation by Proteome Discoverer. It is also noteworthy that FMI is estimated after the cross-link results are filtered at a conventional FDR threshold ('1% FDR' in the

current study). Additionally, it should be pointed out that similar to the conventional FDR calculations[27], FMI calculations can also be sensitive to drastic differences in sizes of the proteome database of the unrelated organism. We utilized the following equation adapted from Fischer and Rappsilber[7] to account for difference in database sizes and observed similar trend to that of uncorrected FMI across all three data sets analyzed in the current study (Extended data Fig.5).

$$\text{FMI}_{\text{corrected}}(\%) = \frac{\text{TD}+\text{DD}\left(1 - \frac{TD_{DB}}{DD_{DB}}\right)}{\text{TT}} \times 100 \tag{3}$$

where, TT is the number of target-target matches, DD is the number of decoy-decoy matches, and TD is number of target-decoy and decoy-target matches. $TD_{DB}$ is the number of all possible unique target-decoy and decoy-target peptide pairs and $DD_{DB}$ is the number of all possible unique decoy-decoy peptide pairs.

### Fraction of interprotein cross-links from known interactions (FKI)

FKI for proteome-wide XL-MS studies can be defined as the fraction of the identified interprotein cross-links from previously known protein-protein interactions. It can be derived using the following equation:

$$\text{FKI}(\%) = \frac{\text{Number of true positives}}{\text{Total number of postives}} X\ 100 \tag{4}$$

where, "positives" refer to all the identified interprotein cross-links, and "true positives" refer to interprotein cross-links from known protein-protein interactions. If a given interprotein cross-link represents multiple potential interactions and at least one of those potential interactions was mapped to the list of known protein-protein interactions, it was counted as a "true positive". We compiled the known protein-protein interactions for *E. coli* (24,745), Mouse (40,527) and Human (336,033) from seven primary interaction databases. These databases include IMEx[28] partners IntAct[29], MINT[30], and DIP[31]; IMEx observer BioGRID[32]; and additional sources HPRD[33], MIPS[34], and iRefWeb[35]. Furthermore, iRefWeb combines interaction data from CORUM[36], BIND[37], MPPI[34] and OPHID[38]. We converted all gene identifiers in each database to Entrez gene IDs and then mapped to Uniprot IDs.

We would like to point out that FSI and FKI are calculated using similar denominator, conceptually. For FSI, the dominator consists of all interprotein XLs with at least one of the two peptides mapped to the reference structure. In case of FKI, the denominator consists of all the interprotein XLs in the dataset. Even though FKI's equation does not explicitly require all the XLs to have at least one of the two proteins to be present in the reference interactome database, we expect that almost all interprotein XLs to satisfy this criterion. Moreover, we analyzed all the datasets from the current study and noted that all the datasets have more than 97% of all their interprotein XLs with at least one of the proteins in the reference interactome database. We acknowledge that someone who has a smaller reference database might not note the same observation. However, we argue that such case would lead

to underestimation of FKI (i.e., overestimation of error rate), thereby making FKI more stringent.

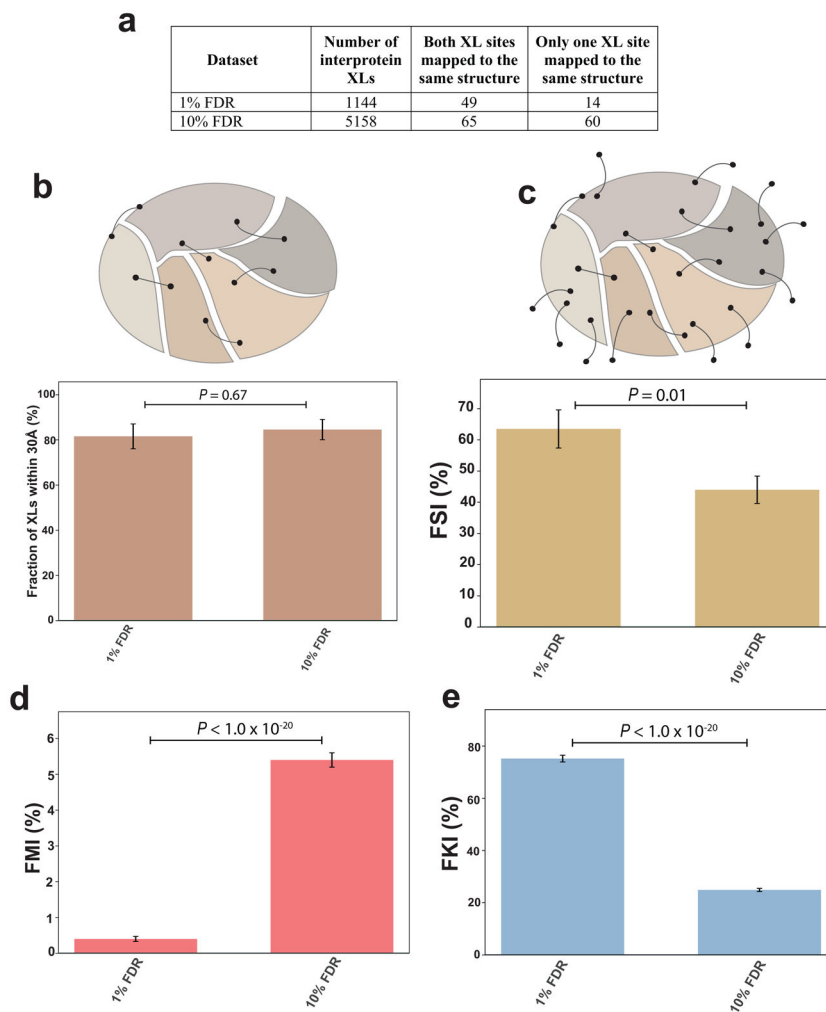### Fraction of validated novel interactions using orthogonal experiment, namely protein complementation assay (PCA)

