

## Mutation Impact on mRNA Versus Protein Expression across Human Cancers --Manuscript Draft--

<b>Manuscript Number:</b>	GIGA-D-24-00168R1	
<b>Full Title:</b>	Mutation Impact on mRNA Versus Protein Expression across Human Cancers	
<b>Article Type:</b>	Research	
<b>Funding Information:</b>	National Institute of General Medical Sciences (R35GM138113)	Dr Kuan-lin Huang
	American Cancer Society (RSG-22-115-01-DMC)	Dr Kuan-lin Huang
<b>Abstract:</b>	<p>Cancer mutations are often assumed to alter proteins, thus promoting tumorigenesis. However, how mutations affect protein expression has rarely been systematically investigated. We conduct a comprehensive analysis of mutation impacts on mRNA- and protein-level expressions of 953 cancer cases with paired genomics and global proteomic profiling across six cancer types. Protein-level impacts are validated for 47.2% of the somatic expression quantitative trait loci (seQTLs), including mutations from likely “long-tail” driver genes. Devising a statistical pipeline for identifying somatic protein-specific QTLs (spsQTLs), we reveal several gene mutations, including NF1 and MAP2K4 truncations and TP53 missenses showing disproportional influence on protein abundance not readily explained by transcriptomics. Cross-validating with data from massively parallel assays of variant effects (MAVE), TP53 missenses associated with high tumor TP53 proteins were experimentally confirmed as functional. Our study demonstrates the importance of considering protein-level expression to validate mutation impacts and identify functional genes and mutations.</p>	
<b>Corresponding Author:</b>	Kuan-lin Huang, PhD Icahn School of Medicine at Mount Sinai New York, NY UNITED STATES	
<b>Corresponding Author Secondary Information:</b>		
<b>Corresponding Author's Institution:</b>	Icahn School of Medicine at Mount Sinai	
<b>Corresponding Author's Secondary Institution:</b>		
<b>First Author:</b>	Yuqi Liu	
<b>First Author Secondary Information:</b>		
<b>Order of Authors:</b>	Yuqi Liu	
	Abdulkadir Elmas	
	Kuan-lin Huang, PhD	
<b>Order of Authors Secondary Information:</b>		
<b>Response to Reviewers:</b>	<p>Authors: We would like to express our gratitude to the editor and the reviewers for the valuable feedback on our manuscript titled “Mutation Impact on mRNA Versus Protein Expression across Human Cancers” (GIGA-D-24-00168). We have carefully considered the comments, particularly regarding the correlation between mRNA and protein expression, and have conducted additional analyses/edits to address each of the concerns listed by reviewers. We are pleased to submit a revised version of the manuscript for your consideration. Below is a detailed response to all reviewers’ comments:</p> <p>Reviewer #1: Despite the fact that it is already well known that proteomics is important and provides a unique angle to studying cancer, this paper contributes to such knowledge from an interesting angle with the use of published data. The paper can benefit from having further descriptions on the metrics used to measure performance,</p>	

and should discuss more thoroughly alternative metrics and shortcomings of the current ones. Figures should be better prepared (e.g. Figure 1 can be enlarged or table extracted; Figure 3 Legend is truncated; )

Authors: We thank Reviewer #1 for applauding our novel approach and the feedback. We have expanded the Methods section to provide a more comprehensive description of the statistical metrics used, "spsQTL identification

We combined two complementary statistical methods to identify spsQTLs. In the first method adopted from Battle et al.4, we compared the following two nested linear models using likelihood ratio test (LRT) with the "anova" function in R:

$$p = \mu + [\beta]_0 g + [\beta]_1 r$$

$$p = \mu + [\beta]_2 r$$

where g is the genotype, r represents RNA level, and p is the protein level. By comparing these models using LRT and filtering results with an FDR less than 0.05, we identified candidate spsQTLs where the genotype (mutation) has a disproportionate impact on protein abundance independent of mRNA expression.

In the second method adopted from Mirauta et al.22, we selected QTLs where the spQTL FDR was less than 0.05 but the corresponding seQTL FDR was greater than 0.05 as candidate spsQTLs, to specifically identify mutations that affect protein levels without influencing mRNA. We then overlapped these two lists of candidate spsQTLs obtained from two complementary methods to identify the final list of spsQTLs for downstream analyses."

We also added more discussions of alternative approaches and the limitations of our current methods in Discussion, "This study has several limitations... . Fourth, our regression models assumes a linear relationship between mutations (one gene at a time), confounders, and expression, which may not capture more complex, nonlinear effects of mutations on multiple mRNA or protein expression. Future studies could explore non-linear regression models or neural network approaches to better account for these effects. Fifth, we employed two complementary methods to confidently identify spsQTLs that represent true protein-specific regulatory events. However, the reliance on FDR thresholds could still limit the detection of spsQTLs with subtle effects. Alternative approaches, such as Bayesian models that account for prior biological knowledge or hierarchical modeling, could be considered in future analyses to improve the specificity of spsQTL detection. Additionally, while our method focuses on cis-acting mutations, potential trans-acting effects could be missed, a limitation that should be explored in larger datasets or by incorporating network-based analyses."

We also have revised the figures as suggested. Figure 1 has been enlarged for clarity, and the legend for Figure 3 has been corrected.

Reviewer #2: The manuscript "Mutation Impact on mRNA Versus Protein Expression across Human Cancers" investigates how somatic mutations affect mRNA and protein expression using data from 953 cancer cases across six types. The study identifies that 47.2% of mutations impacting mRNA levels (seQTLs) also affect protein levels, validating their broader impact. A novel statistical method uncovers 83 protein-specific QTLs (spsQTLs), primarily truncating mutations, significantly affecting protein abundance. Functional validation confirms TP53 missense mutations with high protein levels are functional. However, my main concern is the relationship between mRNA expression and protein expression. The low correlation between these two levels may undermine the analysis, suggesting different regulatory mechanisms. If low correlation is observed, the overlap between seQTL and spQTL may lack biological significance. Also, truncating mutations reducing protein expression seems straightforward, but this does not fully address the complex regulation mechanisms. Therefore, I suggest that the authors first compare the correlation between mRNA and protein expression and select cancer types that show high correlation for subsequent analyses. This approach would provide a more robust biological foundation for the study.

Authors: We greatly appreciate Reviewer #2's insightful comments on the low

	<p>correlation between mRNA and protein expression and their suggestion to focus on cancer types with higher correlation for further analyses. We like to highlight that the low/moderate mRNA-protein correlation is one of the main motivations for our analyses, whereby mutations found to have mRNA effects (more known) may differ from those showing protein expression impacts (less studied). Genomics or eQTL studies in the field often neglect these potential discrepancies in their assumption.</p> <p>The added analyses and discussion are added to the main text,  “One possible source of spsQTLs is the imperfect correlation between mRNA and protein expression in the affected genes. Additional statistical analyses revealed that this mRNA-protein correlations range widely across genes and cancer types (Figure S5). While genes harboring spsQTLs have lower mRNA-protein correlations in general than genes with concordant eQTL and pQTL, this is not the case for several discordant genes, including MAP2K4 in BRCA and PBRM1 in CCRCC (Table S7). Based on the number of mutations and genes identified, CRC and UCEC reached statistically significant differences between concordant and all other expressed genes (Wilcoxon rank-sum tests, <math>p = 0.0056</math> and <math>p = 0.022</math>, respectively); in CRC, mRNA-protein correlations also showed significant differences between discordant and all other expressed genes (<math>p = 0.013</math> and <math>p = 0.29</math>, respectively); other cancer types likely did not reach statistical significance likely due to sufficient mutations identified. The imperfect correspondence between gene mRNA-protein correlations and mutation impacts further stresses the need to analyze and consider protein-specific impacts of mutations. Table S7 provides complete mRNA-protein correlation data for all concordant/discordant eQTL/pQTLs in their respective cancer type for in-depth examination.”</p> <p>As the reviewer also pointed out, truncating mutations that reduce protein expression (likely through NMD) seem straightforward but may not fully capture complex regulatory mechanisms. To clarify this, we had added to our discussion other potential post-transcriptional processes, including the role of translation efficiency and context-specific regulatory factors, that may explain the observed discordant effects between mRNA and protein levels,  “This study has several limitations. First, our findings do not distinguish between several potential mechanisms that could lead to discordant effects of mutations on gene and protein expression. One possibility is that the mutation affects the efficiency of translation, leading to changes in protein levels that are not reflected in mRNA levels. For example, accumulating evidence in recent years suggests that NMD is closely tied to the termination of translation<sup>23</sup>, which may explain instances where some truncations afford much stronger associations with protein levels in our findings. But, in many cases, the mechanisms of how mutations may affect protein abundance may be context- and gene-specific and remain to be elucidated. For example, certain mutations may influence the binding of RNA binding proteins and the efficiency of translation, whereas others may alter post-translational modifications, such as phosphorylation or ubiquitination, which can impact protein stability or degradation without affecting transcription or translation rates.”</p>
<b>Additional Information:</b>	
<b>Question</b>	<b>Response</b>
Are you submitting this manuscript to a special series or article collection?	No
<b>Experimental design and statistics</b>	Yes
Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our <a href="#">Minimum Standards Reporting Checklist</a> . Information essential to interpreting the data presented should be made available in the figure legends.	

<p>Have you included all the information requested in your manuscript?</p>	
<p><b>Resources</b></p> <p>A description of all resources used, including antibodies, cell lines, animals and software tools, with enough information to allow them to be uniquely identified, should be included in the Methods section. Authors are strongly encouraged to cite <a href="#">Research Resource Identifiers</a> (RRIDs) for antibodies, model organisms and tools, where possible.</p> <p>Have you included the information requested as detailed in our <a href="#">Minimum Standards Reporting Checklist</a>?</p>	<p>Yes</p>
<p><b>Availability of data and materials</b></p> <p>All datasets and code on which the conclusions of the paper rely must be either included in your submission or deposited in <a href="#">publicly available repositories</a> (where available and ethically appropriate), referencing such data using a unique identifier in the references and in the “Availability of Data and Materials” section of your manuscript.</p> <p>Have you have met the above requirement as detailed in our <a href="#">Minimum Standards Reporting Checklist</a>?</p>	<p>Yes</p>

# 1 **Mutation Impact on mRNA Versus Protein Expression across Human Cancers**

2

3 Yuqi Liu<sup>1\*</sup>, Abdulkadir Elmas<sup>1\*</sup>, Kuan-lin Huang<sup>1#</sup>

4

5 <sup>1</sup> Department of Genetics and Genomic Sciences, Department of Artificial Intelligence  
6 and Human Health, Center for Transformative Disease Modeling, Tisch Cancer Institute,  
7 Icahn Genomics Institute, Icahn School of Medicine at Mount Sinai, New York, NY  
8 10029, USA.

9 \* These authors contributed equally to this work.

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

31

32

33

#Corresponding Author:

34

Kuan-lin Huang, Ph.D.

35

Departments of Genetics and Genomic Sciences &amp; Artificial Intelligence and Human

36

Health

37

Icahn School of Medicine at Mount Sinai

38

New York, NY 10029

39

Email: [kuan-lin.huang@mssm.edu](mailto:kuan-lin.huang@mssm.edu)

40 **ABSTRACT**

41 Cancer mutations are often assumed to alter proteins, thus promoting tumorigenesis.  
42 However, how mutations affect protein expression has rarely been systematically  
43 investigated. We conduct a comprehensive analysis of mutation impacts on mRNA- and  
44 protein-level expressions of 953 cancer cases with paired genomics and global proteomic  
45 profiling across six cancer types. Protein-level impacts are validated for 47.2% of the  
46 somatic expression quantitative trait loci (seQTLs), including mutations from likely “long-  
47 tail” driver genes. Devising a statistical pipeline for identifying somatic protein-specific  
48 QTLs (spsQTLs), we reveal several gene mutations, including *NF1* and *MAP2K4*  
49 truncations and *TP53* missenses showing disproportional influence on protein abundance  
50 not readily explained by transcriptomics. Cross-validating with data from massively  
51 parallel assays of variant effects (MAVE), *TP53* missenses associated with high tumor  
52 *TP53* proteins were experimentally confirmed as functional. Our study demonstrates the  
53 importance of considering protein-level expression to validate mutation impacts and  
54 identify functional genes and mutations.

