

# Enabling cross-indication protein expression analysis using a curated pan-cancer dataset and a tailored workflow

Received: 30 October 2024

Accepted: 16 March 2026

Published online: 23 March 2026

Cite this article as: Wang J., Tian X., Yu W. *et al.* Enabling cross-indication protein expression analysis using a curated pan-cancer dataset and a tailored workflow. *Sci Rep* (2026). <https://doi.org/10.1038/s41598-026-44872-z>

Jixin Wang, Xiaowen Tian, Wen Yu, Benjamin S. Pullman, John Bullen, Elaine Hurt & Wenyan Zhong

We are providing an unedited version of this manuscript to give early access to its findings. Before final publication, the manuscript will undergo further editing. Please note there may be errors present which affect the content, and all legal disclaimers apply.

If this paper is publishing under a Transparent Peer Review model then Peer Review reports will publish with the final article.

ARTICLE IN PRESS

**Enabling cross-indication protein expression analysis using a curated pan-cancer dataset  
and a tailored workflow**

Jixin Wang<sup>1</sup>, Xiaowen Tian<sup>2</sup>, Wen Yu<sup>3</sup>, Ben Pullman<sup>4</sup>, John Bullen, Jr.<sup>5</sup>, Elaine Hurt<sup>5</sup>, Wenyan  
Zhong<sup>\*6</sup>

<sup>1</sup>Oncology Data Science, AstraZeneca, Gaithersburg, MD, USA; <sup>2</sup>Statistical Innovation,  
Oncology R&D, AstraZeneca, Gaithersburg, MD, USA; <sup>3</sup>Data Science and AI,  
BioPharmaceuticals R&D, AstraZeneca, Gaithersburg, MD, USA; <sup>4</sup>Centre for Genomics  
Research, Discovery Sciences, R&D, AstraZeneca, Gaithersburg, MD, USA; <sup>5</sup>Oncology  
Targeted Delivery, Research and Early Development, Oncology R&D, AstraZeneca,  
Gaithersburg, MD, USA; <sup>6</sup>Oncology Data Science, Oncology R&D, AstraZeneca, New York,  
NY, USA

**Running title:** Computational approaches for pan-cancer protein expression comparison

**Keywords:** CPTAC, proteomics, iBAQ, normalization, differential analysis, TCGA

Jixin Wang and Xiaowen Tian are co-first authors.

**Corresponding author:**

Wenyan Zhong

AstraZeneca, 430 East 29th Street, New York, NY, USA 10016

Phone: (301) 398-4209

Fax: (301) 398-9000

E-mail: [wenyan.zhong@astrazeneca.com](mailto:wenyan.zhong@astrazeneca.com)

## 1 **Abstract**

2 The National Cancer Institute's Clinical Proteomic Tumor Analysis Consortium (CPTAC)  
3 recently generated harmonized genomic, transcriptomic, proteomic, and clinical data for over  
4 1,000 tumors across 10 cohorts to facilitate pan-cancer discovery research. However, protein  
5 expression comparison across CPTAC cohorts remains challenging due to non-uniform missing  
6 data and varying protein expression distribution patterns across tumor types. To enable the  
7 cancer research community to conduct robust cross-cohort protein expressions, we present a  
8 curated and normalized pan-cancer protein expression dataset derived from the CPTAC pan-  
9 cancer study. Our workflow integrates systematic filtering, various missing data handling and  
10 normalization strategies. We developed a novel algorithm to select robustly expressed proteins in  
11 tumors within any CPTAC cohort; applied a cohort hybrid imputation approach to protein  
12 abundance values from FragPipe within each cohort, based on protein expression distribution  
13 patterns; and calculated intensity-based absolute quantification using protein abundance values  
14 and applied both global and smooth quantile normalization methods. Our analysis demonstrates  
15 that global quantile normalization surpasses both smooth quantile normalization and no  
16 normalization, as evidenced by its higher rank correlation across cancer cohorts between CPTAC  
17 and TCGA for selected proteins. The findings suggest that combining cohort hybrid imputation  
18 with global quantile normalization is an effective method for creating a normalized CPTAC pan-  
19 cancer protein dataset, which can facilitate the study of protein expression across different cancer  
20 types and accelerate cancer research.

21

## 22 **Introduction**

23 The identification and prioritization of therapeutic targets are pivotal in both drug development  
24 and biological understanding. Traditionally, target prioritization has been predominantly tailored  
25 to specific cohorts. However, recent developments in oncology underscore the significance of  
26 indication prioritization, where a single target holds relevance across multiple cancer types,  
27 thereby streamlining drug development and expediting the creation of novel therapies.

28 The Clinical Proteomic Tumor Analysis Consortium (CPTAC), established by the National  
29 Cancer Institute <sup>1</sup>, stands as a comprehensive resource of global mass-spectrometry based  
30 proteomics experiments. These initiatives enable high-throughput biological understanding by  
31 elucidating the relative and absolute expression of nearly the entire proteome across over 10  
32 tumor types, inclusive of both tumor tissue and normal adjacent tissue (NAT). The data  
33 generated by CPTAC has proven to be instrumental in target discovery and prioritization within  
34 tissue cohorts, leveraging a methodology to estimate absolute protein abundance from relative  
35 quantitation data <sup>2-4</sup>.

36 Pan-cancer proteomic efforts across various studies (e.g., ProteomicsDB and Expression Atlas)  
37 have been conducted, providing valuable resources in indication assessment<sup>5,6</sup>. However, a  
38 substantial challenge remains in comparing protein expression levels across cohorts due to a  
39 diverse array of samples, each exhibiting distinct missing data and expression patterns. In this  
40 study, we focus on the generation of a standardized, filtered, imputed, and normalized pan-  
41 cancer protein expression dataset from the CPTAC resource. Building on a systematic workflow  
42 that incorporates evaluation of multiple data processing approaches, we deliver a processed  
43 dataset ready for cross-cohort analyses, alongside workflow details to empower researchers in

44 adapting and refining our methodology to suit their specific biological questions. By making this  
45 standardized resource available, we aim to enable robust pan-cancer protein expression  
46 comparisons and support cancer treatment development and biological insight generation.

## 47 **Methods**

### 48 **Evaluation of CPTAC pan-cancer reprocessed proteomics data**

49 CPTAC pan-cancer samples were reprocessed using the FragPipe computational platform  
50 (version 15) with MSFragger (version 3.2)<sup>7</sup> and Philosopher (version 3.4.13)<sup>8</sup> for protein  
51 identification and quantitation as well as intensity-based absolute quantification (iBAQ)  
52 derivation and uniform FragPipe/MSFragger workflow with gene-level grouping and 1% FDR  
53 control, at the fragment and reporter ions level, both the normalization to the pool channel, and  
54 the median centering (polishing) can reduce the batch effects due to the sample preparation  
55 within the cohort to improve the relative quantitation, and iBAQ/riBAQ to stabilize length- and  
56 detectability-related variation as described in Wang et al.<sup>3</sup> FragPipe was used to assess protein  
57 abundance in tumor and normal tissues from 10 cancer indications in TMT labeled CPTAC  
58 dataset. Log<sub>2</sub>-transformed protein abundance was used for assessing missing value patterns and  
59 developing an imputation strategy. “Cohort” was defined as the unique combination of tissue  
60 type and indication, and “missing rate” was defined as the percentage of samples with missing  
61 values in each cohort.