The ORFs of novel protein-protein interactions in pDONR223 vector was innoculated from hORFeome v8.1 library[39]. In each of the categories, namely '1% FDR with XlinkX Score 50', '1% FDR', and '10% FDR', 93 protein pairs were randomly picked without any overlaps between categories. The Gateway LR reactions were performed to clone the individual bait and prey protein of each protein pair into the expression plasmids containing the complementation fragments of the fluorescent protein Venus. To perform the assay, the HEK293T cells were prepared in Dulbecco's Modified Eagle Medium (DMEM) supplemented with 10% fetal bovine serum (FBS) (ATCC) in black 96-well flat-bottom plates (Costar) with 5% $CO_2$ at 37°C. When reaching 60-70% confluency, the cells were co-transfected with both plasmids containing the Venus fragments-tagged bait and prey ORF (100 ng for each) which were pre-mixed and incubated with polyethylenimine (PEI) (Polysciences Inc.) and OptiMEM (Gibco). For positive and negative controls, the set containing previously published 92 positive reference pairs and 92 negative reference pairs were simultaneously examined[19, 40]. After 58 hours, the fluorescence intensity of the transfected cells was measured and recorded using Infinite M1000 microplate reader (Tecan) (excitation = $514 \pm 5$ nm / emission = $527 \pm 5$ nm). The PCA experiments were performed and analyzed in triplicate. We performed a statistical power analysis (Using in-built R v3.6.3 functions and Python 2.7) and confirmed that using 92 interactions would give us >97% power to detect the difference for the 'Positive Reference Set (PRS)' versus 'Random Reference Set (RRS)', '1% FDR with XlinkXScore 50' versus the '1% FDR' and '10% FDR' datasets. The effect sizes (Cohen's d) were calculated from the means and pooled standard deviation of given two groups under comparison (all effect sizes were large i.e., d > 0.8). The results are provided in Supplementary Table 1. Additionally, a short discussion on the utility of PCA to validate interactions from large-scale XL-MS studies on cell organelles and different organisms is provided in Supplementary Note 5.

## Data Availability Statement

The human K562 XL-MS raw files [122 raw files (97 HILIC and 25 SCX fractions) from our recent proteome-wide human K562 XL-MS study[2]] analyzed in this study have been deposited to the ProteomeXchange Consortium via the PRIDE[41] partner repository with the dataset identifier PXD018771. Raw data from our PCA experiments are available from the corresponding author upon request. Protein sequences were obtained from Uniprot database (https://www.uniprot.org/). Residue level mapping was performed using data from SIFTS database (https://www.ebi.ac.uk/pdbe/docs/sifts/index.html). Protein three dimensional structures utilized in this study were obtained from PDB (https://www.rcsb.org/ ; Accession codes: 5GJQ, 1EUC, 1T9G, 5LNK, 1ZOY, 1NTM, 1V54, 5MY1, 5ADY, 5ME0, 2RDO, 2VRH, 4JK2, 4YLN, 4YLO, 4XO2, 4YFH and 4YF0).