55

## 56 INTRODUCTION

57 Cancer arises from the acquisition of mutations that confer selective advantages. The  
58 majority of these mutations are thought to affect cellular functions by regulating the  
59 expression of gene products. For example, truncations can result in nonsense-mediated  
60 decay (NMD)<sup>1,2</sup>, which protects eukaryotic cells through degrading premature termination  
61 codon (PTC) bearing mRNA<sup>3</sup>. Additionally, a fraction of cancer mutations may uniquely  
62 affect protein abundance but not mRNA expression. However, previous studies  
63 characterizing genomic mutations affecting mRNA vs. protein levels have focused on  
64 germline variants as expression quantitative trait loci (eQTL)<sup>4-6</sup>. While other cancer  
65 studies have characterized the effect of somatic mutations on mRNA expression levels<sup>7-  
66 9</sup>, it remains unclear how somatic mutations may affect protein abundance. The gap of  
67 knowledge is critical given that mRNA and protein levels are only moderately correlated<sup>10-  
68 13</sup>. A myriad of factors, including cell state transition, signal delay, translation on demand,  
69 and cellular energy constraint, can lead to discrepancies between mRNA and protein  
70 levels<sup>14</sup>. Understanding protein-level consequences of cancer mutations is critical in  
71 identifying functionally important mutations and revealing their downstream mechanisms.

72 In recent years, advances in mass spectrometry (MS) technologies have generated a  
73 wealth of global proteomic profiles of primary tumor cohorts, many of which also have  
74 concurrent genomic and transcriptomic profiling<sup>15-20</sup>. These proteogenomic datasets  
75 present ample opportunities to validate somatic mutations that show concordant impacts  
76 on downstream mRNA and protein levels. On the other hand, protein abundance may  
77 also be uniquely influenced by the efficiency of protein translation, transport, and  
78 degradation. Thus, proteogenomic analyses can reveal mutations that disproportionately  
79 impact protein abundances that may not be found using genomic analyses alone.

80 Herein, we conducted a systematic analysis to decode the relationship between somatic  
81 mutations vs. mRNA and protein levels using data from nearly a thousand cases across  
82 six cancer types in prospective and retrospective cohorts from the Clinical Proteomic  
83 Tumor Analysis Consortium (CPTAC). We identified mutations showing concordant  
84 effects at both mRNA and protein expression levels *in cis*, as well as those that showed  
85 protein-specific effects. We further examined how mutations associated with expression

86 changes may predict *in vitro* and *in vivo* functional effects measured by a massively  
87 parallel assays of variant effects (MAVE) of TP53<sup>21</sup>. Our results highlight the importance  
88 of pairing genomic and proteomic analyses to prioritize functionally important mutations.

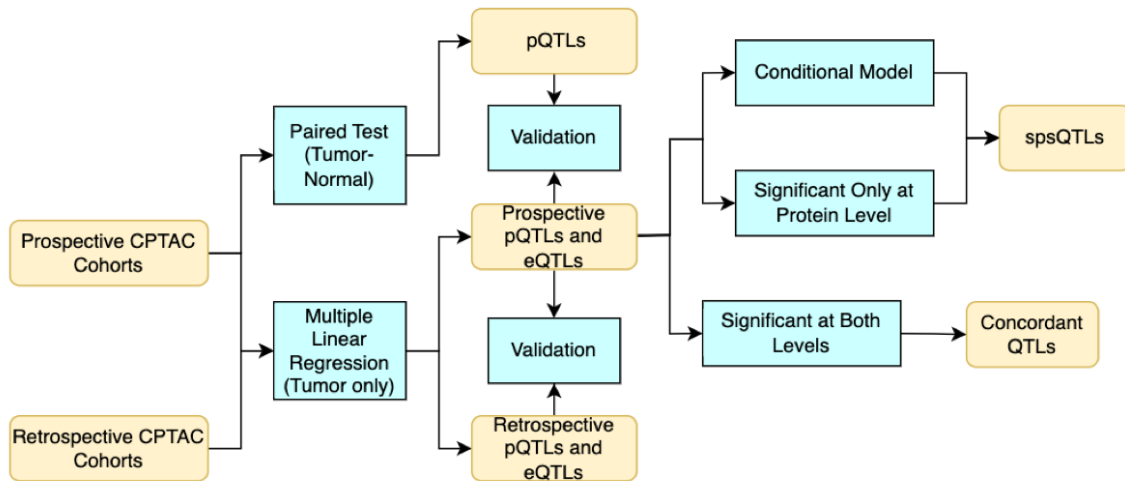
## 89 **RESULTS**

### 90 **Mutation impacts on the mRNA and protein levels**

91 Following the study workflow (**Figure 1A**), we first sought to identify somatic mutations  
92 that may impact the corresponding gene's mRNA expression (somatic eQTL, termed  
93 seQTL below) and protein abundance (somatic pQTL, termed spQTL below) in primary  
94 tumor tissue samples. We performed a multiple regression analysis adjusted for age,  
95 gender, ethnicity, and TMT batch using the prospective CPTAC datasets that included  
96 matched DNA-Seq, RNA-Seq, and mass spectrometry (MS) global proteomics data of  
97 primary tumor samples across six cancer types (**Methods, Figure 1B**), including 115  
98 breast cancer (BRCA)<sup>19</sup>, 95 colorectal cancer (CRC)<sup>16</sup>, 110 clear cell renal cell carcinoma  
99 (CCRCC)<sup>15</sup>, 109 lung adenocarcinoma (LUAD)<sup>17</sup>, 84 ovarian cancer (OV)<sup>20</sup>, and 97  
100 uterine corpus endometrial carcinoma (UCEC)<sup>18</sup>, as well as proteogenomic datasets for  
101 additional, retrospective BRCA<sup>11</sup>, CRC<sup>13</sup>, and OV<sup>12</sup> cohorts from CPTAC for validation  
102 (**Figure S1A**). We focused on coding mutations given the coverage of the whole-exome  
103 sequencing (WES) data used in CPTAC studies; the analyses were further stratified for  
104 truncations, missense, and synonymous mutations given their likely different mechanisms  
105 of action in affecting levels of the mutated gene product.



A



B

Cancer Type	Breast Cancer	Clear Cell Renal Cell Carcinoma	Colorectal Cancer	Lung Adenocarcinoma	Ovarian Cancer	Uterine Corpus Endometrial Carcinoma
Abbreviation	<b>BRCA</b>	<b>CCRCC</b>	<b>CRC</b>	<b>LUAD</b>	<b>OV</b>	<b>UCEC</b>
Data Source	Krug et al. 2020 (PMID: 33212010)	Clark et al. 2019 (PMID: 31675502)	Vasaikar et al. 2019 (PMID: 31031003)	Gillette et al. 2020 (PMID: 32649874)	McDermott et al. 2020 (PMID: 32529193)	Dou et al. 2020 (PMID: 32059776)
Sample Size (Tumors/Normals)	T: 115 N: 18	T: 110 N: 84	T: 95 N: 100	T: 109 N: 102	T: 84 N: 19	T: 97 N: 20
Female %	100%	25.2%	57.4%	34.6%	100%	100%
Average Onset (yr)	60.4	60.6	65.2	62.7	59.1	63.7
Tumor Stage	1: 3% 2: 60.4% 3: 26.9% NA: 9.7%	1: 41.8% 2: 15.5% 3: 34.5% 4: 8.2%	1: 10.2% 2: 40.6% 3: 41.1% 4: 8.1%	1: 53.5% 2: 27.5% 3: 18.5% 4: 0.5%	1: 1% 2: 1% 3: 72.8% 4: 15.5% NA: 9.7%	1: 76.1% 2: 6.8% 3: 14.5% 4: 2.6%

106

107

108

109

110

111

112

113

114

115

116

117

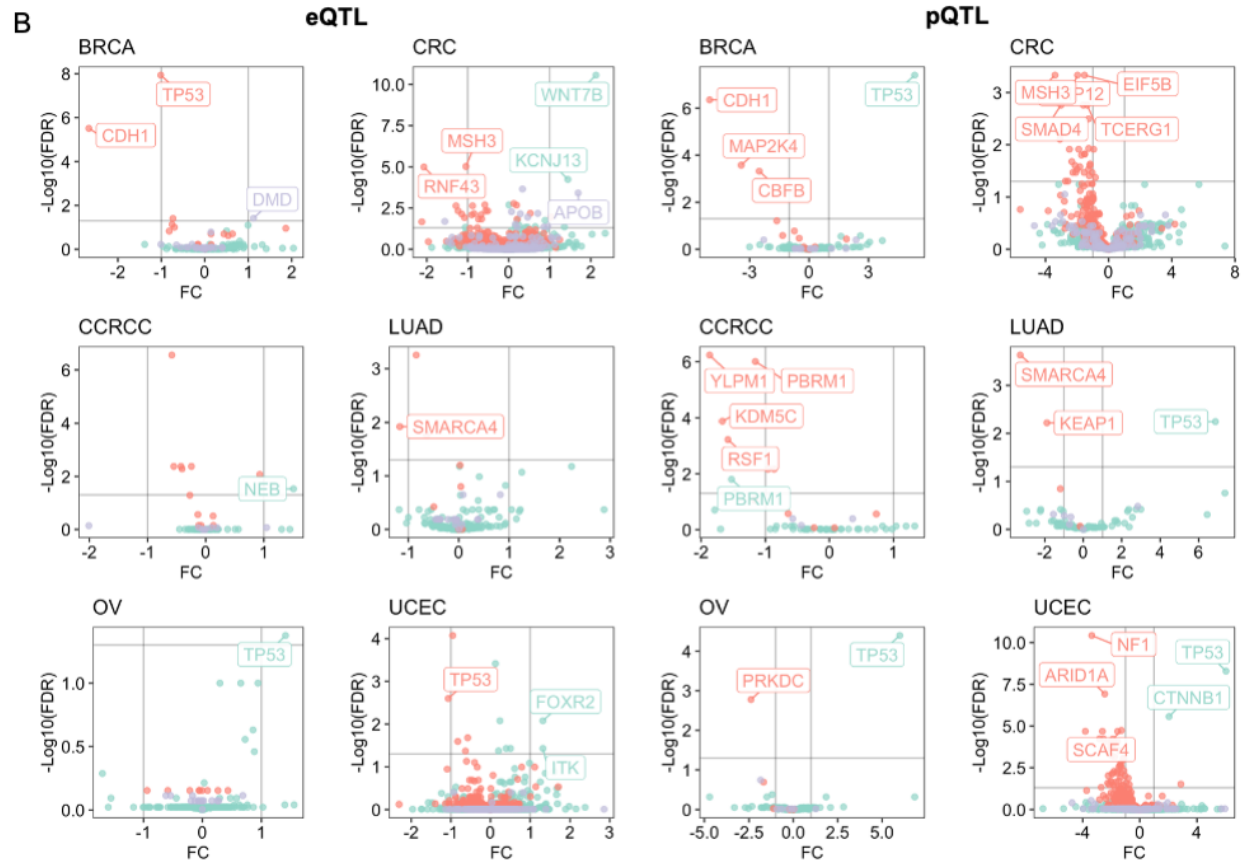
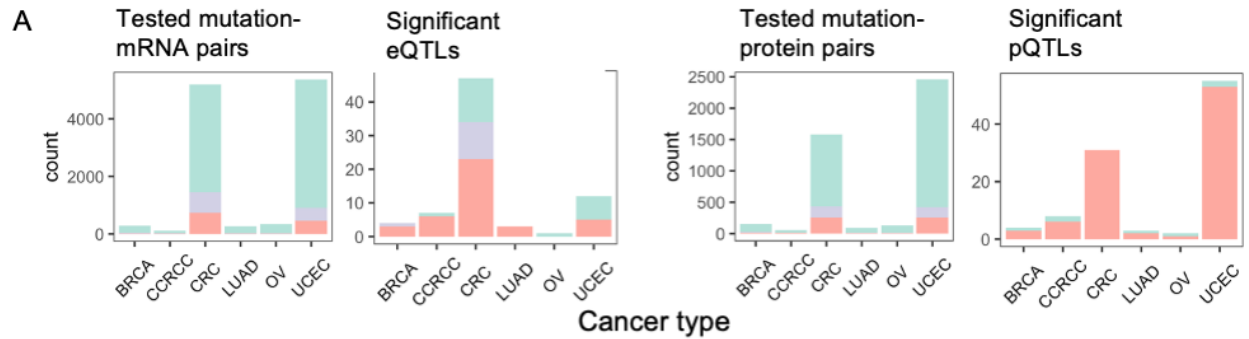
118