### 62 **Algorithm for selection of robustly expressed proteins**

63 An algorithm was developed to select robustly expressed proteins in the samples as follows. For  
64 protein  $j$  in tumor samples from indication  $i$ , the percentage of samples with missing abundance

65 values was denoted as  $M_{ij}$  and the median abundance from non-missing observations was  
 66 denoted as  $A_{ij}$ . The corresponding percentile rank of  $A_{ij}$  within each cohort  $i$  was  
 67 obtained and denoted as  $F_i(A_{ij})$ . The protein  $j$  was kept in all cohorts if there was an  
 68 indication  $i$  such that  $F_i(A_{ij}) > 50\%$  and  $M_{ij} < 25\%$ . In other words, the protein was  
 69 kept if the median abundance percentile rank was  $>50\%$  and the corresponding missing rate was  
 70  $<25\%$  in at least one indication among tumor samples. The determination of these cutoffs was  
 71 empirical, aimed at achieving a balance between discarding an excessive number of proteins and  
 72 retaining those with a high rate of missing protein expression measurements.

### 73 **iBAQ conversion and normalization after cohort hybrid imputation**

74 The R package `imputeLCMD`<sup>9</sup> was used for imputation after filtering proteins that did not pass  
 75 the above criteria. This package employs a method selection approach that first classifies the  
 76 missingness mechanism for each protein.<sup>9</sup> It assumes that the mean protein abundances within  
 77 the cohort would follow a normal distribution if fully observed. The algorithm estimates the  
 78 parameters of the complete data distribution through quantile regression on observed feature  
 79 means. It then sets a censoring threshold by comparing the empirical cumulative distribution  
 80 function (CDF) of observed data against the theoretical CDF of this estimated complete  
 81 distribution. A point of maximum positive divergence is determined by computing an objective  
 82 function designed to quantify the relative difference between these two CDFs, particularly in the  
 83 lower intensity range. The objective function has the form of  $(\frac{F_e - F_o + 1}{F_o + 1} - 1) * 10$  where  $F_e$  and  $F_o$   
 84 represent the theoretical and empirical CDF, respectively. The intensity value at which this  
 85 objective function attains its maximal positive value statistically identifies the point of greatest  
 86 differences of observed cumulative probability compared to the theoretical, thereby establishing

87 the censoring threshold. Proteins with mean observed values at or below this censoring threshold  
88 are classified as Missing Not At Random (MNAR), while those above it are designated Missing  
89 At Random (MAR) or Missing Completely At Random (MCAR). If the determined missing  
90 mechanism was MNAR, quantile regression imputation of left-censored data was applied.  
91 Otherwise, k-nearest neighbor (KNN) imputation was applied. Previous studies have shown that  
92 KNN imputation performs better in TMT proteomics data<sup>10,11</sup>. Imputation was performed on  
93 robustly expressed proteins as determined by the method described above, in log<sub>2</sub> scale. iBAQ  
94 values were then calculated by dividing the raw imputed protein abundance by the number of  
95 theoretically observable peptides of the protein from FragPipe, as described by Wang et al.<sup>3</sup> The  
96 iBAQ values were then normalized using one of two quantile normalization methods: (1) global  
97 quantile normalization, with quantile normalization<sup>12</sup> applied to all samples across indications  
98 and tissue types together; or (2) smooth quantile normalization<sup>13</sup> with the assumption that the  
99 statistical distribution of each sample was the same within indications or tissue types, but  
100 allowing for differences between groups. Additionally, relative iBAQ (riBAQ)<sup>14</sup> and riBAQ-  
101 derived copy number<sup>15</sup> were calculated according to Wang et al.<sup>3</sup> The iBAQ values obtained  
102 from the imputed protein abundance was used to derive riBAQ and copy number values without  
103 applying additional normalization methods.

#### 104 **Differential protein expression analysis with normalized iBAQ**

105 To determine whether our missing data imputation and normalization strategy affected  
106 downstream analyses, we compared the fold change in differentially expressed proteins (DEP)  
107 between tumor and matched NAT samples for each indication, using non-normalized, global  
108 quantile-normalized, and smooth quantile-normalized protein iBAQ values.

109 FragPipe protein abundance–derived iBAQ data from CPTAC tumor and NAT samples were  
110 used for DEP analysis using R (version 4.1.1), with empirical Bayes statistics on protein-wise  
111 linear models with limma<sup>16</sup> embedded in the DEP package (version 1.16.0).<sup>17</sup> The correlations  
112 between the log<sub>2</sub> fold change from DEP analysis were examined between each comparison using  
113 Pearson correlation.

#### 114 **Comparison of the indication ranks of selected proteins between CPTAC and TCGA**

115 We identified proteins whose protein and RNA expression were highly correlated across CPTAC  
116 cohorts prior to normalization. Within each indication, we calculated the Pearson correlation  
117 between protein and RNA expression using tumor tissues. Then, proteins with a correlation  
118 greater than 0.5 across the majority of indications were kept for subsequent analysis. We then  
119 compared their protein expression rank across CPTAC cohorts with their RNA expression rank  
120 across corresponding cohorts from The Cancer Genome Atlas (TCGA).<sup>22</sup> Specifically, the  
121 median log<sub>2</sub>(iBAQ) before and after normalization, the median log<sub>2</sub>(riBAQ) of CPTAC, and the  
122 median log<sub>2</sub>(TPM) of TCGA were calculated for those proteins within each indication. Next,  
123 indications were ranked by the median log<sub>2</sub>(iBAQ) before and after normalization, the median  
124 log<sub>2</sub>(riBAQ) and the median log<sub>2</sub>(TPM) of CPTAC and TCGA, respectively. A weighted rank  
125 correlation approach<sup>23</sup> was used to measure rank agreement between each comparison, where the  
126 correlation between the protein and RNA expression in CPTAC was used as the weight.

#### 127 **Sensitivity Analysis**

128 To further delineate the impact of imputation and normalization separately, we performed the  
129 following sensitivity analyses.

130 **Comparison of normalization with normalization/imputation impact on analyses for**  
131 **understanding transcriptional differences between tumor and normal tissues**

132 We performed differential expression analysis between tumor and normal tissues for ccRCC and  
133 LSCC, using datasets processed through normalization alone as well as normalization combined  
134 with imputation. For both indications, we conducted gene ontology enrichment analysis and  
135 gene set enrichment analysis (GSEA). Only differentially expressed proteins with fold  $\geq 2$  and  
136 FDR  $\leq 0.01$  were used for gene ontology enrichment and GSEA analysis. Gene ontology  
137 analysis was performed with enrichR<sup>18</sup> package and database GO\_Biological\_Process\_2023 was  
138 used for query. GSEA analysis was performed with fgsea<sup>19,20</sup> package, and MSigDB<sup>21</sup> hallmark  
139 50 gene sets were assessed.