## Extended Data
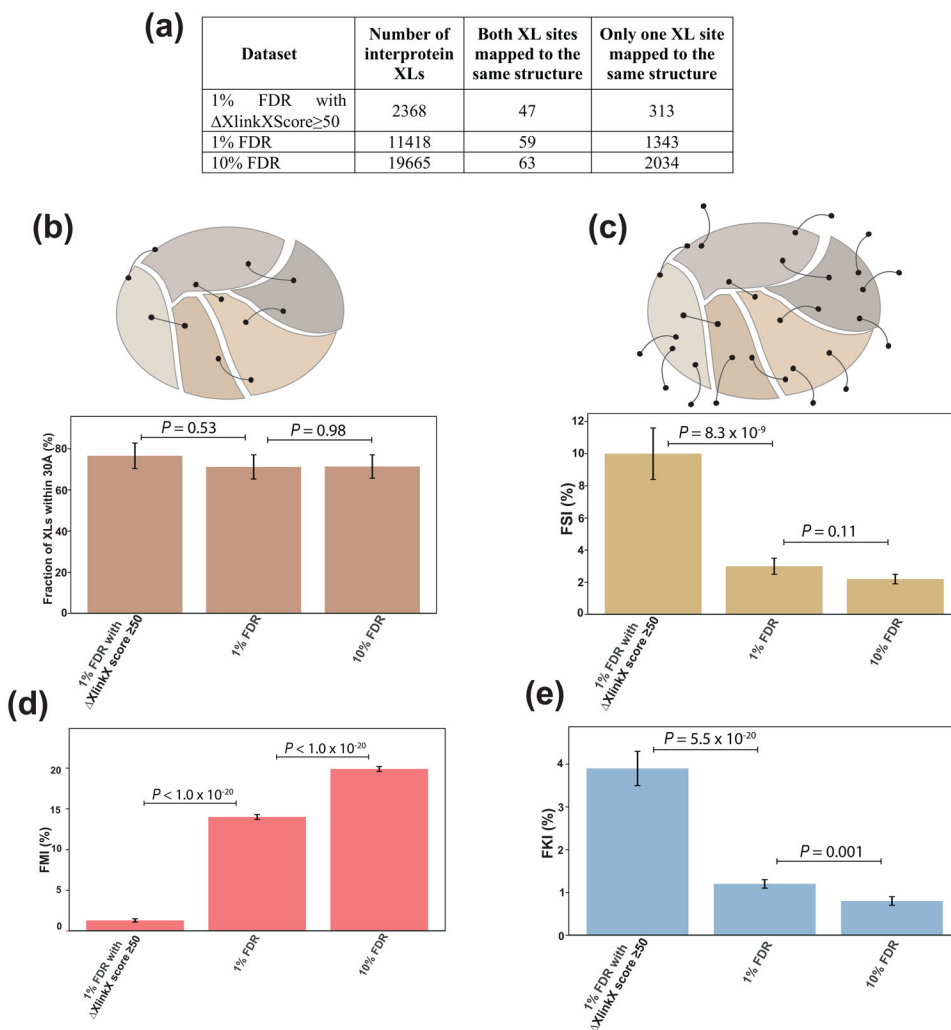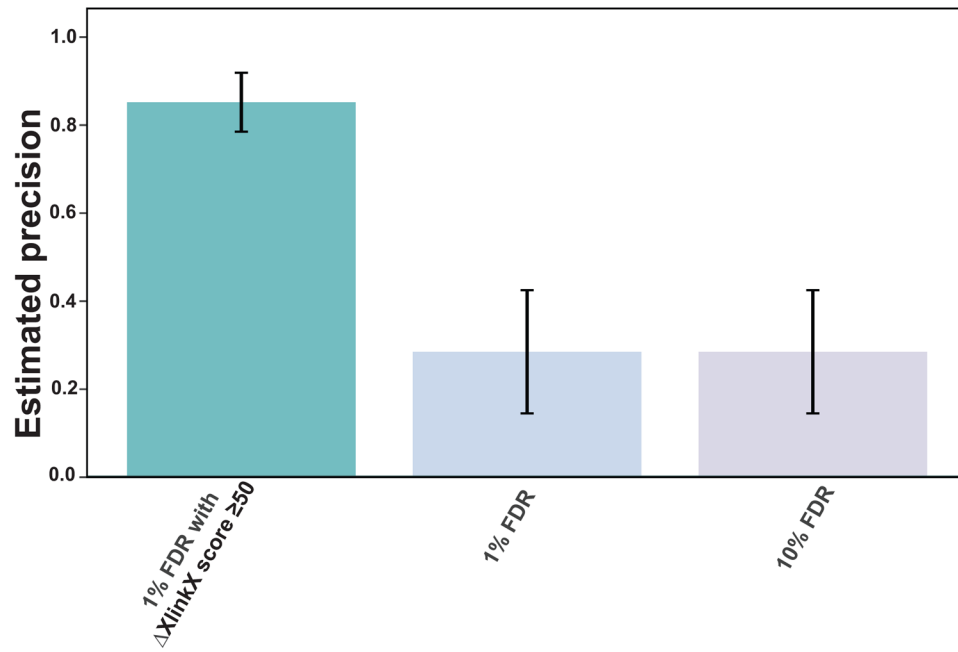


**a**

| Dataset | Number of interprotein XLs | Both XL sites mapped to the same structure | Only one XL site mapped to the same structure |
|---|---|---|---|
| 1% FDR | 1144 | 49 | 14 |
| 10% FDR | 5158 | 65 | 60 |

**Extended Data Fig. 1. Analysis of the human proteome-wide XL-MS dataset using MaXLinker software.**

**(a)** Table showing the number of cross-links obtained at different filtering criteria, and upon mapping to a representative 3D structure of a human 26S proteasome (PDB id: 5GJQ). **(b)** Comparison of the fraction of validated cross-links using the conventional structure-based approach (n = 49 XLs for '1% FDR'; n = 65 XLs for '10% FDR). **(c)** Comparison using the fraction of structure-corroborating identifications (FSI) (n = 63 XLs for '1% FDR'; n = 125 XLs for '10% FDR). **(d)** Comparison using the fraction of mis-identifications (FMI) (n = 8127 XLs for '1% FDR'; n = 15110 XLs for '10% FDR). **(e)** Comparison using the fraction of interprotein cross-links from known interactions (FKI) (n = 1144 XLs for '1% FDR'; n = 5158 XLs for '10% FDR). for **(b-e)**, the P values were calculated using a two-sided Z-test and the error bars indicate +/− SE of proportion.