**Figure 1. Overview of the study workflow and proteogenomic cohorts.** (A) Study workflow to identify eQTLs, pQTLs, concordant QTLs (between mRNA and protein levels), and spsQTLs showing disproportional effects on protein expression. (B) Summary of the prospective CPTAC proteogenomic cohorts used for the discovery analyses, including cancer type abbreviation, data source, sample size of tumor (T) and normal (N) tissues, female percentage, average onset age in years, and tumor stage distribution.

Based on the statistical power achieved by these cohort sizes and to reduce false positives, we focused on genes with three or more samples affected by mutations in each functional class of missense, truncation, and synonymous within the cancer cohort, including 134, 13, and 15 genes tested in BRCA; 1360, 318, and 226 genes tested in CRC; 55, 12, and 4 genes tested in CCRCC; 94, 4, and 8 genes tested in LUAD; 134, 5,

119 and 8 genes tested in OV; 2243, 273, and 196 genes tested in UCEC. We sought to  
120 identify their seQTLs affecting *cis*-expression, i.e., expression of the mutation-affected  
121 genes. Using the multiple regression model (**Methods**), we identified 74 gene-cancer  
122 seQTL pairs (FDR < 0.05), including 4 in BRCA, 47 in CRC, 7 in CCRCC, 3 in LUAD, 1  
123 in OV, and 12 in UCEC (**Figure 2A, Table S1**). Separated by the functional classes of  
124 mutations, 22 of those seQTLs are missense mutations, 12 are synonymous, and 40 are  
125 truncating. Top seQTLs showing up-regulation of gene expression are primarily  
126 missenses, including *SMARCA4* in LUAD, *WNT7B* in CRC, *TP53* in OV, and *FOXR2* in  
127 UCEC. Top candidates showing down-regulation of gene expression include *TP53* and  
128 *CDH1* truncations in BRCA, as well as *TP53* truncations in OV (**Figure 2B**).

129



Mutation type ■ missense ■ truncating ■ synonymous

130  
 131 **Figure 2. Gene mutations identified as *cis* seQTLs and spQTLs across six adult cancer types.** (A)  
 132 Overview of the somatic mutation QTLs identified in different cancer types and mutation types, including  
 133 missense (green), truncating (orange), and synonymous (purple) mutations. For both eQTLs and pQTLs,  
 134 the panel on the left shows the counts of the mutation-gene pairs included in analyses, and the figure on  
 135 the right shows the counts of the significant eQTLs and pQTLs. (B) Volcano plots showing seQTLs  
 136 associations in the six cancer types (left) and volcano plots showing spQTLs associations (right), where  
 137 each dot denotes a gene-cancer pair included in the analysis. Top associated genes were further labeled.  
 138 FC: mRNA/protein expression log fold change. FDR: false discovery rate.

139

140 Using a similar multiple regression but modeling protein abundance as the dependent  
141 variable, we identified 103 significant gene-cancer spQTL pairs (FDR < 0.05), including 4  
142 in BRCA, 31 in CRC, 8 in CCRCC, 3 in LUAD, 2 in OV, and 55 in UCEC (**Figure 2A**,  
143 **Table S2**). Compared to the proportion of gene-mutation type evaluated in each cancer  
144 type, spQTLs showed significant enrichment for truncations (Fisher exact test p-value <  
145 0.05; **Figure 2A**), highlighting the persistent and more profound effect of truncations on  
146 protein abundance compared to mRNA levels. Among the identified spQTLs across  
147 cancer, 7 are missense and 96 are truncating. For example, truncating mutations of *NF1*  
148 and *ARID1A* in UCEC, and *YLPM1* in CCRCC are each associated with reduced protein  
149 level of the corresponding gene (**Figure 2B**). Notably, *TP53* missenses in OV, BRCA,  
150 LUAD, and UCEC are each significantly associated with increased protein expression in  
151 mutation carriers (**Figure 2B**).

152

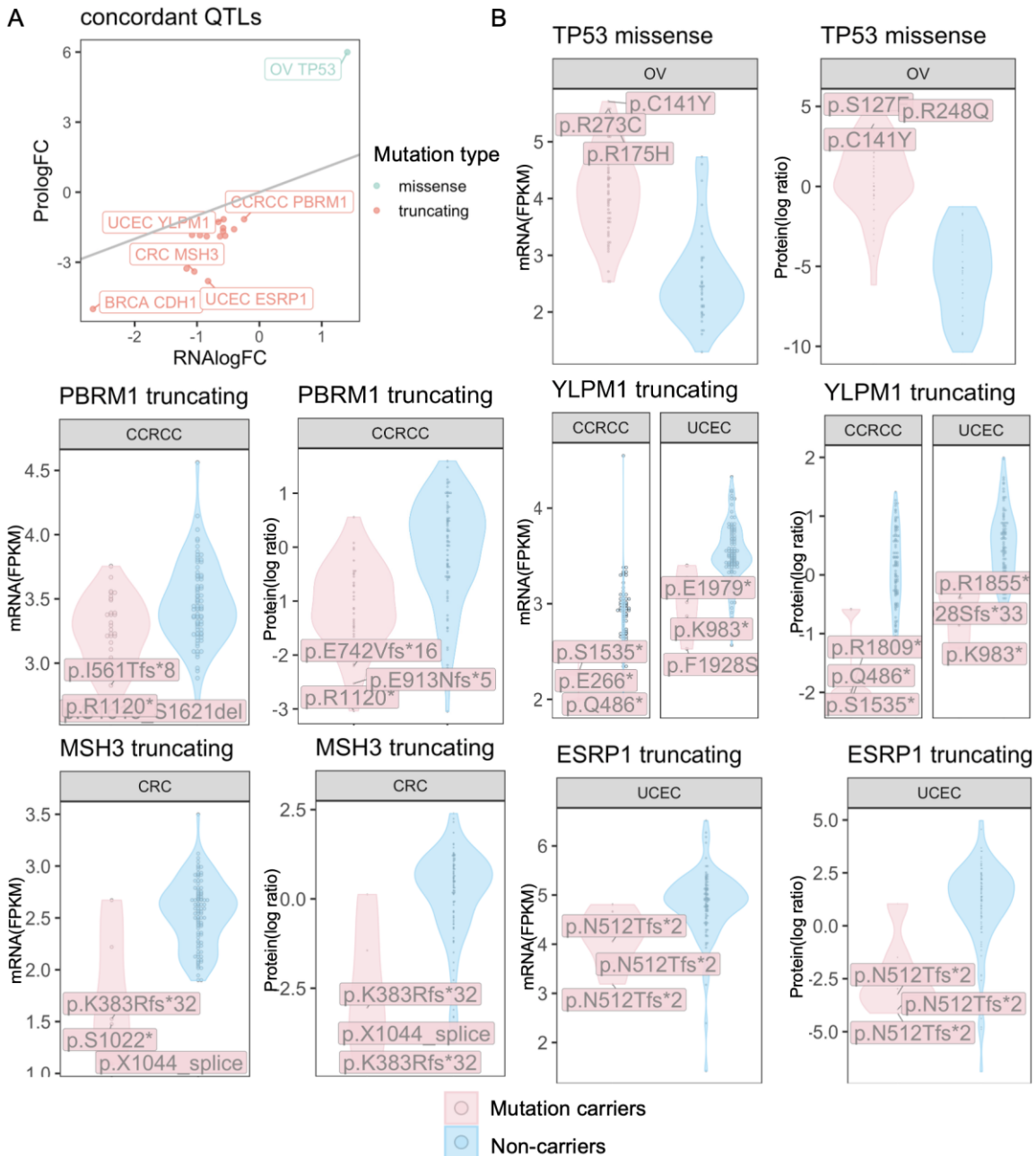
153 To verify these discoveries, we applied the same seQTL and spQTL analyses using  
154 retrospective CPTAC data (**Figure S1A**) that included independent cohorts of BRCA<sup>11</sup>,  
155 CRC<sup>13</sup>, and OV<sup>12</sup> primary tumors. While these cohorts afforded smaller sample sizes, 8  
156 seQTLs and 5 spQTLs were detected in both retrospective and prospective sets. The  
157 gene-cancer spQTL pairs showing strong validation in both datasets include *TP53*  
158 missense mutations and *CDH1* truncations in BRCA, and *TP53* truncations in CRC  
159 (**Figure S1B**).

160

### 161 **Mutations showing concordant effects at mRNA and protein levels**

162 We next examined the concordance of seQTL and spQTL associations for each gene-  
163 cancer type pair. As expected, for most (88.9%) of the significant seQTLs whose genes  
164 had sufficient observations at both the mRNA and protein levels, the identified  
165 associations showed the same directionality. However, we only identified 17 seQTLs  
166 (47.2%) that are also significant spQTLs at an FDR < 0.05, which we show as concordant  
167 QTLs (**Figure 3A, Table S3**). The effect sizes (in log fold change) of these gene-cancer  
168 pairs showing concordant seQTLs and spQTLs showed a high correlation between  
169 mRNA and protein (Pearson  $r = 0.90$ , p-value <  $7.51E-7$ ).

170



171  
 172 **Figure 3. Gene mutations showing concordant impacts on gene and protein expression levels.** (A)  
 173 Overview of concordant QTLs as shown by their effect sizes in log[Fold Change (FC)], where the gray line  
 174 shows when the protein logFC equals RNA logFC. Some of the top concordant QTLs were further labeled  
 175 by cancer type and gene name. (B) Examples of QTL with concordant effects at mRNA and protein  
 176 expression levels. For each gene, the plot on the left shows the corresponding mRNA levels of mutation  
 177 carriers vs. non-carriers in FPKM, and the plot on the right shows protein level comparison in log ratio (MS  
 178 TMT measurements) in the respective cancer type labeled on top of each of the violin plots. The labeled  
 179 mutations are the three mutations whose carriers show the highest absolute expression differences of the  
 180 mutated gene product compared to the non-carriers.