140 **Subset analysis on global quantile normalization**

141 Global quantile normalization was evaluated using a subset-based approach. After imputation,  
142 normalization was applied to all 252 five-cohort subsets constructed from 10 indications.  
143 Sample-level concordance was assessed by computing Pearson correlations of normalized  
144  $\log_2(\text{iBAQ})$  values for overlapping samples across five-cohort subsets. Consistency of  
145 downstream analyses was assessed by computing  $\log_2$  tumor-normal fold changes for COAD  
146 and LSCC across all relevant five-cohort subsets and calculating pairwise correlations among  
147 fold-change vectors. Differential expression analyses were performed within each indication  
148 across the relevant five-cohort subsets. Differentially expressed proteins were defined using FDR  
149  $< 0.01$  and absolute fold change  $> 2$ , and DEP count distributions were summarized for each  
150 indication.

151

## 152 **Results**

### 153 **Evaluation of CPTAC pan-cancer reprocessed proteomics data**

154 We used the FragPipe computational platform to assess protein abundance in tumor tissue and  
155 normal tissue samples from 10 cancer indications in CPTAC. Log<sub>2</sub>-transformed protein  
156 abundance was used to identify missing value patterns and develop an imputation strategy.

157 The workflow used in this study is outlined in Fig. 1. We developed an algorithm to select  
158 robustly expressed proteins across cohorts as described in the Methods section. A total of 15,762  
159 proteins were identified in the union of all cohorts, with 8,419 commonly detected proteins  
160 across all indications. The protein identifications per cohort at 1% protein-level false discovery  
161 rate are shown in Supplementary Table 1. In addition to the 8,419 proteins identified in all  
162 cohorts, we included 1,718 proteins that, while not present in all cohorts, are considered robustly  
163 expressed in certain indications by our robustly expressed proteins selection algorithm defined in  
164 the Methods section. Finally, we selected 10,137 proteins for further analysis, and the selected  
165 proteins per cohort are shown in Supplementary Table 2. Moreover, when we assessed the total  
166 iBAQ values for each sample, excluding the filtered proteins versus including them, our  
167 algorithm showed a mere 1% reduction (see Supplementary Fig. 1).

### 168 **Evaluation of missing values pattern**

169 Our evaluation of the missing values showed both missing-at-random (MAR) and MNAR  
170 patterns within and across cohorts (Fig. 2a, 2b). As shown in Fig. 2a, the COAD cohort had the  
171 highest missing values among all cohorts. The most prominent categories of proteins with  
172 missing values are those with a missing rate of 0–25% in the MAR pattern and those with a

173 missing rate of 75% to 100% in the MNAR pattern. Interestingly, BRCA exhibited a markedly  
174 higher number of proteins with the MNAR pattern in normal tissue than in tumor tissue (Fig. 2b).  
175 Furthermore, proteins with higher missing rates tended to have lower expression levels  
176 (Supplementary Fig. 2), suggesting an MNAR pattern. Because it was critical to apply the  
177 appropriate missing value imputation method before performing any further downstream  
178 analysis, we chose a cohort hybrid algorithm to impute missing values <sup>9</sup>.

### 179 **Assessment of normalization methods**

180 Global quantile normalization ensured identical distribution across cohorts for both tumor and  
181 normal tissues, whereas smooth quantile normalization preserved variability across cohorts (Fig.  
182 3a, 3b). The PCA analysis of protein expression across pan-cancer indications revealed that  
183 cohort separation based on protein expression was more distinct than tissue types, both before  
184 and after normalization. The most noticeable tissue separation occurred following global quantile  
185 normalization (Supplementary Fig. 3 and Supplementary Fig. 4).

### 186 **Differential protein expression analysis with normalized and imputed iBAQ values**

187 To determine whether our missing data imputation and normalization strategy affected  
188 downstream analyses, we compared the fold change in DEP between tumor tissue and matched  
189 normal tissue of selected cohorts (ccRCC, COAD, LSCC, and LUAD), using non-normalized,  
190 riBAQ, global quantile normalized, and smooth quantile normalized protein iBAQ values. Our  
191 results (Pearson  $r$  values) demonstrate a strong correlation in fold change between riBAQ  
192 normalized data and non-normalized data (ccRCC,  $r = 0.99823$ ; COAD,  $r = 0.99791$ ; LUAD,  $r =$   
193  $0.99749$ ; LSCC,  $r = 0.99813$ ) (Fig. 4a), and global quantile normalized data and non-normalized  
194 data (ccRCC,  $r = 0.97262$ ; COAD,  $r = 0.97882$ ; LUAD,  $r = 0.99105$ ; LSCC,  $r = 0.99142$ ) (Fig.

195 4b). Similar results were observed when we compared smooth quantile normalized data with  
196 non-normalized data (ccRCC,  $r = 0.99863$ ; COAD,  $r = 0.99993$ ; LUAD,  $r = 0.99819$ ; LSCC,  $r =$   
197  $0.99996$ ) (Fig. 4c). The Pearson  $r$  values tended to decrease as more variability was removed  
198 from the dataset; for example, the smooth quantile normalization resulted in  $r = 0.99863$ ,  
199 whereas the global quantile normalized resulted in  $r = 0.97262$  in ccRCC. However, with global  
200 quantile normalization, the correlations were still greater than 0.97 in the four selected  
201 indications, indicating that the global quantile normalization method retained biological  
202 differences between tumor tissues and matched normal tissues within cohorts. The estimated  
203 log<sub>2</sub>-fold change, p-value, and p-values adjusted for multiple testing using Benjamini-Hochberg  
204 procedure <sup>24</sup>(FDR) are presented in Supplementary Tables S3-S6.

205 To benchmark the impact of imputation and quantile normalization, we assessed changes in  
206 DEPs identified in LUAD cohort. Applying a threshold of log<sub>2</sub>-fold change > 1 and FDR < 0.01,  
207 we compared DEPs identified using different versions of data: raw data, imputed data where  
208 only missing not at random (MNAR) values are imputed, imputed data where both MNAR and  
209 missing at random (MAR) values are imputed, imputation of MNAR followed by quantile  
210 normalization, and imputation of both MNAR and MAR followed by quantile normalization  
211 (Supplementary Table S7). Prior to normalization, imputing MAR values in addition to MNAR  
212 resulted in the identification of 23 additional DEPs. However, the power increase from imputing  
213 MAR is limited: of the 23 additional DEPs found after KNN imputation, 10 have a log<sub>2</sub>-fold  
214 change above 0.8 prior to KNN imputation is applied. As expected, more extensive data  
215 processing leads to a reduction in the overlap of DEPs shared between raw data and processed  
216 datasets, dropping from 98.1% (raw versus MAR imputation only) to 75.8% (raw versus full  
217 imputation with quantile normalization) as shown in supplementary figure 5. Among the 67

218 DEPs identified exclusively after imputation and normalization (relative to raw data), 52 (78%)  
219 still showed a log<sub>2</sub>-fold change greater than 0.8 in the raw dataset (Supplementary Fig. 5).  
220 However, we observed that the imputed dataset yielded four DEPs with notably large log<sub>2</sub>-fold  
221 change values (ranging from -3.7 to -7.7). Therefore, external validation is necessary to  
222 distinguish whether these observations reflect true biological signals or artifacts of the  
223 imputation process. While the overarching goal of cross-indication analysis necessitates retaining  
224 all proteins, it is particularly important to interpret results for cohorts with completely missing  
225 protein abundances with caution.