**(a)**

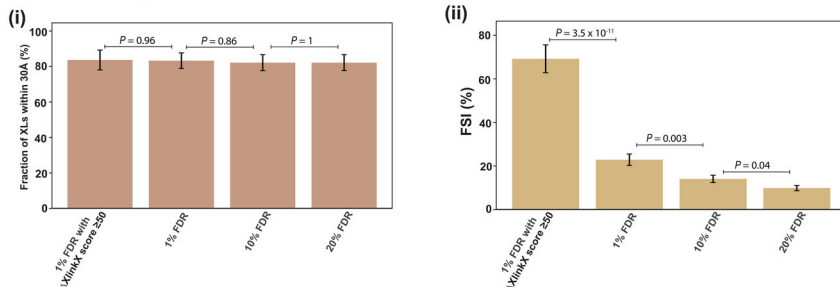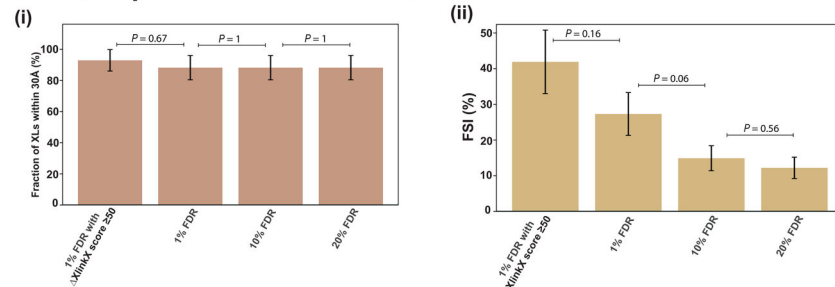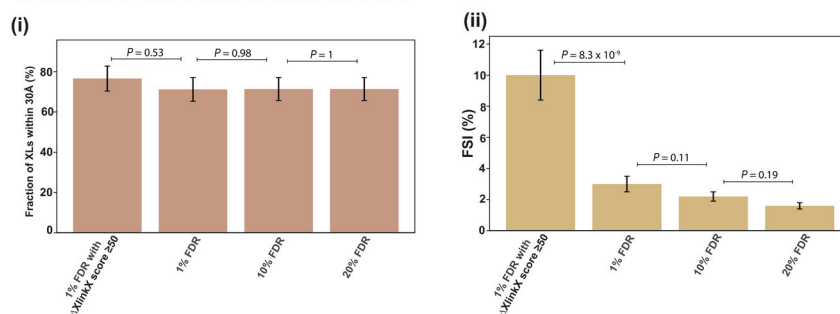| Dataset | Number of interprotein XLs | Both XL sites mapped to the same structure | Only one XL site mapped to the same structure |
|---|---|---|---|
| 1% FDR with ΔXlinkXScore≥50 | 2368 | 47 | 313 |
| 1% FDR | 11418 | 59 | 1343 |
| 10% FDR | 19665 | 63 | 2034 |

**(b)**



**(c)**



**(d)**



**(e)**



**Extended Data Fig. 2. Demonstration of the utility of our comprehensive set of validation metrics on a publicly available mouse mitochondrial XL-MS dataset.**

**(a)** Table showing the number of cross-links obtained at different filtering criteria, and upon mapping to representative 3D structures. **(b)** Conventional structure-based validation (n = 47 XLs for '1% FDR with ΔXlinkX score ≥50'; n = 59 XLs for '1% FDR'; n = 63 XLs for '10% FDR'). **(c)** Fraction of structure-corroborating identifications (FSI) (n = 360 XLs for '1% FDR with ΔXlinkX score ≥50'; n = 1402 XLs for '1% FDR'; n = 2097 XLs for '10% FDR'). **(d)** Fraction of mis-identifications (FMI) (n = 4814 XLs for '1% FDR with ΔXlinkX score ≥50'; n = 15323 XLs for '1% FDR'; n = 24317 XLs for '10% FDR'). **(e)** Fraction of interprotein cross-links from known interactions (FKI) (n = 2368 XLs for '1% FDR with ΔXlinkX score ≥50'; n = 11418 XLs for '1% FDR'; n = 19665 XLs for '10% FDR'). P values in **(b-e)** were calculated using a two-sided Z-test and the error bars indicate +/− SE of proportion.

**Extended Data Fig. 3. Estimated precision using PCA experiments for the three datasets of different quality from our human K562 proteome-wide XL-MS study. Derived from Fig.1g.**
(n = 3 independent experiments; See Methods). The error bars indicate +/− SE of proportion (see Supplementary Note 2 for a detailed description of the methodology).

Yugandhar et al. Page 12

## a. Human proteome-wide XL-MS



## b. *E. coli* proteome-wide XL-MS



## c. Mouse mitochondrial XL-MS



**Extended Data Fig. 4. Structure-based mapping analysis at 20% FDR, extension to the analysis shown in Fig.1, Fig.2, and Extended Data Fig.2.**
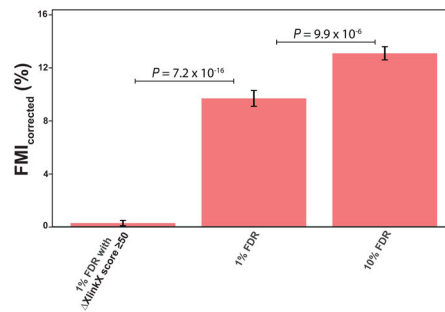
**a.** Human proteome-wide XL-MS study: (i) Conventional structure-based validation (n = 43 XLs for '1% FDR with  XlinkX score  50'; n = 72 XLs for '1% FDR'; n = 73 XLs for '10% FDR'; n = 73 XLs for '20% FDR'). (ii) Fraction of structure-corroborating identifications (FSI) (n = 52 XLs for '1% FDR with  XlinkX score  50'; n = 262 XLs for '1% FDR'; n = 426 XLs for '10% FDR'; n = 605 XLs for '20% FDR'). **b.** *E. coli* proteome-wide XL-MS study: (i) Conventional structure-based validation (n = 14 XLs for '1% FDR with  XlinkX score  50'; n = 17 XLs for '1% FDR'; n = 17 XLs for '10% FDR'; n = 17 XLs for '20% FDR'). (ii) Fraction of structure-corroborating identifications (FSI) (n = 31 XLs for '1% FDR with  XlinkX score  50'; n = 55 XLs for '1% FDR'; n = 101 XLs for '10% FDR'; n = 123 XLs for '20% FDR').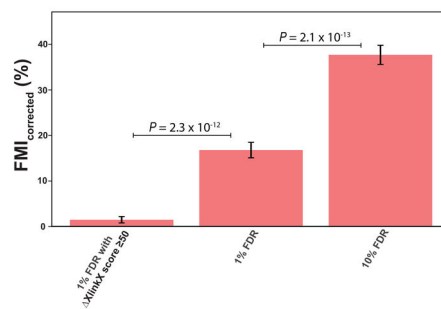 **c.** Mouse mitochondrial XL-MS study: (i) Conventional structure-based validation (n = 47 XLs for '1% FDR with  XlinkX score  50'; n = 59 XLs for '1% FDR'; n = 63 XLs for '10% FDR'; n = 63 XLs for '20% FDR'). (ii) Fraction of structure-corroborating identifications (FSI) (n = 360 XLs for '1% FDR with  XlinkX score  50'; n = 1402 XLs for '1% FDR'; n = 2097 XLs for '10% FDR'; n = 2751