181

182 In different cancer types, genes whose mutation impacts on gene and protein expressions  
183 are concordant include well-known drivers of the disease, including *TP53* missense  
184 mutations in OV, *CDH1* truncations in BRCA, and *MSH3* truncations in CRC. Up-  
185 regulation of mutated *TP53* in OV is the only association found for genes affected by  
186 missense mutations. The 16 other concordant se/spQTLs are all truncations associated  
187 with reduced expression and highlight some “long-tail” driver genes, including *PBRM1* in  
188 CCRCC, *YLPM1* in CCRCC/UCEC, and *ESRP1* in UCEC (**Figure 3B**). The concordant  
189 QTLs with truncating mutation can likely be explained by NMD, which reduces gene  
190 expression and in turn diminishes the expression of the corresponding proteins<sup>3</sup>.  
191 Compared to the substantially higher counts of seQTL associations (**Figure 2A-B**), these  
192 concordant se/spQTL effects validate mutation impacts on the gene product.

193

194

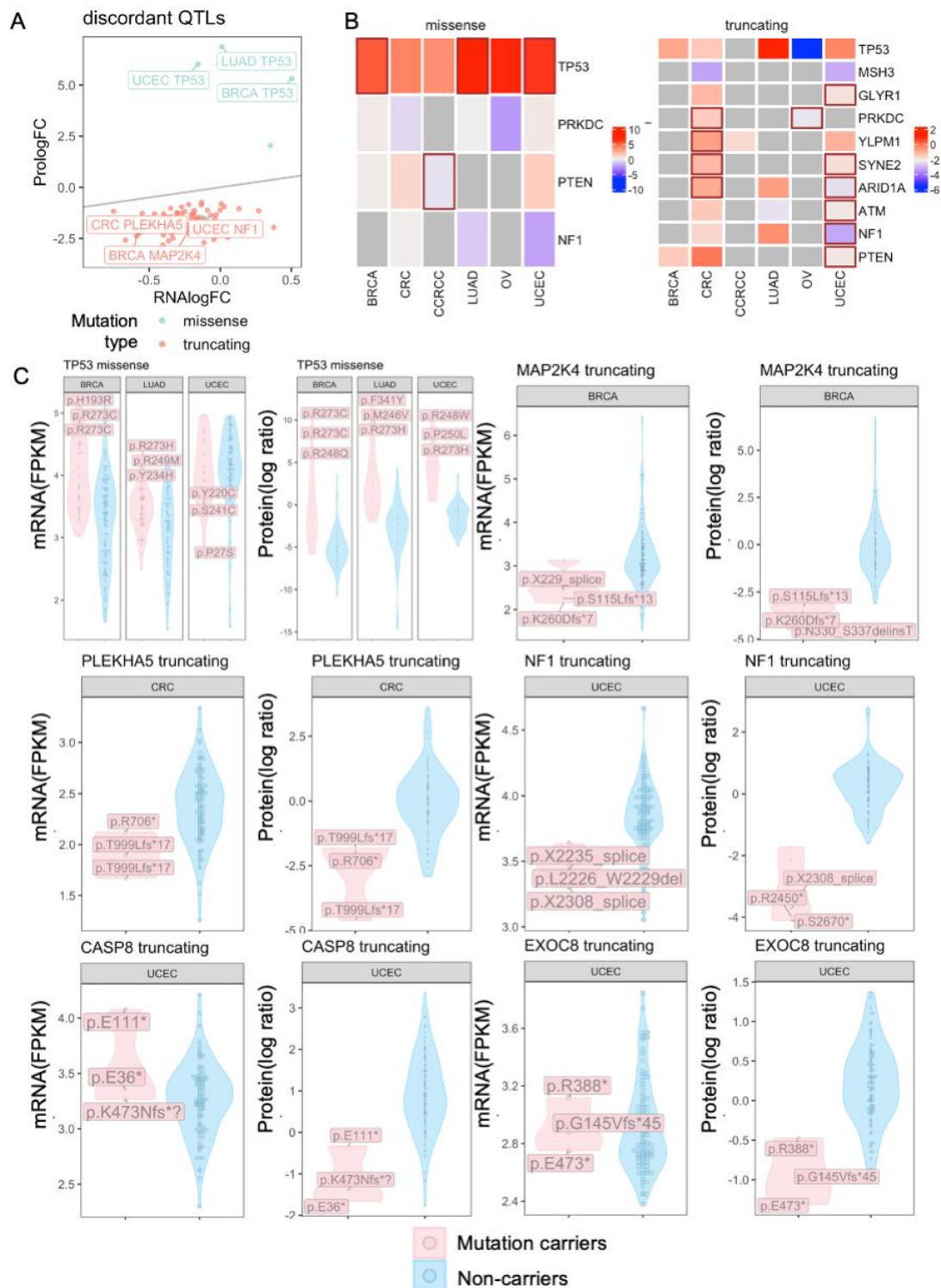
#### 195 **Protein-specific mutation impacts not observed at mRNA levels**

196 While most seQTLs and spQTLs show concordance, we postulate that certain mutations  
197 may uniquely affect protein abundance but not mRNA levels, which we term somatic  
198 protein-specific QTLs (spsQTLs). To identify spsQTLs, we applied two methods to  
199 stringently retain QTLs with discordant effects at mRNA and protein levels. First, applying  
200 a likelihood ratio test (LRT) between two regression models of protein level being  
201 predicted by mRNA level with or without the mutation term (**Methods**)<sup>4</sup>, 96 candidate  
202 spsQTLs (FDR < 0.05) were identified. Second, complementing this LRT test with an  
203 approach filtering for gene-cancer pair showing significant spQTL (FDR < 0.05) but not  
204 seQTLs (**Methods**)<sup>22</sup>, 86 candidate spsQTLs (FDR < 0.05) were identified.

205

206 By overlapping candidate spsQTLs identified by both methods, we retained 83 spsQTLs,  
207 the majority (92.8%) of which are truncating mutations (**Figure 4A, Table S4**). Top  
208 spsQTLs associated with diminished protein expression include *NF1* truncations in UCEC,  
209 *PLEAHK5* truncations in CRC, and *MAP2K4* truncations in BRCA. The only spsQTLs that  
210 increase protein expression include *TP53* missense mutations in BRCA, LUAD, and  
211 UCEC. (**Figure 4B**). We further examined the discordance in mutation impacts on gene

212 and protein expression levels (**Figure 4C**). While some of these truncations, such as *NF1*  
213 in UCEC and *MAP2K4* in BRCA, were often accompanied by lower-than-median mRNA  
214 expression in their respective tumor cohorts, their impacts were strikingly observed at  
215 diminished protein expression levels. We highlighted in **Figure S2A** spsQTLs where the  
216 affected gene's protein showed negative protein log fold-change (logFC) whereas the  
217 mRNA logFC is non-negative, including *CASP8* truncations in UCEC, *ARID1A* truncations  
218 in CRC and UCEC, and *ATM* truncations in LUAD and UCED. We also identified a set of  
219 spsQTLs truncations, where the logFC associated with a reduction in proteins is 15 times  
220 greater than mRNAs logFC (**Figure S2B**). These results suggest that NMD associated  
221 with these gene truncations are closely tied to the terminated translation but may not  
222 affect mRNA expression to the same degree<sup>23</sup>.



223

224 **Figure 4. Gene mutations showing discordant impacts on gene and protein expression levels. (A)**

225 Overview of discordant QTLs identified by our statistical pipeline as shown by their effect sizes in log[Fold

226 Change (FC)], where the gray line shows when the protein logFC equals RNA logFC. (B) Heatmaps of

227 QTLs that are significant as either seQTL or spQTL and that are shared across at least two cancer types.



228 Brown box indicates significant spsQTLs, and color indicates the effect size in log[Fold Change (FC)],  
229 average protein expression of mutation carriers in log ratio from the MS TMT quantifications. (C) Examples  
230 of QTL with discordant effects at mRNA vs. protein levels. For each gene, the plot on the left shows the  
231 corresponding mRNA levels of mutation carriers vs. non-carriers in FPKM, and the plot on the right shows  
232 protein level comparison in log ratio (MS TMT measurements) in the respective cancer type labeled on top  
233 of each of the violin plots. The labeled mutations are the three mutations whose carriers show the highest  
234 absolute expression differences of the mutated gene product compared to the non-carriers.

235

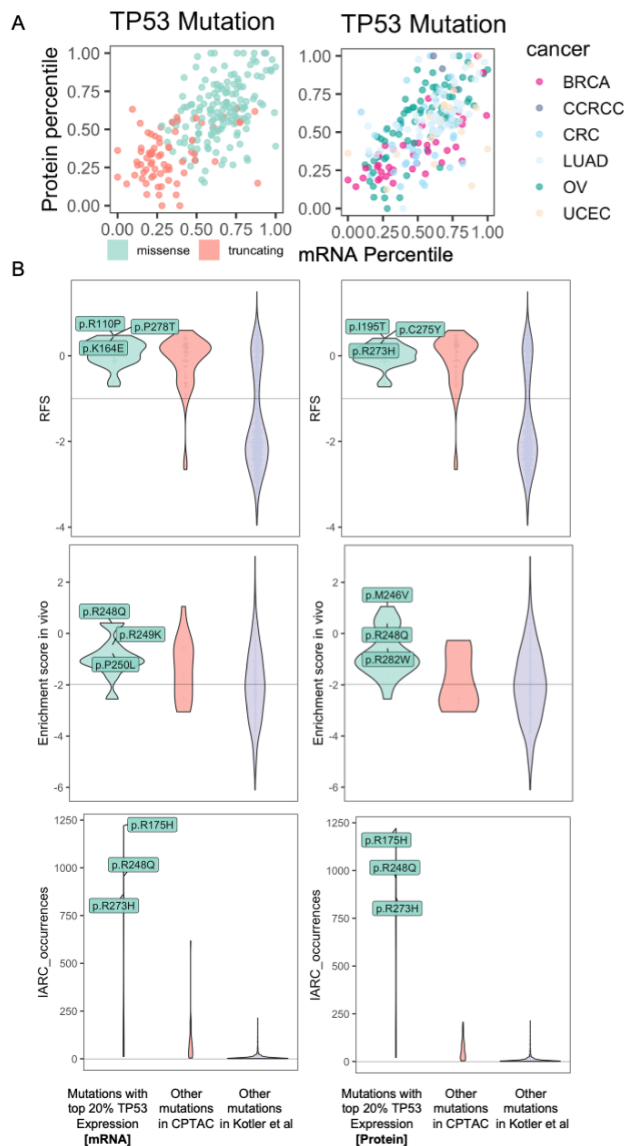
236 To complement the cross-tumor analyses, we also utilized the CPTAC samples with  
237 paired tumor-normal tissues to conduct paired differential expression tests for both protein  
238 and mRNA expression (**Figure 1A**). The paired sample sizes with proteomic data include  
239 17 in BRCA, 17 in UCEC, 84 in CCRCC, 100 in LUAD, 29 in CRC, and 10 in OV (**Figure**  
240 **1B**). Covariates including age at diagnosis, ethnicity, race, and sequencing operator are  
241 adjusted in the analysis. While this analysis had varied statistical power due to different  
242 normal tissue availabilities across cancer types, it served as an independent validation of  
243 spQTLs (**Table S5**). This paired tumor-normal analysis validated the protein-level impacts  
244 of several discordant spsQTLs (**Figure S3A**) as well as some concordant se/spQTLs  
245 (**Figure S3B**). For example, the validated discordant spsQTLs include truncations of  
246 *SMAD4* and *SCRIB* in CRC as well as *NF1*, *GLYR1*, and *RASA1* in UCEC (**Figure S3A**).  
247 The validated concordant se/spQTLs include truncations of *YLPM1* and *PBRM1* in  
248 CCRCC, *SMARCA4* and *KEAP1* in LUAD, and *ESRP1* as well as *JAK2* in UCEC (**Figure**  
249 **S3B**).