#### 226 **Comparison of the indication ranks of selected proteins between CPTAC and TCGA**

227 We identified proteins whose protein and RNA expression were highly correlated ( $r > 0.5$ )  
228 across CPTAC cohorts: ERAP2, CA9, GSTM3, MX1, and STAT1. We then compared their  
229 median protein expression rank across CPTAC cohorts with their median RNA expression rank  
230 across corresponding TCGA cohorts. A weighted rank correlation approach<sup>23</sup> was used to  
231 measure rank agreement between each comparison (e.g., denoted as “v” in Supplementary Fig.  
232 5). OV and glioblastoma multiforme GBM were excluded from the weighted rank correlation  
233 calculation because they had a small number of proteins with  $r > 0.5$ , and those proteins did not  
234 overlap with those selected from other indications. Global quantile normalization had a higher  
235 rank correlation (weighted rank correlation of 0.597–0.931) than smooth quantile normalization  
236 (weighted rank correlation of 0.168–0.76) or no normalization (weighted rank correlation of  
237 0.168–0.76) (Supplementary Fig. 8).

#### 238 **Comparison of normalization with normalization/imputation impact on analyses for** 239 **understanding transcriptional differences between tumor and normal tissues**

240 Differential expression analysis results between tumor and normal tissues for both ccRCC and  
241 LSCC, using datasets processed through normalization alone as well as normalization combined  
242 with imputation. Results from the comprehensive abundance analyses covering both  
243 normalization +imputation and normalization-only approaches are available in Supplementary  
244 Table 8-11.

245 For ccRCC and LUAD indications, gene ontology enrichment analysis identified highly similar  
246 top biological process terms from the normalized+imputed datasets and the normalization-only  
247 datasets (Supplementary Fig. 6 and Supplementary Fig. 7). Likewise, GSEA highlighted  
248 comparable MSigDB signaling pathways as enriched in both data processing approaches  
249 (Supplementary Fig. 6 and Supplementary Fig. 7). At an FDR threshold of 0.01, there are 223  
250 overlapping biological processes between the normalized/imputed and normalization-only  
251 datasets for ccRCC. For LSCC, 161 biological processes are overlapped at the same threshold  
252 (Supplementary Fig. 6 and Supplementary Fig. 7). These results suggest that imputation has  
253 minimal impact on the identification of relevant biological processes, indicating that our analysis  
254 workflow preserves the underlying biology while enabling comparisons across indications.

### 255 **Subset analysis on global quantile normalization**

256 To further assess the potential impact of global quantile normalization on biological variation,  
257 we applied the normalization procedure to all 252 possible 5-cohort subsets derived from the 10  
258 available indications after imputation. The consistency of normalized data resulting from  
259 different subsets was examined in two ways.

260 First, we assessed the correlations between normalized  $\log_2(\text{iBAQ})$  values for overlapping  
261 samples present in different pairs of 5-cohort subsets. For each pair of 5-cohort subsets, we

262 identified samples that appeared in both subsets and compared their  $\log_2(\text{iBAQ})$  values. Due to  
263 computational constraints, we did not evaluate all possible pairs of subsets; however, a  
264 representative analysis using samples from cohorts (BRCA, ccRCC, COAD, GBM, HNSCC),  
265 paired with all possible 5-cohort subsets, showed a minimum Pearson correlation coefficient of  
266 0.9987. The high level of agreement indicates substantial stability of the normalization process  
267 and suggests that downstream analyses based on this normalized data will be highly consistent.

268 Second, we evaluated the consistency of downstream analyses by comparing fold changes  
269 between tumor and matched normal tissues in COAD and LSCC across all relevant 5-cohort  
270 subsets. There are 126 possible of 5-cohort subsets that include samples from COAD (or LSCC),  
271 and each subset generates a set of fold change estimates. We calculated the pairwise correlations  
272 among these  $\log_2$ -transformed fold changes and visualized using heatmaps (Supplementary Fig.  
273 9). As expected, the heatmaps demonstrated a consistently high level of agreement in fold  
274 change estimates across subsets, further supporting the consistency of the downstream analyses.

275 Additionally, we have conducted differential expression analysis between tumor and normal  
276 tissues using 5-cohort subsets. Specifically, for each indication, we considered all relevant 5-  
277 cohort subset combinations (126 subsets per indication), allowing us to quantify the number of  
278 DEPs identified in each case. DEPs were defined consistently using the same criteria  
279 ( $\text{FDR} < 0.01$  and absolute fold change  $> 2$ ). We have summarized the distribution of DEP counts  
280 for each indication in supplementary figure 10. For most indications, the number of DEPs  
281 remains relatively stable across different 5-cohort subsets, suggesting a consistent effect of data  
282 processing strategy. We observed notably greater variability in ccRCC, LSCC, and UCEC, with  
283 interquartile ranges of 61, 71, and 101, respectively. This increased variability may be

284 attributable to quantile normalization moderately compressing or enlarging biological effects in  
285 these indications, subsequently affecting the number of DEPs detected.

286 Taken together, these results suggest that global quantile normalization can substantially  
287 minimize technical variation while maintaining underlying biological differences, indicating that  
288 the principal biological conclusions are generally robust to the specific composition of cohorts  
289 included in the normalization.

## 290 **Discussion**

291 CPTAC recently generated harmonized genomic, transcriptomic, proteomic, and clinical data for  
292 over 1,000 tumors in 10 cohorts to facilitate pan-cancer discovery research.<sup>2</sup> To use these data  
293 for prioritizing cancer surface antigens and selecting indications in the discovery and  
294 development of drug targets, protein expression levels often need to be compared and ranked  
295 across cohorts.<sup>25,26</sup> Efforts to do so are hindered, however, by non-uniform missing data and  
296 varying protein expression distribution patterns across tumor types. To evaluate various  
297 strategies for missing data handling and normalization for the generation of a normalized pan-  
298 cancer protein expression dataset, we built a computational workflow that incorporates the  
299 selection of robustly expressed proteins, cohort hybrid imputation, and quantile normalization.<sup>27</sup>  
300 We then evaluated our strategy by comparing indication ranking between CPTAC and TCGA,  
301 using a set of proteins with high correlation between protein and RNA expression (Fig. 1). To  
302 our knowledge, this work represents the first cross-cohort normalized CPTAC pan-cancer data  
303 containing estimated absolute protein expression quantification derived from mass spectrometry-  
304 based proteomic data.

305 There are considerable challenges in harmonizing proteomics data across the 10 CPTAC cohorts  
306 because they have been generated over multiple years and analyzed by various laboratories using  
307 different processing pipelines. To avoid potential confounding factors, it was necessary to  
308 reanalyze TMT global proteomics data with the same pipeline. We first estimated the absolute  
309 protein abundance in reprocessed CPTAC pan-cancer data as described previously.<sup>3</sup> In addition,  
310 to create a unified dataset, it was necessary to define a set of proteins that maximizes the  
311 identification of proteins across all 10 cohorts with high quality while minimizing the occurrence  
312 of missing values. We designed a protein selection algorithm to effectively balance these  
313 requirements. Our robustly expressed proteins selection algorithm defined 10,137 proteins,  
314 which indicated that more than half of the human proteome (assuming 20,000 protein  
315 products)<sup>28,29</sup> is robustly expressed.