XLs for '20% FDR'). P values in all the panels were calculated using a two-sided Z-test and the error bars indicate +/− SE of proportion.
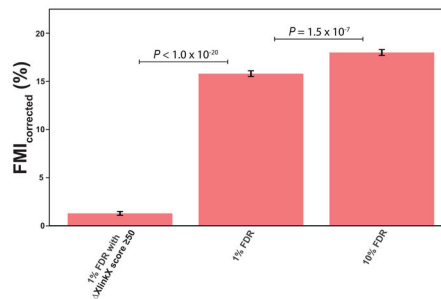
## a. Human proteome-wide XL-MS



## b. *E. coli* proteome-wide XL-MS



## c. Mouse mitochondrial XL-MS



**Extended Data Fig. 5. Corrected FMI for the three datasets analyzed in the study (Utilizing equation 3 from Methods section).**
**(a)** Human proteome-wide XL-MS (n = 668 XLs for '1% FDR with XlinkX score 50'; n = 3029 XLs for '1% FDR'; n = 4957 XLs for '10% FDR'). **(b)** *E. coli* proteome-wide XL-MS (n = 340 XLs for '1% FDR with XlinkX score 50'; n = 553 XLs for '1% FDR'; n = 755 XLs for '10% FDR'). **(c)** Mouse mitochondrial XL-MS (n = 4814 XLs for '1% FDR with XlinkX score 50'; n = 15323 XLs for '1% FDR'; n = 24317 XLs for '10% FDR'). P values in all the panels were calculated using a two-sided Z-test and the error bars indicate +/− SE of proportion.

## Supplementary Material

Refer to Web version on PubMed Central for supplementary material.

## Acknowledgements

## References

1. Yu C & Huang L Cross-Linking Mass Spectrometry: An Emerging Technology for Interactomics and Structural Biology. Analytical Chemistry 90, 144–165 (2018). [PubMed: 29160693]

2. Yugandhar K et al. MaXLinker: Proteome-wide Cross-link Identifications with High Specificity and Sensitivity. Molecular & Cellular Proteomics 19, 554 (2020). [PubMed: 31839598]

3. Iacobucci C, Götze M & Sinz A Cross-linking/mass spectrometry to get a closer view on protein interaction networks. Current Opinion in Biotechnology 63, 48–53 (2020). [PubMed: 31891863]

4. Ferber M et al. Automated structure modeling of large protein assemblies using crosslinks as distance restraints. Nature Methods 13, 515–520 (2016). [PubMed: 27111507]

5. Karaca E, Rodrigues JPGLM, Graziadei A, Bonvin AMJJ & Carlomagno T M3: an integrative framework for structure determination of molecular machines. Nature Methods 14, 897 (2017). [PubMed: 28805795]

6. Hauri S et al. Rapid determination of quaternary protein structures in complex biological samples. Nature Communications 10, 192 (2019).

7. Fischer L & Rappsilber J Quirks of Error Estimation in Cross-Linking/Mass Spectrometry. Analytical Chemistry 89, 3829–3833 (2017). [PubMed: 28267312]

8. O'Reilly FJ & Rappsilber J Cross-linking mass spectrometry: methods and applications in structural, molecular and systems biology. Nature Structural & Molecular Biology 25, 1000–1008 (2018).

9. Klykov O et al. Efficient and robust proteome-wide approaches for cross-linking mass spectrometry. Nature Protocols 13, 2964–2990 (2018). [PubMed: 30446747]

10. Liu F, Lössl P, Rabbitts BM, Balaban RS & Heck AJR The interactome of intact mitochondria by cross-linking mass spectrometry provides evidence for coexisting respiratory supercomplexes. Molecular & Cellular Proteomics 17, 216 (2018). [PubMed: 29222160]

11. Keller A, Chavez JD, Felt KC & Bruce JE Prediction of an Upper Limit for the Fraction of Interprotein Cross-Links in Large-Scale In Vivo Cross-Linking Studies. Journal of Proteome Research (2019).

12. Bartolec TK et al. Cross-linking Mass Spectrometry Analysis of the Yeast Nucleus Reveals Extensive Protein–Protein Interactions Not Detected by Systematic Two-Hybrid or Affinity Purification-Mass Spectrometry. Analytical Chemistry 92, 1874–1882 (2020). [PubMed: 31851481]

13. Liu F, Lössl P, Scheltema R, Viner R & Heck AJR Optimized fragmentation schemes and data analysis strategies for proteome-wide cross-link identification. Nature Communications 8, 15473 (2017).

14. Chen Z-L et al. A high-speed search engine pLink 2 with systematic evaluation for proteome-scale identification of cross-linked peptides. Nature Communications 10, 3404 (2019).