250

### 251 **Functional evidence of *TP53* missenses associated with high protein expression**

252 Notably, *TP53* missenses are associated with higher protein expression in multiple cancer  
253 cohorts, in addition to the expected reduction in expression associated with truncations  
254 (**Figure 5A**). Such cis-effect of functional *TP53* missense mutations had previously been  
255 observed through immunohistochemistry (IHC<sup>24</sup>) or MS global proteomics experiments<sup>25</sup>.  
256 Here, we hypothesized that functional *TP53* missense mutations are more likely to show  
257 high levels of concurrent protein-level expression in the mutated tumor sample. To test  
258 this hypothesis, we compared gene and protein-level *TP53* expression from CPTAC with  
259 *TP53* mutation-level functional data from the *in vitro* and *in vivo* MAVe experiment

260 conducted by Kotler et al<sup>21</sup>, where they designed a p53 variants library to study the  
 261 functional impact of those mutations.



262  
 263 **Figure 5. Functional verification of *TP53* mutation associated with high mRNA or protein levels**  
 264 **using *in vitro* and *in vivo* data from a MAVE experiment.** (A) Percentile of averaged expression  
 265 associated with a given *TP53* mutation at the mRNA (x-axis) and protein (y-axis) levels in the respective  
 266 cancer cohort. *TP53* mutations are color coded by mutation type (left) and observed cancer type (right),  
 267 respectively. (B) Violin plots comparing the *in vitro* functional score (RFS, top), *in vivo* enrichment score  
 268 (middle), and IARC occurrences (bottom) for *TP53* mutations in the three groups defined by (1) *TP53*  
 269 mutations with top 20% mRNA (left) or protein (right) expression in the prospective CPTAC cohorts, (2) the  
 270 other *TP53* mutations observed across all CPTAC samples, and (3) the rest of the assayed *TP53* mutations  
 271 from Kotler et al<sup>21</sup>.

272

273 We divided the *TP53* missense mutations from Kotler et al. into three categories: (1) *TP53*  
274 mutations with top 20% mRNA or protein expression in the prospective CPTAC cohorts,  
275 (2) the other *TP53* mutations observed across all CPTAC samples, and (3) the rest of the  
276 assayed *TP53* mutations from Kotler et al. For *in vitro* data, the number of tested  
277 mutations by each category is 32, 78, and 1,033, respectively. For *in vivo* data, the  
278 number of tested mutations by each category is 19, 10, and 381, respectively. We first  
279 compared the relative fitness score (RFS) measured from the *in vitro* assays<sup>17</sup>. While  
280 there may be a trend, we did not observe a significant difference between all the other  
281 mutations versus *TP53* missenses associated with either top 20% expression based on  
282 either mRNA (p-value = 0.090, Wilcoxon rank-sum test) or protein expression (p-value =  
283 0.720).

284

285 We next compared the *in vivo* enrichment scores across the same categories, and found  
286 that *TP53* missenses associated with top 20% protein expression showed significantly  
287 higher enrichment score *in vivo* compared to that of other *TP53* missenses found in  
288 CPTAC (p-value = 0.016) or other experimentally-measured *TP53* mutations (p-value =  
289 3.23E-5, **Figure 5B, Table S6**). In comparison, *TP53* missenses associated with top 20%  
290 mRNA expression did not show a significant *in vivo* score difference to that of other *TP53*  
291 missenses found in CPTAC (p-value = 0.170). Kotler et al. observed that there was no  
292 significant correlation between enrichment score *in vivo* and RFS *in vitro*, which is  
293 consistent with our observations and may be explained by the different selective  
294 pressures between these settings *in vivo* and *in vitro*<sup>21</sup>. Finally, *TP53* missenses  
295 associated with top 20% protein expression (p-value = 5.91E-7) or top 20% mRNA  
296 expression (p-value = 2.38E-2) showed significantly higher prevalence than other CPTAC  
297 mutations based on counts from the International Agency for Research on Cancer (IARC)  
298 database<sup>21</sup> (**Figure 5B, Table S6**). Overall, these analyses suggested that protein-level  
299 consequences from primary tumor samples can aid the identification of functional  
300 mutations.

301

302

## 303 DISCUSSION

304

305 Herein, we analyzed how somatic mutations affect mRNA and protein levels using  
306 matched genomic, transcriptomic, and global proteomic data from 953 cases across six  
307 solid cancer types. We first investigated the mutation impacts at the mRNA level and  
308 protein level, finding that although most seQTLs have the same direction of effect as  
309 spQTLs, less than half of them are also significant at the protein level. We also studied  
310 the concordant or discordant relationship between seQTL versus spQTLs, finding several  
311 spsQTLs that have disproportional effects on protein. Finally, we conducted analyses to  
312 provide functional validation<sup>21</sup> for our findings of TP53 missenses associated with high  
313 protein expression.

314

315 Integrating protein-level data identified nearly 47.2% seQTLs as concordant, significant  
316 spQTLs. The result demonstrates the capacity of proteomic data to validate genomic  
317 findings and potentially filter out noises that may arise for example due to the more  
318 transient nature of transcription compared to translation. In addition to well-known tumor  
319 suppressors like *TP53* and *MSH3*, other gene mutations with concordant effects may also  
320 be “long tail” driver genes that will otherwise require large cohort sample sizes to discover.  
321 For example, *PBRM1*, which we found in CCRCC, is a subunit of the PBAF chromatin  
322 remodeling complex thought to be a tumor suppressor gene whose mutations may confer  
323 synthetic lethality to DNA repair inhibitors<sup>26</sup>. *ESRP1*, found in UCEC, is crucial in  
324 regulating alternative splicing and the translation of some genes during organogenesis<sup>27</sup>.  
325 Other less-studied genes we identified include *YLPM1* truncations associated with  
326 concordantly reduced *YLPM1* mRNA and protein expression levels in both CCRCC and  
327 UCEC. Analyzing the distribution of these gene mutations on NCI’s Genome Data  
328 Commons, we observed many other recurrent truncations (**Figure S4**), suggesting these  
329 mutations may represent some of the “long tail” driver mutations that warrant further  
330 investigation<sup>28,29</sup>.

331

332 By devising a specific pipeline to detect spsQTLs, our results showed that apart from  
333 mutations that influence protein level mediated by changes in mRNA level, many

334 mutations are associated with disproportional aberrations at the protein level compared  
335 to mRNA changes, indicating post-transcriptional regulation. SpsQTLs were found to  
336 affect known driver genes such as *TP53* missenses, and truncations in *NF1*<sup>30</sup> and  
337 *MAP2K4*<sup>31</sup>. In most cases, protein molecules are more direct mediators of cellular  
338 functions and phenotypes than mRNAs<sup>32</sup>. Thus, the discordant effect between mRNA  
339 level and protein level discovered in our study highlights the importance of exploring  
340 disease mechanisms and developing treatments at the protein level.

341  
342 One possible source of spsQTLs is the imperfect correlation between mRNA and protein  
343 expression in the affected genes. Additional statistical analyses revealed that this mRNA-  
344 protein correlations range widely across genes and cancer types (**Figure S5**). While  
345 genes harboring spsQTLs have lower mRNA-protein correlations in general than genes  
346 with concordant eQTL and pQTL, this is not the case for several discordant genes,  
347 including *MAP2K4* in BRCA and *PBRM1* in CCRCC (**Table S7**). Based on the number of  
348 mutations and genes identified, CRC and UCEC reached statistically significant  
349 differences between concordant and all other expressed genes (Wilcoxon rank-sum tests,  
350  $p = 0.0056$  and  $p = 0.022$ , respectively); in CRC, mRNA-protein correlations also showed  
351 significant differences between discordant and all other expressed genes ( $p = 0.013$  and  
352  $p = 0.29$ , respectively); other cancer types likely did not reach statistical significance likely  
353 due to sufficient mutations identified. The imperfect correspondence between gene  
354 mRNA-protein correlations and mutation impacts further stresses the need to analyze and  
355 consider protein-specific impacts of mutations. **Table S7** provides complete mRNA-  
356 protein correlation data for all concordant/discordant eQTL/pQTLs in their respective  
357 cancer type for in-depth examination.

358  
359 This study has several limitations. First, our findings do not distinguish between several  
360 potential mechanisms that could lead to discordant effects of mutations on gene and  
361 protein expression. One possibility is that the mutation affects the efficiency of translation,  
362 leading to changes in protein levels that are not reflected in mRNA levels. For example,  
363 accumulating evidence in recent years suggests that NMD is closely tied to the  
364 termination of translation<sup>23</sup>, which may explain instances where some truncations afford

365 much stronger associations with protein levels in our findings. But, in many cases, the  
366 mechanisms of how mutations may affect protein abundance may be context- and gene-  
367 specific and remain to be elucidated. For example, certain mutations may influence the  
368 binding of RNA binding proteins and the efficiency of translation, whereas others may  
369 alter post-translational modifications, such as phosphorylation or ubiquitination, which  
370 can impact protein stability or degradation without affecting transcription or translation  
371 rates. Second, the proteogenomic tumor cohorts used herein, while being some of the  
372 largest studies to date, still are limited in sample sizes and preclude sufficient statistical  
373 power to identify pQTLs at a single mutation level or reveal *trans* effects. Third, given the  
374 limitation of current omic technology and data, our findings do not resolve mutation impact  
375 on proteins at the temporal, spatial, or single-cell resolution, but provide candidate  
376 mutations to be investigated in future studies. Fourth, our regression models assumes a  
377 linear relationship between mutations (one gene at a time), confounders, and expression,  
378 which may not capture more complex, nonlinear effects of mutations on multiple mRNA  
379 or protein expression. Future studies could explore non-linear regression models or  
380 neural network approaches to better account for these effects. Fifth, we employed two  
381 complementary methods to confidently identify spsQTLs that represent true protein-  
382 specific regulatory events. However, the reliance on FDR thresholds could still limit the  
383 detection of spsQTLs with subtle effects. Alternative approaches, such as Bayesian  
384 models that account for prior biological knowledge or hierarchical modeling, could be  
385 considered in future analyses to improve the specificity of spsQTL detection. Additionally,  
386 while our method focuses on cis-acting mutations, potential trans-acting effects could be  
387 missed, a limitation that should be explored in larger datasets or by incorporating network-  
388 based analyses.

389  
390 Finally, using *TP53* missense mutations as an example, we showed that protein-level  
391 expression can serve as an effective strategy to prioritize functional mutations. As DNA-  
392 Seq become ever more commonplace, many rare mutations are being identified and it  
393 remains challenging to accurately classify their functional impacts. Our data  
394 demonstrated that *TP53* missenses associated with high protein expression show  
395 significantly higher functional scores, particularly those measured *in vivo*. This protein-

396 expression-based prioritization strategy can be particularly powerful when combined with  
397 high-throughput functional assays like using MAVE model systems that are typically *in*  
398 *vitro*. Considering that both MAVE and proteogenomic datasets of tumor cohorts are both  
399 expanding quickly in the next few years<sup>33,34</sup>, the combined approaches can help  
400 effectively pinpoint functional mutations for mechanistic and clinical characterization. The  
401 prioritized mutations based on protein-level consequences may also guide the selection  
402 of targeted therapy to advance precision medicine.

## 403 **METHODS**

### 404 **Proteogenomic datasets**

405  
406 The prospective CPTAC data were downloaded and processed as described in the  
407 Method section of the work of Elmas et al<sup>35</sup>. The overview table in **Figure 1A** of the  
408 dataset describes, for each cancer cohort, the sample size, female patient percentage,  
409 average cancer onset age, and tumor stage. Samples are normalized by their median  
410 absolute deviations (MAD), so that the MAD of all samples in the dataset is 1. Protein  
411 markers with high fractions (greater than 20%) of missing values are filtered out. For the  
412 corresponding RNA-seq data, we used the log<sub>2</sub> normalization on the FPKM (fragments  
413 per kilobase of exon per million mapped fragments)-normalized RNA-seq counts and  
414 genes have no expression in at least 90% of the samples were filter out.