316 Proteomics data often contain missing values due to a variety of technical factors, such as protein  
317 abundance below the instrument limit of detection, poor ionization efficiency, and low signal-to-  
318 noise ratio.<sup>30-32</sup> When merging the 10 cohorts, we faced similar issues. Although many  
319 algorithms have been developed for missing value imputation,<sup>33</sup> it is critical to select the ones  
320 that are best suited to the specific characteristics of the data and downstream analysis needs. An  
321 analysis of the missing values in the CPTAC cohort uncovered patterns indicative of both MAR  
322 and MNAR. Based on these observations, a cohort-hybrid imputation approach was applied. Our  
323 study introduces a general framework to manage missing data when comparing protein  
324 expression across cohorts. For other specific analyses, one can adopt this model, integrating  
325 multiple imputation methods to measure the uncertainty linked to imputation, thereby obtaining  
326 more accurate estimations.

327 Due to the lack of a definitive list of proteins with known expression levels across different  
328 tumor types, we chose to use proteins that exhibit a strong correlation between their RNA and  
329 protein expression levels. We then assessed our approach by comparing the rank order of these  
330 proteins' expression in our normalized dataset to their RNA expression rank in TCGA as an  
331 indirect method of evaluation. Furthermore, we implemented a weighted rank correlation method  
332 that accounts for the correlation between protein and RNA expression levels. The results of our  
333 validation method showed a significantly enhanced rank correlation with the application of  
334 global cohort normalization as opposed to using a dataset without normalization (see  
335 Supplementary Fig. 8). Further, when comparing the fold change of tumor versus matched  
336 normal tissues using the normalized dataset versus the original dataset, there was significant  
337 agreement (Fig. 4), indicating that our normalization method did not inadvertently eliminate  
338 biological variation.

339 Our assessment indicates that a combination of cohort hybrid imputation and global quantile  
340 normalization is a reasonable approach to generate a normalized CPTAC pan-cancer protein  
341 dataset that could be leveraged to compare protein expression across different cancer types.  
342 Nonetheless, the limited number of proteins and indications used for validation in our study  
343 represents an important limitation. While our approach was tested using selected highly  
344 correlated proteins and rank alignment with TCGA RNA-Seq data, we recommend further  
345 validation using a larger number of proteins with expression levels measured across cohorts with  
346 other orthogonal experimental methods such as targeted proteomics to enhance the robustness of  
347 our approach.

348 Quantile normalization assumes that, across runs or samples, the majority of proteins are drawn  
349 from the same underlying distribution—i.e., most proteins are unchanged and present at

350 comparable concentrations—so observed global differences primarily reflect technical rather  
351 than biological effects. In a pan-cancer, multi-tissue, multi-cohort TMT design, this assumption  
352 can be challenged because tumor types and adjacent normal tissues often exhibit genuine global  
353 proteomic shifts<sup>34,35</sup>. Our preprocessing is not intended to bypass this assumption; it mitigates  
354 known technical heterogeneity before cross-cohort harmonization.

355 We have established a computational workflow for cross-indication comparison from proteomics  
356 data, thereby providing a unique data resource to interrogate protein expression across different  
357 cancer types. Despite our efforts to normalize the CPTAC pan-cancer protein dataset across  
358 different cancer cohorts, batch effects may still persist due to upstream technical variations, such  
359 as differences in laboratory practices and platforms<sup>36</sup>. Our results demonstrate that missing data  
360 imputation and normalization strategies do not affect downstream analyses. The methodology  
361 and the data sources used in this study can serve as valuable resources for cancer research.  
362 Despite evaluating multiple normalization strategies (including global quantile normalization,  
363 smooth quantile normalization, and riBAQ), we did not implement median normalization in the  
364 primary analysis. We recognize that this more comprehensive normalization can result in greater  
365 alteration to the original data structure, which may affect subtle biological signals. To address  
366 this concern, we conducted downstream sensitivity analyses demonstrating that key biological  
367 patterns and signals remain preserved following quantile normalization. A comprehensive  
368 benchmarking of all potential normalization strategies, including alternative approaches beyond  
369 those implemented here, remains an important direction for future research and will be critical  
370 for optimizing data integration and downstream analyses in large proteomics studies.

371

372 **References**

- 373 1 Lindgren, C. M. *et al.* Simplified and Unified Access to Cancer Proteogenomic Data.  
374 *Journal of Proteome Research* **20**, 1902-1910, doi:10.1021/acs.jproteome.0c00919  
375 (2021).
- 376 2 Li, Y. *et al.* Proteogenomic data and resources for pan-cancer analysis. *Cancer Cell* **41**,  
377 1397-1406, doi:10.1016/j.ccell.2023.06.009 (2023).
- 378 3 Wang, J. *et al.* Pan-Cancer Proteomics Analysis to Identify Tumor-Enriched and Highly  
379 Expressed Cell Surface Antigens as Potential Targets for Cancer Therapeutics. *Mol Cell*  
380 *Proteomics* **22**, 100626, doi:10.1016/j.mcpro.2023.100626 (2023).
- 381 4 Wang, J. *et al.* Abstract LB012: Evaluating computational approaches for CPTAC pan-  
382 cancer cross-cohort protein expression comparison. *Cancer Research* **84**, LB012-LB012,  
383 doi:10.1158/1538-7445.Am2024-lb012 (2024).
- 384 5 Moreno, P. *et al.* Expression Atlas update: gene and protein expression in multiple  
385 species. *Nucleic Acids Res* **50**, D129-d140, doi:10.1093/nar/gkab1030 (2022).
- 386 6 Schmidt, T. *et al.* ProteomicsDB. *Nucleic Acids Res* **46**, D1271-d1281,  
387 doi:10.1093/nar/gkx1029 (2018).
- 388 7 Yu, F. *et al.* Analysis of DIA proteomics data using MSFragger-DIA and FragPipe  
389 computational platform. *Nature Communications* **14**, 4154, doi:10.1038/s41467-023-  
390 39869-5 (2023).
- 391 8 Kong, A. T., Lprevost, F. V., Avtonomov, D. M., Mellacheruvu, D. & Nesvizhskii, A. I.  
392 MSFragger: ultrafast and comprehensive peptide identification in mass spectrometry-  
393 based proteomics. *Nat. Methods* **14**, 513-520, doi:10.1038/nmeth.4256 (2017).