15. Götze M, Iacobucci C, Ihling CH & Sinz A A Simple Cross-Linking/Mass Spectrometry Workflow for Studying System-wide Protein Interactions. Analytical Chemistry 91, 10236–10244 (2019). [PubMed: 31283178]

16. Yu H et al. High-Quality Binary Protein Interaction Map of the Yeast Interactome Network. Science 322, 104 (2008). [PubMed: 18719252]

17. Vo Tommy V. et al. A Proteome-wide Fission Yeast Interactome Reveals Network Evolution Principles from Yeasts to Human. Cell 164, 310–323 (2016). [PubMed: 26771498]

18. Nyfeler B, Michnick SW & Hauri H-P Capturing protein interactions in the secretory pathway of living cells. Proceedings of the National Academy of Sciences of the United States of America 102, 6350 (2005). [PubMed: 15849265]

19. Braun P et al. An experimentally derived confidence score for binary protein-protein interactions. Nature Methods 6, 91 (2008). [PubMed: 19060903]

20. Makowski MM, Willems E, Jansen PWTC & Vermeulen M Cross-linking immunoprecipitation-MS (xIP-MS): Topological Analysis of Chromatin-associated Protein Complexes Using Single Affinity Purification. Molecular & Cellular Proteomics 15, 854 (2016). [PubMed: 26560067]

21. Beveridge R, Stadlmann J, Penninger JM & Mechtler K A synthetic peptide library for benchmarking crosslinking-mass spectrometry search engines for proteins and protein complexes. Nature Communications 11, 742 (2020).

22. Rual J-F et al. Towards a proteome-scale map of the human protein–protein interaction network. Nature 437, 1173–1178 (2005). [PubMed: 16189514]

23. Makowski MM, Willems E, Jansen PWTC & Vermeulen M Cross-linking immunoprecipitation-MS (xIP-MS): Topological Analysis of Chromatin-associated Protein Complexes Using Single Affinity Purification. Molecular & Cellular Proteomics 15, 854 (2016). [PubMed: 26560067]

## Methods-only References

24. The UniProt Consortium UniProt: the universal protein knowledgebase. Nucleic Acids Research 45, D158–D169 (2017). [PubMed: 27899622]

25. Dana JM et al. SIFTS: updated Structure Integration with Function, Taxonomy and Sequences resource allows 40-fold increase in coverage of structure-based annotations for proteins. Nucleic Acids Research 47, D482–D489 (2018).

26. Altschul SF, Gish W, Miller W, Myers EW & Lipman DJ Basic local alignment search tool. Journal of Molecular Biology 215, 403–410 (1990). [PubMed: 2231712]

27. Gupta N, Bandeira N, Keich U & Pevzner PA Target-Decoy Approach and False Discovery Rate: When Things May Go Wrong. Journal of The American Society for Mass Spectrometry 22, 1111–1120 (2011). [PubMed: 21953092]

28. Orchard S et al. Protein interaction data curation: the International Molecular Exchange (IMEx) consortium. Nature Methods 9, 345 (2012). [PubMed: 22453911]

29. Kerrien S et al. The IntAct molecular interaction database in 2012. Nucleic Acids Research 40, D841–D846 (2012). [PubMed: 22121220]

30. Licata L et al. MINT, the molecular interaction database: 2012 update. Nucleic Acids Research 40, D857–D861 (2012). [PubMed: 22096227]

31. Salwinski L et al. The Database of Interacting Proteins: 2004 update. Nucleic Acids Research 32, D449–D451 (2004). [PubMed: 14681454]

32. Chatr-aryamontri A et al. The BioGRID interaction database: 2015 update. Nucleic Acids Research 43, D470–D478 (2015). [PubMed: 25428363]

33. Keshava Prasad TS et al. Human Protein Reference Database—2009 update. Nucleic Acids Research 37, D767–D772 (2009). [PubMed: 18988627]

34. Pagel P et al. The MIPS mammalian protein–protein interaction database. Bioinformatics 21, 832–834 (2005). [PubMed: 15531608]

35. Turner B et al. iRefWeb: interactive analysis of consolidated protein interaction data and their supporting evidence. Database 2010, baq023–baq023 (2010).

36. Ruepp A et al. CORUM: the comprehensive resource of mammalian protein complexes—2009. Nucleic Acids Research 38, D497–D501 (2010). [PubMed: 19884131]

37. Alfarano C et al. The Biomolecular Interaction Network Database and related tools 2005 update. Nucleic Acids Research 33, D418–D424 (2005). [PubMed: 15608229]

38. Brown KR & Jurisica I Online Predicted Human Interaction Database. Bioinformatics 21, 2076–2082 (2005). [PubMed: 15657099]

39. Yang X et al. A public genome-scale lentiviral expression library of human ORFs. Nature Methods 8, 659 (2011). [PubMed: 21706014]

40. Venkatesan K et al. An empirical framework for binary interactome mapping. Nature Methods 6, 83 (2008). [PubMed: 19060904]