415  
416 The proteomics data used for validation were downloaded from the NCI CPTAC portal.  
417 The dataset overview table in **Figure S1A** describe for each cancer cohort, the sample  
418 size, female patient percentage, average cancer onset age, and tumor stage. The  
419 validation data are processed in the same way as prospective data. The RNA-seq data  
420 sets of the three retrospective CPTAC cohorts were downloaded from the NCI CPTAC  
421 DCC portal. The RNA expression was measured in FPKM and was further normalized by  
422 log<sub>2</sub>(FPKM+1).

### 423 424 **pQTL and eQTL identification**

425  
426 For each cancer cohort, we identified pQTLs and eQTLs using the multiple linear  
427 regression model as implemented in the “limma” R package. We also corrected

428 confounding factors including age, gender, ethnicity, and TMT batch. The false discovery  
429 rate (FDR) was corrected from the p-values with the Benjamini-Hochberg procedure,  
430 ensuring that the identified QTLs are statistically robust. Somatic mutations are grouped  
431 at a gene level in the multiple regression model, similar to that implemented by our  
432 previously developed AeQTL tool<sup>7</sup>. Mutations separated are analyzed by their  
433 mechanisms of action, including nonsynonymous mutations as controls that likely do not  
434 affect expression, missense mutations, and truncating mutations including frameshift and  
435 in-frame indels, nonsense, splice site, and translation start site mutations. To improve  
436 statistical power, we focused our analysis on genes with three or more mutations in each  
437 cancer cohort and analyzed associations of mutations affecting *cis*-expression of the  
438 corresponding mRNA or protein products.

439

#### 440 **spsQTL identification**

441

442 We combined two complementary statistical methods to identify spsQTLs. In the first  
443 method adopted from Battle et al.<sup>4</sup>, we compared the following two nested linear models  
444 using likelihood ratio test (LRT) with the “anova” function in R:

445

$$\square = \square + \square_0 \square + \square_1 \square$$

446

$$\square = \square + \square_2 \square$$

447

448 where  $\square$  is the genotype,  $\square$  represents RNA level, and  $p$  is the protein level. By  
449 comparing these models using LRT and filtering results with an FDR less than 0.05, we  
450 identified candidate spsQTLs where the genotype (mutation) has a disproportionate  
451 impact on protein abundance independent of mRNA expression.

452

453 In the second method adopted from Mirauta et al.<sup>22</sup>, we selected QTLs where the spQTL  
454 FDR was less than 0.05 but the corresponding seQTL FDR was greater than 0.05 as  
455 candidate spsQTLs, to specifically identify mutations that affect protein levels without  
456 influencing mRNA. We then overlapped these two lists of candidate spsQTLs obtained  
457 from two complementary methods to identify the final list of spsQTLs for downstream  
458 analyses.

459

#### 460 **mRNA-Protein correlation:**



461 To investigate the impact of mutations on mRNA and protein expression, we performed  
462 a comparative analysis across the six solid cancer types. For each cancer type, Pearson  
463 correlation coefficients were calculated for individual genes using paired mRNA and  
464 protein expression data. We analyzed three groups of genes we identified as showing  
465 variable impact on mRNA/protein level expressions: Concordant genes (with mutations  
466 showing concordant effects at both mRNA and protein levels in cis), Discordant genes  
467 (showing protein-specific effects), and Other genes (showing no concordant or protein-  
468 specific impact). Our aim was to test the hypothesis whether the mRNA-protein  
469 correlations of the Concordant/Discordant groups differed from the baseline genome-  
470 wide mRNA-protein correlations, indicating biological significance. To assess this, we  
471 employed two-sample Wilcoxon rank-sum test, comparing the mRNA-protein correlations  
472 for the Concordant/Discordant and Other gene groups within each cancer type. Pairwise  
473 comparisons were made between the Concordant and Other gene sets, as well as  
474 between the Discordant and Other gene sets, demonstrating that the correlation  
475 coefficients for these groups were drawn from distinct population distributions with  
476 statistical significance at a p-value threshold of 0.05.

477

#### 478 **Tumor-normal differential expression analysis**

479 We conducted this analysis in the prospective CPTAC cohorts with paired tumor-adjacent  
480 tissue normal samples. For each cancer cohort, we paired the tumor and normal samples  
481 from the same patient and performed a differential protein/mRNA expression analysis to  
482 identify differentially expressed proteins with “limma” package. Demographic factors and  
483 batch effects, including age, ethnicity, race, and sequencing operator are adjusted in the  
484 multiple regression model.

485

486

#### 487 **Supplementary Tables**

488

489 **Table S1. List of expression quantitative trait loci (eQTLs) identified across 6 cancer**  
490 **types.** This table provides details on the gene mutations associated with mRNA expression  
491 levels, including statistical test results, mutation type, p-values (adjusted), and effect sizes.

492

493 **Table S2. List of protein quantitative trait loci (pQTLs) identified across 6 cancer types.**

494 This table provides details on the gene mutations associated with protein abundance levels,  
495 including statistical test results, mutation type, p-values (adjusted), and effect sizes.

496

497 **Table S3. Concordant expression and protein quantitative trait loci (eQTLs and pQTLs)**

498 **identified across 6 cancer types.** This table includes information on the gene mutations,  
499 identified cancer types, and their impact on both mRNA and protein expression levels,  
500 demonstrating loci with consistent effects across both molecular layers.

501

502 **Table S4. Significant somatic protein-specific QTLs (spsQTLs) identified by our**  
503 **statistical pipeline across six cancer types.** This table details the loci with mutations

504 showing significant impacts on protein abundance not explained by mRNA levels, including  
505 summary statistics for eQTL/pQTL tests and the LRT and overlap test results.

506

507 **Table S5. Summary statistics for differentially expressed proteins (DEPs) identified in**  
508 **paired tumor-normal (TN) samples across six cancer types.** This table includes the test

509 statistics of protein expression differences between tumor and normal tissues harboring the  
510 specific mutation.

511

512 **Table S6. Test statistics between the three groups of TP53 mutations.** The tested groups

513 were defined by (1) TP53 mutations with top 20% mRNA (left) or protein (right) expression in  
514 the prospective CPTAC cohorts, (2) the other TP53 mutations observed across all CPTAC  
515 samples, and (3) the rest of the assayed TP53 mutations from Kotler et al. using *TP53*  
516 functional scores from Kotler et al.

517

518 **Table S7. Pearson's correlation coefficient tests between paired mRNA and protein**  
519 **expressions for each concordant and discordant gene, within each cancer cohort.**

520

521 **Supplementary Figures**

522

523 **Supplementary Figure 1. Overview of the retrospective cohorts** (A) Summary of the  
524 retrospective CPTAC proteogenomic cohorts used for the discovery analyses, including  
525 cancer type abbreviation, data source, sample size of tumor (T) and normal (N) tissues,  
526 female percentage, average onset age in years, and tumor stage distribution. (B) Volcano  
527 plots showing seQTLs associations in the six cancer types (left) and volcano plots  
528 showing spQTLs associations (right), where each dot denotes a gene-cancer pair  
529 included in the analysis. Top associated genes were further labeled. FC: log fold change.  
530 FDR: false discovery rate.

531  
532 **Supplementary Figure 2. spsQTLs with strong effects.** (A) Examples of spsQTL  
533 whose effect sizes in mRNA level and protein level are in different direction. For each  
534 gene, the plot on the left shows the corresponding mRNA levels of mutation carriers vs.  
535 non-carriers in FPKM, and the plot on the right shows protein level comparison in log ratio  
536 (MS TMT measurements) in the respective cancer type labeled on top of each of the violin  
537 plots. The labeled mutations are the three mutations whose carriers show the highest  
538 absolute expression differences of the mutated gene product compared to the non-  
539 carriers. (B) Examples of spsQTL with a protein logFC and mRNA logFC ratio greater  
540 than 15

541  
542 **Supplementary Figure 3. Overlapped of significant QTLs in cross-tumor analysis**  
543 **and matched tumor-normal analysis projected onto pQTL volcano plots based on**  
544 **cross-tumor analyses.** The plots were made separately for (A) discordant spsQTLs, and  
545 (B) concordant eQTL/pQTLs.

546  
547 **Supplementary Figure 4. Example lollipop plots showing mutations for two genes that**  
548 **were identified as spsQTLs, including YLPM1 and ESRP1.** The number on each disc  
549 denotes the number of mutations in that position and the color of the disc represents the  
550 mutation type.

551  
552 **Supplementary Figure 5. Correlation coefficients of Concordant vs. Discordant**  
553 **genes.** The violin plots depict the distribution of correlation coefficients between matched

554 mRNA and protein expressions for Concordant (blue), Discordant (red), and Other genes  
555 (gray) across the six cancer types studied. Genes with notable correlations are labeled in  
556 each plot.

557

558

## 559 **DATA AND SOFTWARE AVAILABILITY**

### 560 **Data Availability**

561 Proteomic data for CPTAC-2/3 cohorts can be found on National Cancer Institute (NCI)  
562 Proteomic Data Commons (PDC): <https://cptac-data-portal.georgetown.edu/cptacPublic/>.

563 The studies used in the discovery cohorts and their PDC study IDs are: BRCA  
564 (PDC000120), CRC (PDC000116), CCRCC (PDC000127), LUAD (PDC000153), OV  
565 (JHU: PDC000110; PNNL: PDC000118), UCEC (PDC000125)

566 The studies used in the validation cohorts and their PDC study IDs are: BRCA  
567 (PDC000173), CRC (PDC000111), OV (JHU: PDC000113; PNNL: PDC000114)

568 Genomic data, including DNA mutation and transcriptome profiling for all CPTAC-2/3  
569 cohorts used herein can be found on National Cancer Institute (NCI) Genome Data  
570 Commons (GDC): <https://portal.gdc.cancer.gov/projects/CPTAC-2> (dbGaP Study  
571 Accession #: phs000892) and <https://portal.gdc.cancer.gov/projects/CPTAC-3> (dbGaP  
572 Study Accession #: phs001287)

573 Data for TP53 MAVE assays can be downloaded from the Supplementary Information  
574 from Kotler et al<sup>21</sup>.

575

### 576 **Code Availability**

577 The source code used for all analyses in this article is available at  
578 <https://github.com/Huang-lab/pQTL> under an MIT license.

## 579 **ACKNOWLEDGEMENTS**

580 The authors wish to acknowledge CPTAC and its participating patients and families that  
581 generously contributed the data. This work was supported by NIH NIGMS  
582 R35GM138113, ACS RSG-22-115-01-DMC, and Mount Sinai funds to KH.

583 **DECLARATION OF INTERESTS**

584 K.H. is a co-founder and board member of a non-for-profit 501(c)(3) organization, Open  
585 Box Science, from which he does not receive any compensation and pose no competing  
586 financial interests with this work. All authors declare no competing interests.

587 **CONTRIBUTIONS**

588 K.H. conceived the research; Y.L and K.H. designed the analyses. Y.L. and A.E.  
589 developed the software and conducted the bioinformatics analyses, A.E. curated and  
590 preprocessed the datasets. Y.L., A.E., and K.H. wrote the manuscript. K.H. supervised  
591 the study. All authors read, edited, and approved the manuscript.