- 394 9 Lazar, C., Burger, T. & Wieczorek, S. imputeLCMD: a collection of methods for left-  
395 censored missing data imputation. (2022). <[https://cran.r-](https://cran.r-project.org/web/packages/imputeLCMD/imputeLCMD.pdf)  
396 [project.org/web/packages/imputeLCMD/imputeLCMD.pdf](https://cran.r-project.org/web/packages/imputeLCMD/imputeLCMD.pdf)>.
- 397 10 Gao, Q. *et al.* Integrated Proteogenomic Characterization of HBV-Related Hepatocellular  
398 Carcinoma. *Cell* **179**, 561-577.e522, doi:10.1016/j.cell.2019.08.052 (2019).
- 399 11 Palstrøm, N. B., Matthiesen, R. & Beck, H. C. Data Imputation in Merged Isobaric  
400 Labeling-Based Relative Quantification Datasets. *Methods Mol Biol* **2051**, 297-308,  
401 doi:10.1007/978-1-4939-9744-2\_13 (2020).
- 402 12 Bolstad, B. M., Irizarry, R. A., Astrand, M. & Speed, T. P. A comparison of  
403 normalization methods for high density oligonucleotide array data based on variance and  
404 bias. *Bioinformatics* **19**, 185-193, doi:10.1093/bioinformatics/19.2.185 (2003).
- 405 13 Hicks, S. C. *et al.* Smooth quantile normalization. *Biostatistics* **19**, 185-198,  
406 doi:10.1093/biostatistics/kxx028 (2018).
- 407 14 Shin, J. B. *et al.* Molecular architecture of the chick vestibular hair bundle. *Nat. Neurosci.*  
408 **16**, 365-374, doi:10.1038/nn.3312 (2013).
- 409 15 Milo, R. What is the total number of protein molecules per cell volume? A call to rethink  
410 some published values. *Bioessays* **35**, 1050-1055, doi:10.1002/bies.201300066 (2013).
- 411 16 Smyth, G. K. Linear models and empirical Bayes methods for assessing differential  
412 expression in microarray experiments. *Stat. Appl. Genet. Mol. Biol.* **3**,  
413 doi:doi:10.2202/1544-6115.1027 (2004).
- 414 17 Zhang, X. *et al.* Proteome-wide identification of ubiquitin interactions using UbIA-MS.  
415 *Nat. Protoc.* **13**, 530-550, doi:10.1038/nprot.2017.147 (2018).

- 416 18 Kuleshov, M. V. *et al.* Enrichr: a comprehensive gene set enrichment analysis web server  
417 2016 update. *Nucleic Acids Res* **44**, W90-97, doi:10.1093/nar/gkw377 (2016).
- 418 19 Korotkevich, G. *et al.* Fast gene set enrichment analysis. *bioRxiv*, 060012,  
419 doi:10.1101/060012 (2021).
- 420 20 Subramanian, A. *et al.* Gene set enrichment analysis: a knowledge-based approach for  
421 interpreting genome-wide expression profiles. *Proc Natl Acad Sci U S A* **102**, 15545-  
422 15550, doi:10.1073/pnas.0506580102 (2005).
- 423 21 Liberzon, A. *et al.* Molecular signatures database (MSigDB) 3.0. *Bioinformatics* **27**,  
424 1739-1740, doi:10.1093/bioinformatics/btr260 (2011).
- 425 22 Chang, K. *et al.* The Cancer Genome Atlas Pan-Cancer analysis project. *Nat. Genet.* **45**,  
426 1113-1120, doi:10.1038/ng.2764 (2013).
- 427 23 Plaia, A., Buscemi, S. & Sciandra, M. Consensus among preference rankings: a new  
428 weighted correlation coefficient for linear and weak orderings. *Adv. Data Anal. Classif.*  
429 **15**, 1015-1037, doi:10.1007/s11634-021-00442-x (2021).
- 430 24 Benjamini, Y. & Hochberg, Y. Controlling the False Discovery Rate: A Practical and  
431 Powerful Approach to Multiple Testing. *Journal of the Royal Statistical Society: Series B*  
432 *(Methodological)* **57**, 289-300, doi:<https://doi.org/10.1111/j.2517-6161.1995.tb02031.x>  
433 (1995).
- 434 25 Jarnuczak, A. F. *et al.* An integrated landscape of protein expression in human cancer.  
435 *Sci. Data* **8**, 115, doi:10.1038/s41597-021-00890-2 (2021).
- 436 26 Kosti, I., Jain, N., Aran, D., Butte, A. J. & Sirota, M. Cross-tissue analysis of gene and  
437 protein expression in normal and cancer tissues. *Sci. Rep.* **6**, 24799,  
438 doi:10.1038/srep24799 (2016).

- 439 27 Gong, T. Q. *et al.* Proteome-centric cross-omics characterization and integrated network  
440 analyses of triple-negative breast cancer. *Cell Rep.* **38**, 110460,  
441 doi:10.1016/j.celrep.2022.110460 (2022).
- 442 28 Aebersold, R. *et al.* How many human proteoforms are there? *Nat. Chem. Biol.* **14**, 206-  
443 214, doi:10.1038/nchembio.2576 (2018).
- 444 29 Ponomarenko, E. A. *et al.* The size of the human proteome: the width and depth. *Int. J.*  
445 *Anal. Chem.* **2016**, 7436849, doi:10.1155/2016/7436849 (2016).
- 446 30 McGurk, K. A. *et al.* The use of missing values in proteomic data-independent  
447 acquisition mass spectrometry to enable disease activity discrimination. *Bioinformatics*  
448 **36**, 2217-2223, doi:10.1093/bioinformatics/btz898 (2020).
- 449 31 Jin, L. *et al.* A comparative study of evaluating missing value imputation methods in  
450 label-free proteomics. *Sci. Rep.* **11**, 1760, doi:10.1038/s41598-021-81279-4 (2021).
- 451 32 Gardner, M. L. & Freitas, M. A. Multiple imputation approaches applied to the missing  
452 value problem in bottom-up proteomics. *Int. J. Mol. Sci.* **22**, 9650,  
453 doi:10.3390/ijms22179650 (2021).
- 454 33 Lazar, C., Gatto, L., Ferro, M., Bruley, C. & Burger, T. Accounting for the multiple  
455 natures of missing values in label-free quantitative proteomics data sets to compare  
456 imputation strategies. *J. Proteome Res.* **15**, 1116-1125,  
457 doi:10.1021/acs.jproteome.5b00981 (2016).
- 458 34 Arend, L. *et al.* Systematic evaluation of normalization approaches in tandem mass tag  
459 and label-free protein quantification data using PRONE. *Briefings in Bioinformatics* **26**,  
460 doi:10.1093/bib/bbaf201 (2025).

- 461 35 Maes, E. *et al.* Determination of Variation Parameters as a Crucial Step in Designing  
462 TMT-Based Clinical Proteomics Experiments. *PLOS ONE* **10**, e0120115,  
463 doi:10.1371/journal.pone.0120115 (2015).
- 464 36 Hui, H. W. H., Kong, W. & Goh, W. W. B. Thinking points for effective batch correction  
465 on biomedical data. *Briefings in Bioinformatics* **25**, doi:10.1093/bib/bbae515 (2024).
- 466

ARTICLE IN PRESS

**467 Acknowledgments**

468 We gratefully acknowledge the CPTAC for providing open-source proteomics data and Deborah  
469 Shuman of AstraZeneca for editing the manuscript and formatting the figures. We also  
470 acknowledge the presentation of a portion of this work at the AACR Annual Meeting 2024 (  
471 April 5-10, 2024, San Diego; [https://aacrjournals.org/cancerres/issue/84/7\\_Supplement](https://aacrjournals.org/cancerres/issue/84/7_Supplement)).