41. Perez-Riverol Y et al. The PRIDE database and related tools and resources in 2019: improving support for quantification data. Nucleic Acids Research 47, D442–D450 (2018).
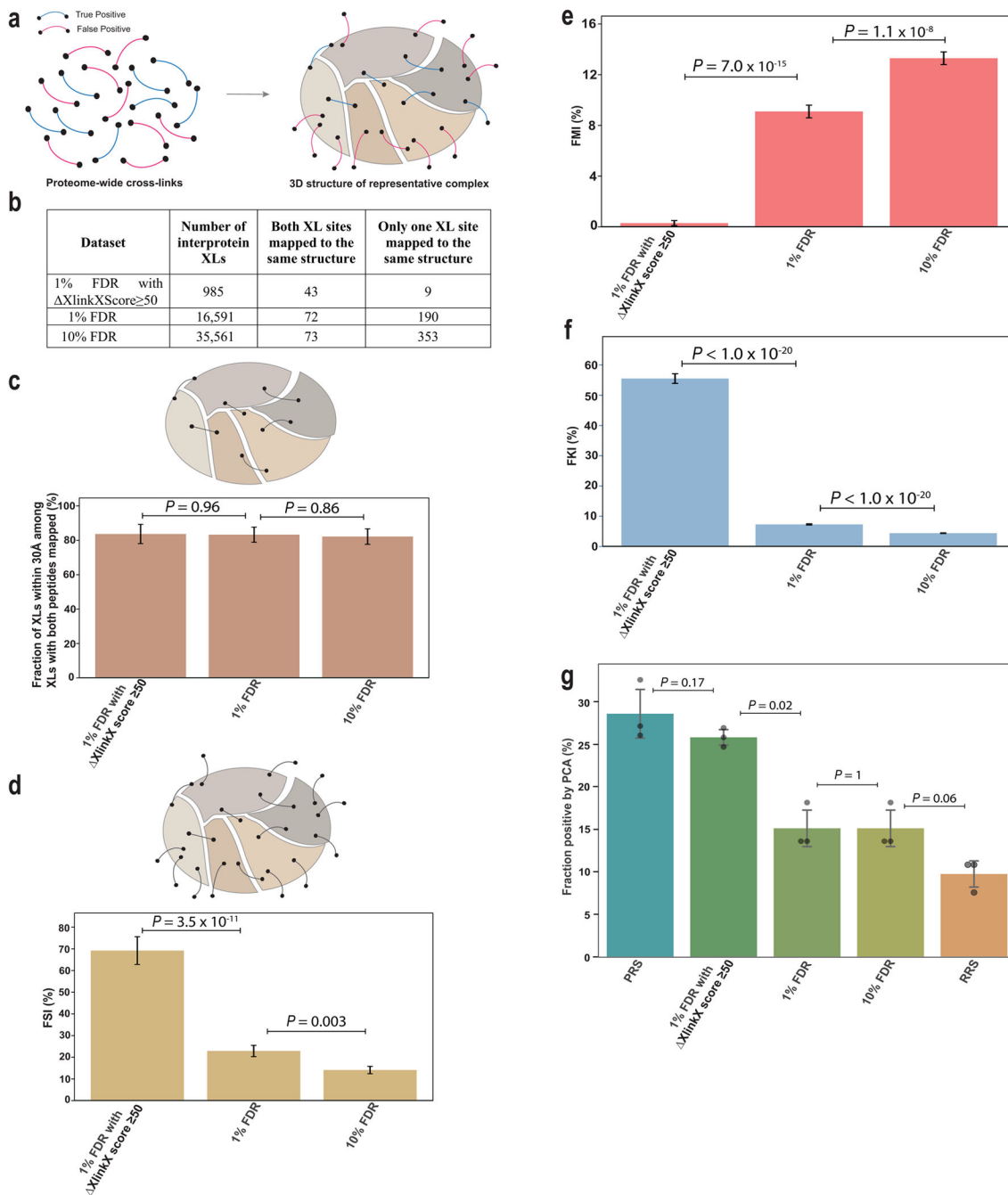
**a** True Positive / False Positive

Proteome-wide cross-links → 3D structure of representative complex

**b**

| Dataset | Number of interprotein XLs | Both XL sites mapped to the same structure | Only one XL site mapped to the same structure |
|---|---|---|---|
| 1% FDR with ΔXlinkXScore≥50 | 985 | 43 | 9 |
| 1% FDR | 16,591 | 72 | 190 |
| 10% FDR | 35,561 | 73 | 353 |