592

593

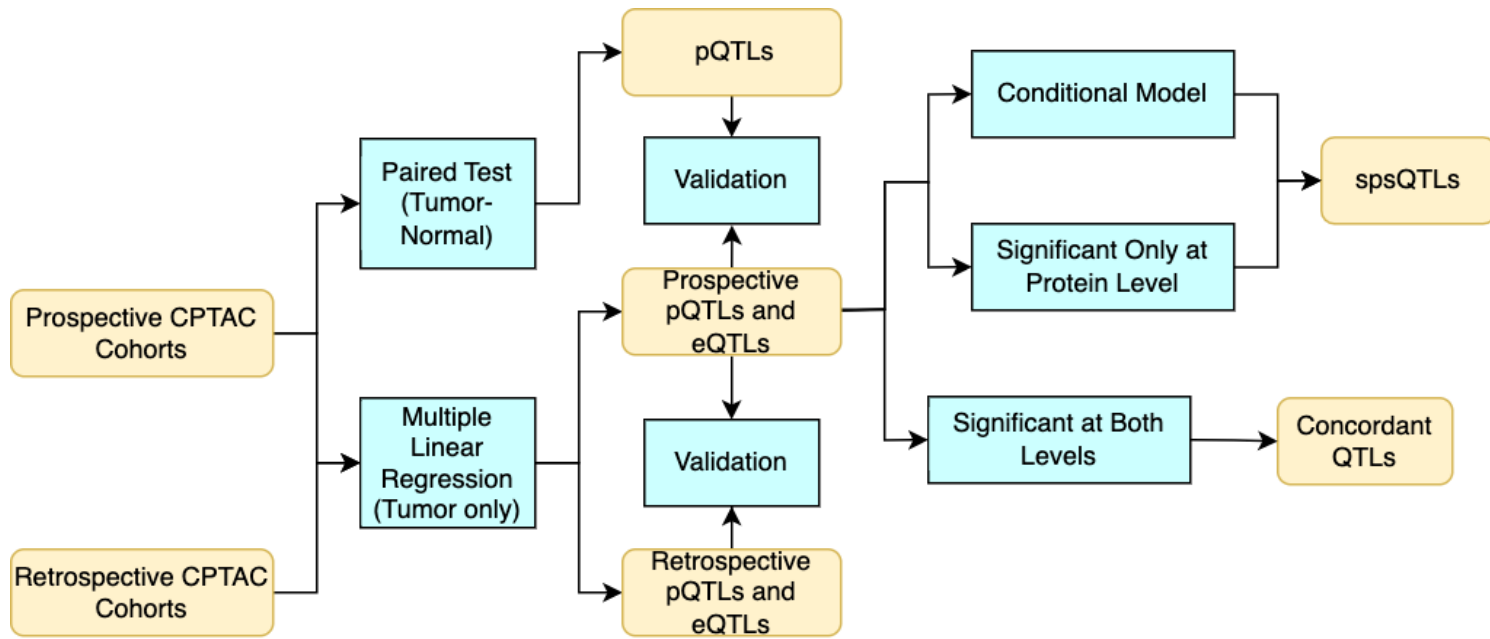
594 **REFERENCES**

- 595 1. Kurosaki, T., Popp, M. W. & Maquat, L. E. Quality and quantity control of gene  
596 expression by nonsense-mediated mRNA decay. *Nature Reviews Molecular Cell*  
597 *Biology* vol. 20 Preprint at <https://doi.org/10.1038/s41580-019-0126-2> (2019).
- 598 2. Wang, Z. *et al.* Non-cancer-related pathogenic germline variants and expression  
599 consequences in ten-thousand cancer genomes. *Genome Med* **13**, (2021).
- 600 3. Lindeboom, R. G. H., Supek, F. & Lehner, B. The rules and impact of nonsense-  
601 mediated mRNA decay in human cancers. *Nat Genet* **48**, (2016).
- 602 4. Battle, A. *et al.* Impact of regulatory variation from RNA to protein. *Science (1979)*  
603 **347**, (2015).
- 604 5. Cenik, C. *et al.* Integrative analysis of RNA, translation, and protein levels reveals  
605 distinct regulatory variation across humans. *Genome Res* **25**, (2015).
- 606 6. Chick, J. M. *et al.* Defining the consequences of genetic variation on a proteome-wide  
607 scale. *Nature* **534**, (2016).
- 608 7. Dong, G., Wendl, M. C., Zhang, B., Ding, L. & Huang, K. L. AeQTL: eQTL analysis  
609 using region-based aggregation of rare genomic variants. *Pac Symp Biocomput* **26**,  
610 (2021).
- 611 8. Rabadán, R. *et al.* Identification of relevant genetic alterations in cancer using  
612 topological data analysis. *Nat Commun* **11**, (2020).
- 613 9. Ding, J. *et al.* Systematic analysis of somatic mutations impacting gene expression in  
614 12 tumour types. *Nat Commun* **6**, (2015).
- 615 10. Arad, G. & Geiger, T. Functional impact of protein-RNA variation in clinical cancer  
616 analyses. *Molecular & Cellular Proteomics* 100587 (2023)  
617 doi:10.1016/J.MCPRO.2023.100587.
- 618 11. Mertins, P. *et al.* Proteogenomics connects somatic mutations to signalling in breast  
619 cancer. *Nature* **534**, (2016).

- 620 12. Zhang, H. *et al.* Integrated Proteogenomic Characterization of Human High-Grade  
621 Serous Ovarian Cancer. *Cell* **166**, (2016).
- 622 13. Zhang, B. *et al.* Proteogenomic characterization of human colon and rectal cancer.  
623 *Nature* **513**, (2014).
- 624 14. Liu, Y., Beyer, A. & Aebersold, R. On the Dependency of Cellular Protein Levels on  
625 mRNA Abundance. *Cell* vol. 165 Preprint at <https://doi.org/10.1016/j.cell.2016.03.014>  
626 (2016).
- 627 15. Clark, D. J. *et al.* Integrated Proteogenomic Characterization of Clear Cell Renal Cell  
628 Carcinoma. *Cell* **179**, (2019).
- 629 16. Vasaikar, S. *et al.* Proteogenomic Analysis of Human Colon Cancer Reveals New  
630 Therapeutic Opportunities. *Cell* **177**, (2019).
- 631 17. Gillette, M. A. *et al.* Proteogenomic Characterization Reveals Therapeutic  
632 Vulnerabilities in Lung Adenocarcinoma. *Cell* **182**, (2020).
- 633 18. Dou, Y. *et al.* Proteogenomic Characterization of Endometrial Carcinoma. *Cell* **180**,  
634 (2020).
- 635 19. Krug, K. *et al.* Proteogenomic Landscape of Breast Cancer Tumorigenesis and  
636 Targeted Therapy. *Cell* **183**, (2020).
- 637 20. McDermott, J. E. *et al.* Proteogenomic Characterization of Ovarian HGSC Implicates  
638 Mitotic Kinases, Replication Stress in Observed Chromosomal Instability. *Cell Rep*  
639 *Med* **1**, (2020).
- 640 21. Kotler, E. *et al.* A Systematic p53 Mutation Library Links Differential Functional  
641 Impact to Cancer Mutation Pattern and Evolutionary Conservation. *Mol Cell* **71**,  
642 (2018).
- 643 22. Mirauta, B. A. *et al.* Population-scale proteome variation in human induced pluripotent  
644 stem cells. *Elife* **9**, (2020).
- 645 23. Karousis, E. D. & Mühlemann, O. Nonsense-mediated mRNA decay begins where  
646 translation ends. *Cold Spring Harb Perspect Biol* **11**, (2019).
- 647 24. Davidoff, A. M., Humphrey, P. A., Dirk Iglehart, J. & Marks, J. R. Genetic basis for  
648 p53 overexpression in human breast cancer. *Proc Natl Acad Sci U S A* **88**, (1991).
- 649 25. Huang, K. lin *et al.* Spatially interacting phosphorylation sites and mutations in cancer.  
650 *Nat Commun* **12**, (2021).
- 651 26. Chabanon, R. M. *et al.* PBRM1 deficiency confers synthetic lethality to DNA repair  
652 inhibitors in cancer. *Cancer Res* **81**, (2021).
- 653 27. Vadlamudi, Y., Dey, D. K. & Kang, S. C. Emerging Multi-cancer Regulatory Role of  
654 ESRP1: Orchestration of Alternative Splicing to Control EMT. *Curr Cancer Drug*  
655 *Targets* **20**, (2020).
- 656 28. Armenia, J. *et al.* The long tail of oncogenic drivers in prostate cancer. *Nat Genet* **50**,  
657 (2018).
- 658 29. Loganathan, S. K. *et al.* Rare driver mutations in head and neck squamous cell  
659 carcinomas converge on NOTCH signaling. *Science* **367**, (2020).
- 660 30. Philpott, C., Tovell, H., Frayling, I. M., Cooper, D. N. & Upadhyaya, M. The NF1  
661 somatic mutational landscape in sporadic human cancers. *Human Genomics* vol. 11  
662 Preprint at <https://doi.org/10.1186/s40246-017-0109-3> (2017).
- 663 31. Xue, Z. *et al.* MAP3K1 and MAP2K4 mutations are associated with sensitivity to  
664 MEK inhibitors in multiple cancer models. *Cell Res* **28**, (2018).

- 665  
666  
667  
668  
669  
670  
671  
672  
673  
674
32. Buccitelli, C. & Selbach, M. mRNAs, proteins and the emerging principles of gene expression control. *Nature Reviews Genetics* vol. 21 Preprint at <https://doi.org/10.1038/s41576-020-0258-4> (2020).
  33. Edwards, N. J. *et al.* The CPTAC data portal: A resource for cancer proteomics research. *J Proteome Res* **14**, (2015).
  34. Kuang, D. *et al.* MaveRegistry: a collaboration platform for multiplexed assays of variant effect. *Bioinformatics* **37**, (2021).
  35. Elmas, A. *et al.* Pan-cancer proteogenomic investigations identify post-transcriptional kinase targets. *Commun Biol* **4**, (2021).