**472 Author contributions**

473 JW and WZ conceived and designed the study. JW, XT, YW, BP, JB, EH, and WZ analyzed and  
474 interpreted the data. All authors contributed to the writing, review, and/or revision of the  
475 manuscript, have approved the final version of the manuscript, and agree to be accountable for  
476 all aspects of the work.

**477 Availability of data and materials**

478 The search output and reports from the FragPipe, Protein-level data without normalization and  
479 without imputation, Protein-level data after imputation (no normalization) and protein-level data  
480 for each normalization method after imputation of protein abundance estimation for the 10  
481 CPTAC indications used in this work and scripts can be downloaded from  
482 <https://doi.org/10.5281/zenodo.17203216>.

**483 Competing interests**

484 All authors are employees of AstraZeneca and may have stock ownership, options, or interests in  
485 the company.

**486 Funding**

487 This study was funded by AstraZeneca.

488

## 489 **Figure legends**

490 **Figure 1.** Workflow for a computational approach that enables the comparison of protein  
491 expression across CPTAC cohorts. Tissue samples from 10 cancer indications in CPTAC were  
492 processed through FragPipe. Protein abundance data were analyzed by iBAQ to derive estimated  
493 absolute protein abundance. A cohort hybrid imputation strategy was developed to address  
494 missing values. Quantile normalization was applied to iBAQ values. Imputed and normalized  
495 iBAQ values were used for validation with differential protein analysis, ranking indications by  
496 protein expression and comparing these ranks with those derived from TCGA RNA-Seq data.

497 **Figure 2.** CPTAC missing values pattern assessment in **a** tumor and **b** normal tissue. “Cohort”  
498 was defined as the unique combination of tissue type and indication. “Missing rate” was defined  
499 as the percentage of samples with missing values in each cohort.

500 **Figure 3.** Quantile normalization. **a** Quantile normalization applied to derived iBAQ values after  
501 imputation in **a** tumor and **b** normal tissue. Two quantile normalization methods, global and  
502 smooth quantile normalization, were used.

503 **Figure 4.** Correlation between fold change from differential expression analysis. **a** riBAQ  
504 normalization and no normalization. **b** Global quantile normalization and no normalization. **c**  
505 Smooth quantile normalization and no normalization.