**c** Fraction of XLs within 30Å among XLs with both peptides mapped (%)

$P = 0.96$    $P = 0.86$

**d** FSI (%)

$P = 3.5 \times 10^{-11}$

$P = 0.003$

**e** FMI (%)

$P = 7.0 \times 10^{-15}$    $P = 1.1 \times 10^{-8}$

**f** FKI (%)

$P < 1.0 \times 10^{-20}$

$P < 1.0 \times 10^{-20}$

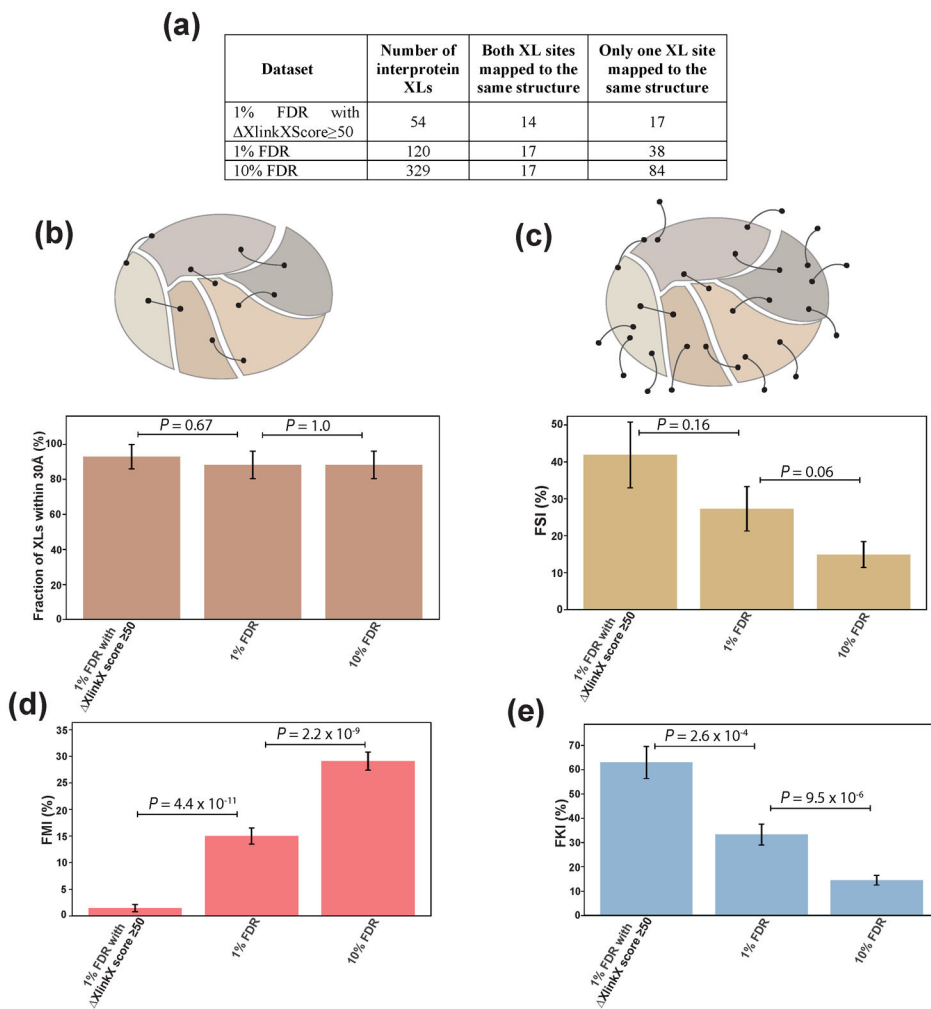**g** Fraction positive by PCA

$P = 0.17$    $P = 0.02$    $P = 1$    $P = 0.06$

**Figure 1.**
Evaluation of the conventional 3D structure-based validation approach for proteome-wide
cross-linking mass spectrometry using human K562 DSSO XL-MS data[2]. (**a**) In the current
structure-mapping approach for validating the false positive cross-link identifications, most
false positive cross-links only have one peptide mapped to the structure and are therefore
ignored. (**b**) Table showing the number of cross-links obtained at different filtering criteria,
and upon mapping to a representative 3D structure of a human 26S proteasome. (**c**) Fraction
of cross-links satisfying the maximum distance constraint ( 30Å) across the three sets,

according to the conventional structure-based validation approach (n = 43 XLs for '1% FDR with ⪰ XlinkX score ⪰ 50'; n = 72 XLs for '1% FDR'; n = 73 XLs for '10% FDR'). **(d)** Fraction of structure-corroborating identifications (FSI) across the three sets (n = 52 XLs for '1% FDR with ⪰ XlinkX score ⪰ 50'; n = 262 XLs for '1% FDR'; n = 426 XLs for '10% FDR'). **(e)** Fraction of misidentifications (FMI) across the three sets (n = 668 XLs for '1% FDR with ⪰ XlinkX score ⪰ 50'; n = 3029 XLs for '1% FDR'; n = 4957 XLs for '10% FDR'; See Methods). **(f)** Fraction of interprotein cross-links from known interactions (FKI) across the three sets (n = 985 XLs for '1% FDR with ⪰ XlinkX score ⪰ 50'; n = 16,591 XLs for '1% FDR'; n = 35,561 XLs for '10% FDR'). For **(c – f)**, *P* values were calculated using a two-sided *Z*-test. The error bars indicate +/− SE of proportion and the centre of the error bars indicate the proportion. **(g)** Orthogonal experimental validation of a random subset of novel interactions from the three sets using protein complementation assay (PCA). [PRS(Positive Reference Set): mean fraction positive: 0.286; RRS (Random Reference Set): mean fraction positive: 0.098; '10% FDR': mean fraction positive: 0.152; '1% FDR': mean fraction positive: 0.152; '1% FDR with ⪰ XlinkX score ⪰ 50': mean fraction positive: 0.258]. The error bars indicate +/− standard deviation and the centre of the error bars indicate mean fraction positive; *P* values were calculated using a two-sided t-test on the log-transformed measurements (n = 3 independent experiments; See Methods); 95% confidence interval; t-statistic 4.04 for '10% FDR' versus RRS, 7.20 for '1% FDR with ⪰ XlinkX score ⪰ 50' versus '1% FDR', 2.13 for PRS versus '1% FDR with ⪰ XlinkX score ⪰ 50'; 2 degrees of freedom.

**(a)**

| Dataset | Number of interprotein XLs | Both XL sites mapped to the same structure | Only one XL site mapped to the same structure |
|---|---|---|---|
| 1% FDR with ΔXlinkXScore≥50 | 54 | 14 | 17 |
| 1% FDR | 120 | 17 | 38 |
| 10% FDR | 329 | 17 | 84 |



**Figure 2.**
Demonstration of our set of validation metrics on a publicly available *E. coli* proteome-wide XL-MS dataset[13]. **(a)** Table showing the number of cross-links obtained at different filtering criteria, and upon mapping to representative 3D structures. **(b)** Fraction of cross-links satisfying the maximum distance constraint ( 30Å) across the three sets, according to the conventional structure-based validation approach (n = 14 XLs for '1% FDR with XlinkX score 50'; n = 17 XLs for '1% FDR'; n = 17 XLs for '10% FDR'). **(c)** Fraction of structure-corroborating identifications (FSI) (n = 31 XLs for '1% FDR with XlinkX score 50'; n = 55 XLs for '1% FDR'; n = 101 XLs for '10% FDR'). **(d)** Fraction of mis-identifications (FMI) (n = 340 XLs for '1% FDR with XlinkX score 50'; n = 553 XLs for '1% FDR'; n = 755 XLs for '10% FDR'). **(e)** Fraction of interprotein cross-links from known interactions (FKI) (n = 54 XLs for '1% FDR with XlinkX score 50'; n = 120 XLs for '1% FDR'; n = 329 XLs for '10% FDR'). For **(b-e),** the *P* values were calculated using a two-sided Z-test and the error bars indicate +/− SE of proportion.