506

507

Figure 1

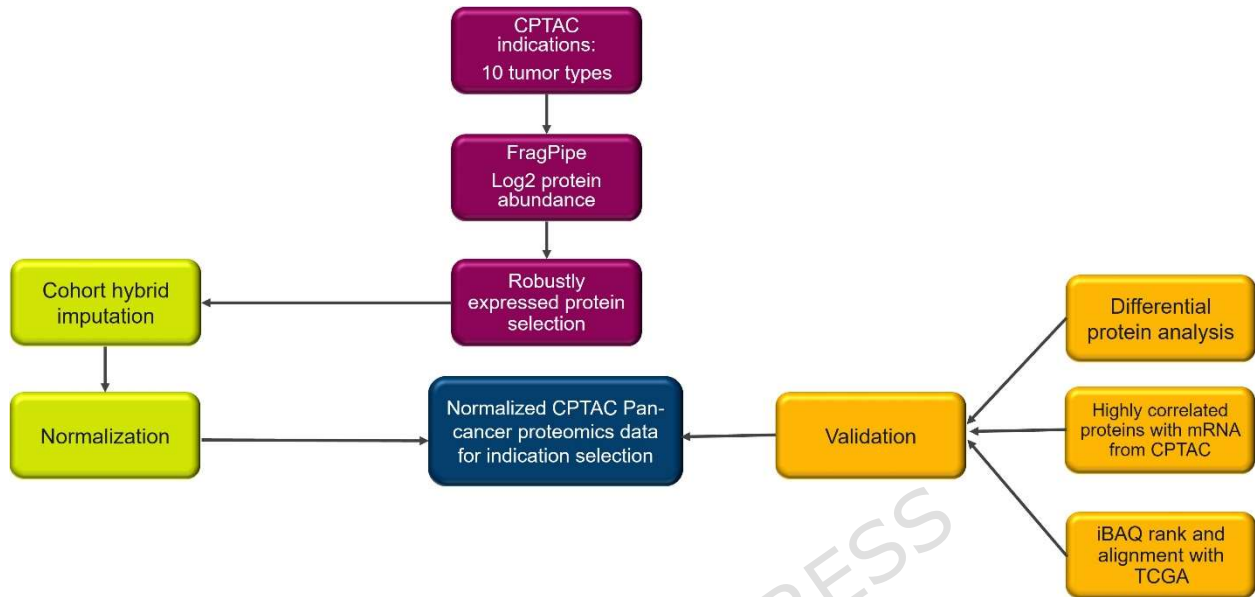


Figure 2

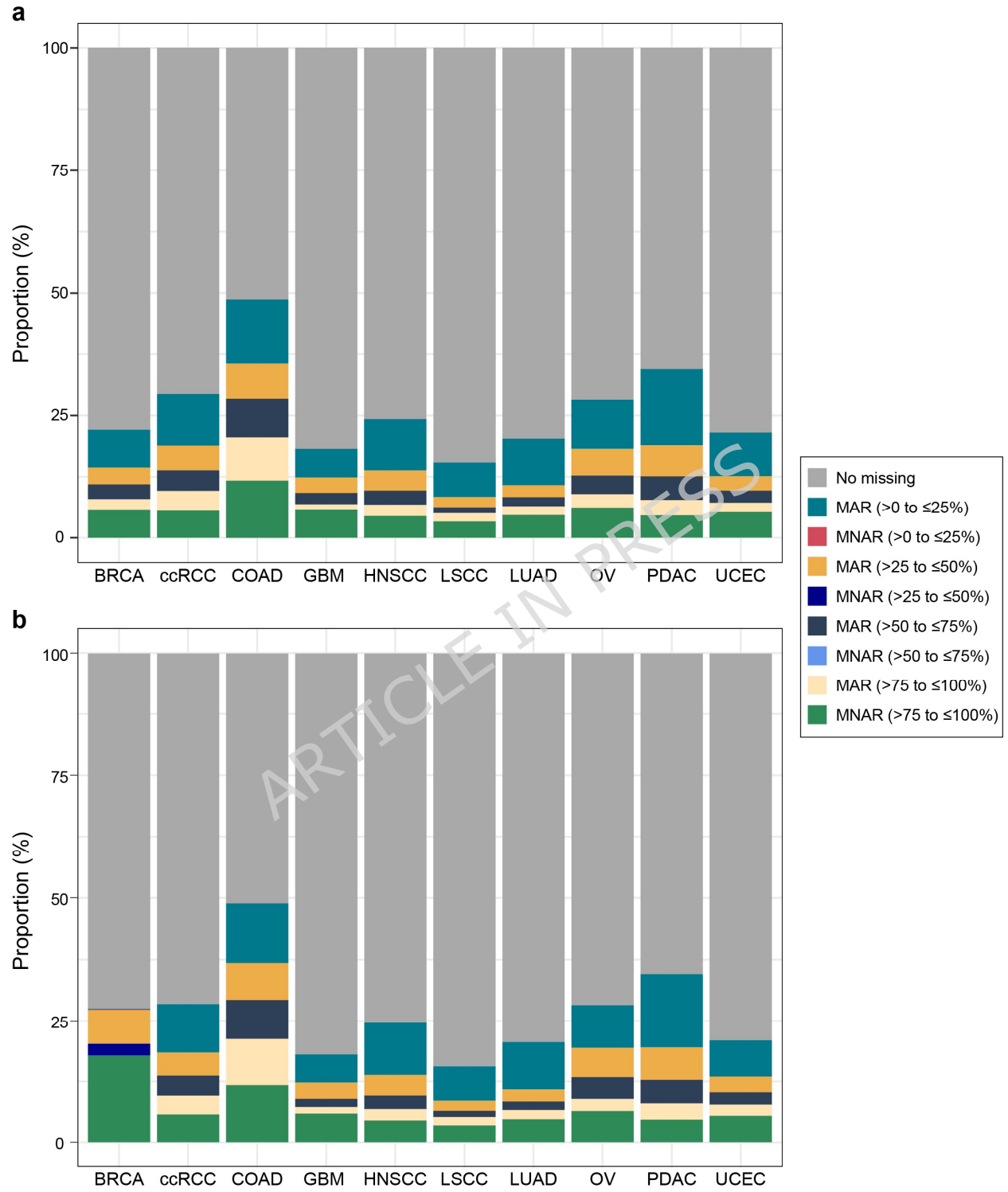
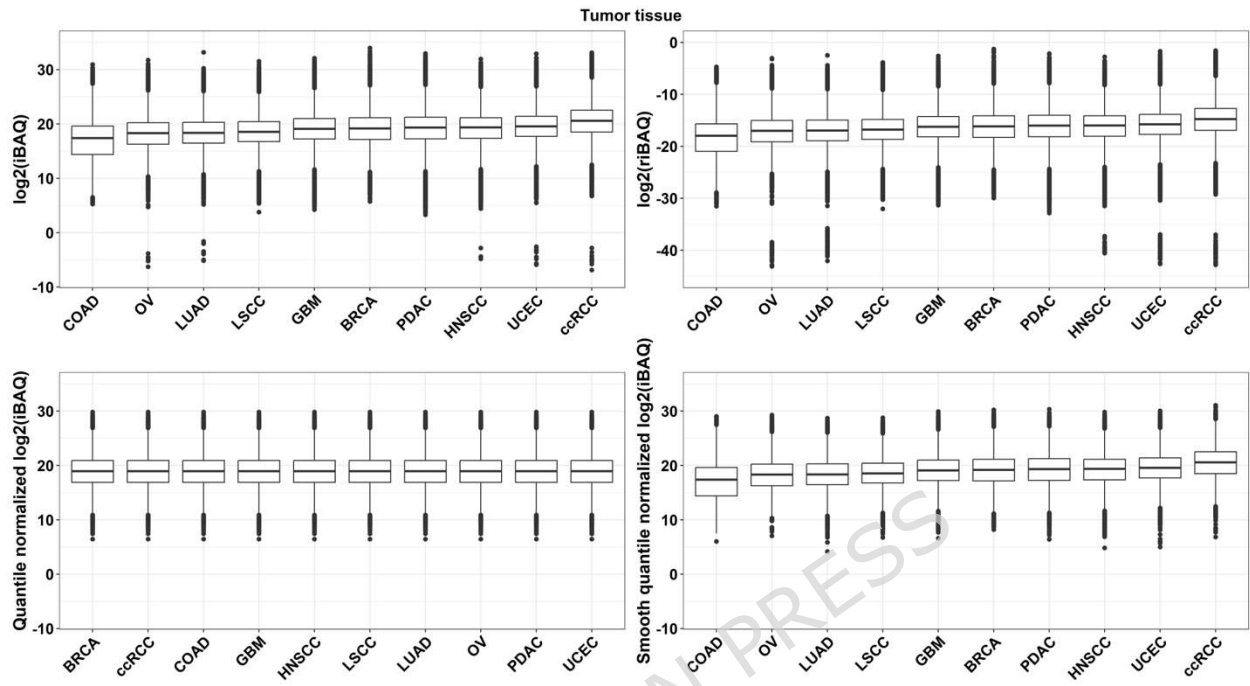


Figure 3

a



b

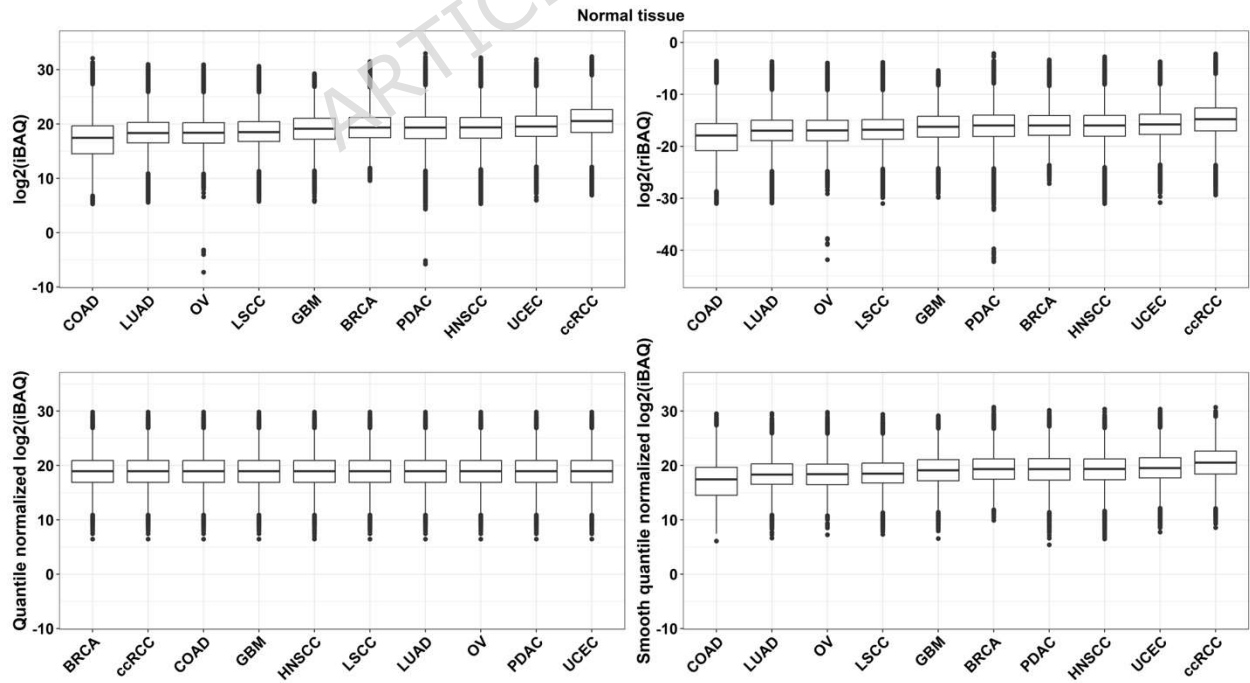


Figure 4