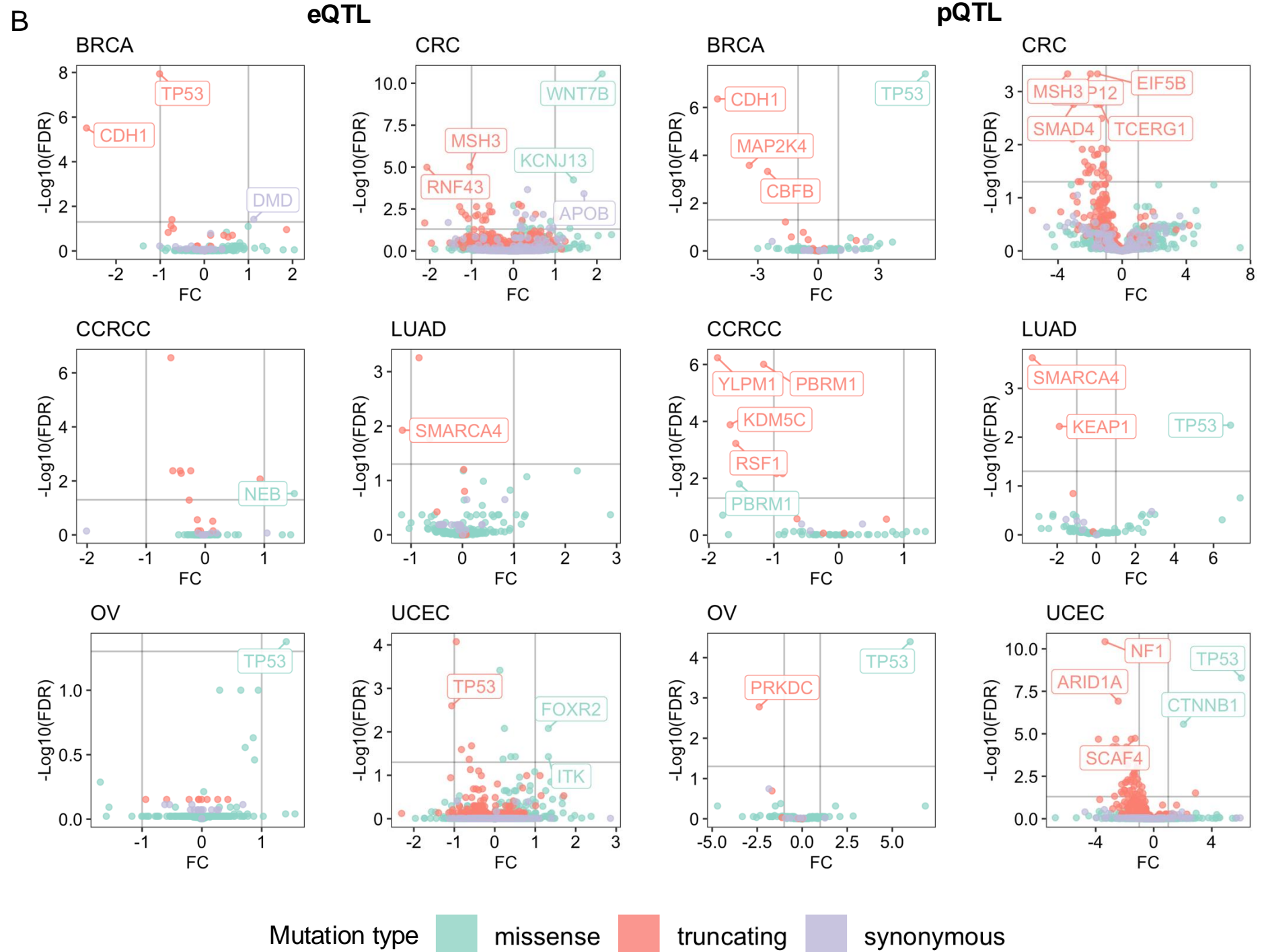
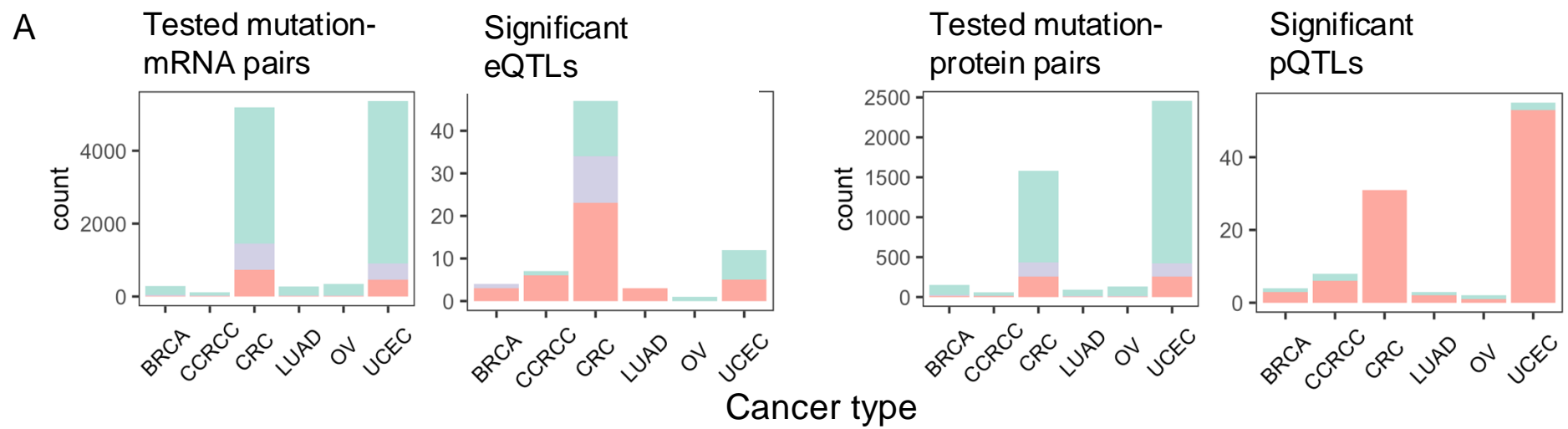
A



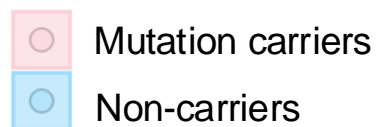
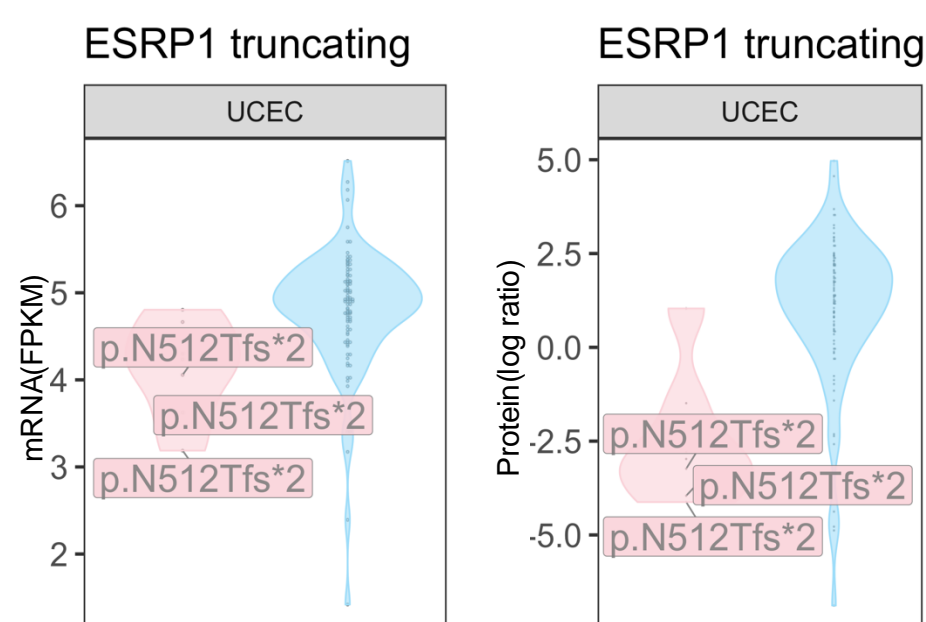
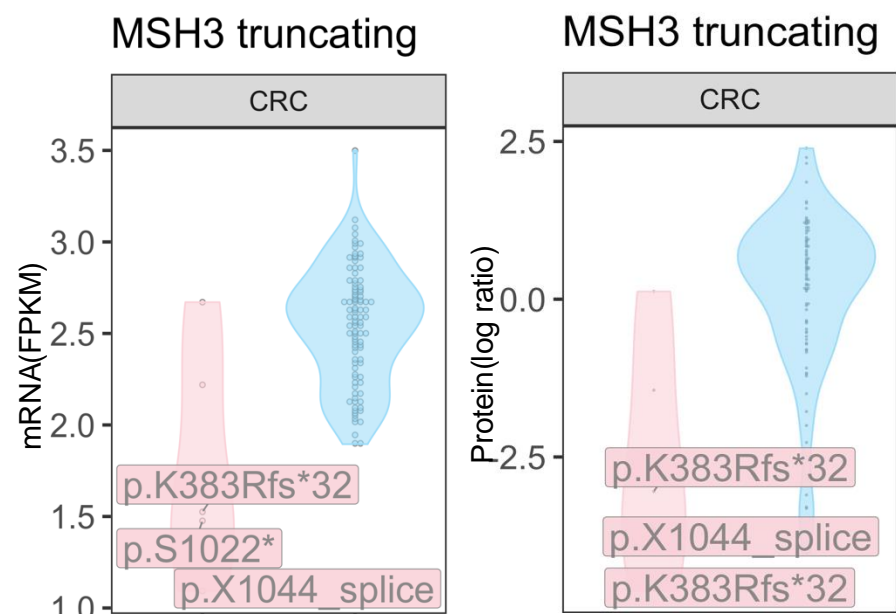
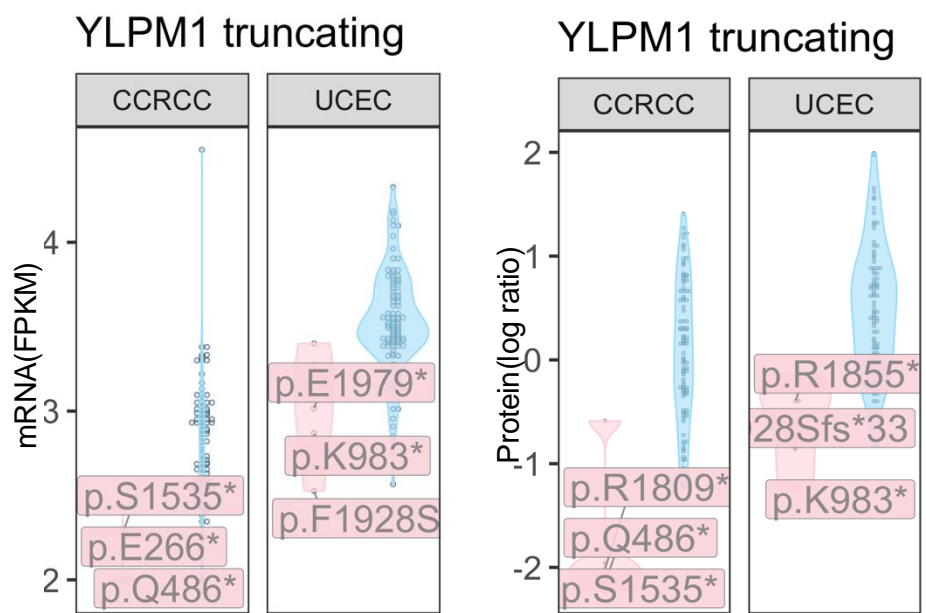
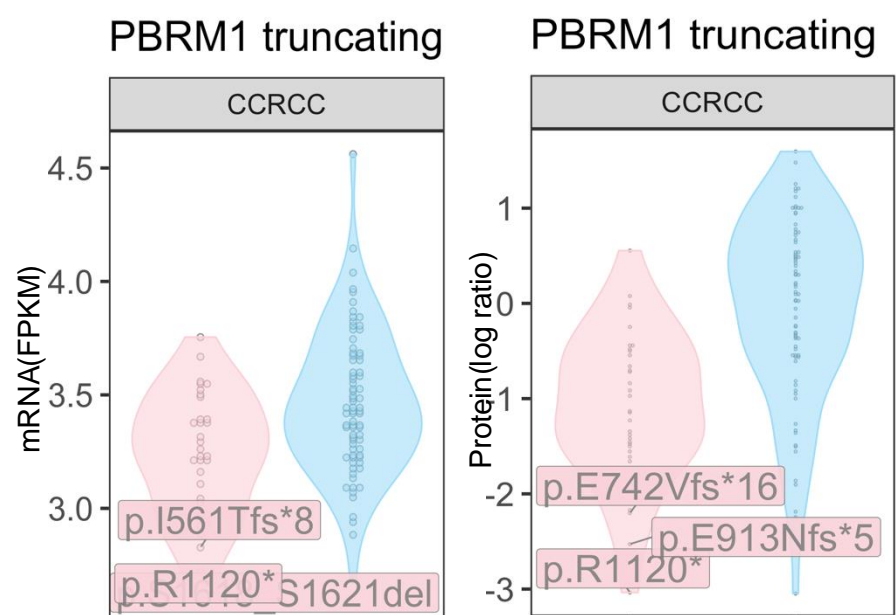
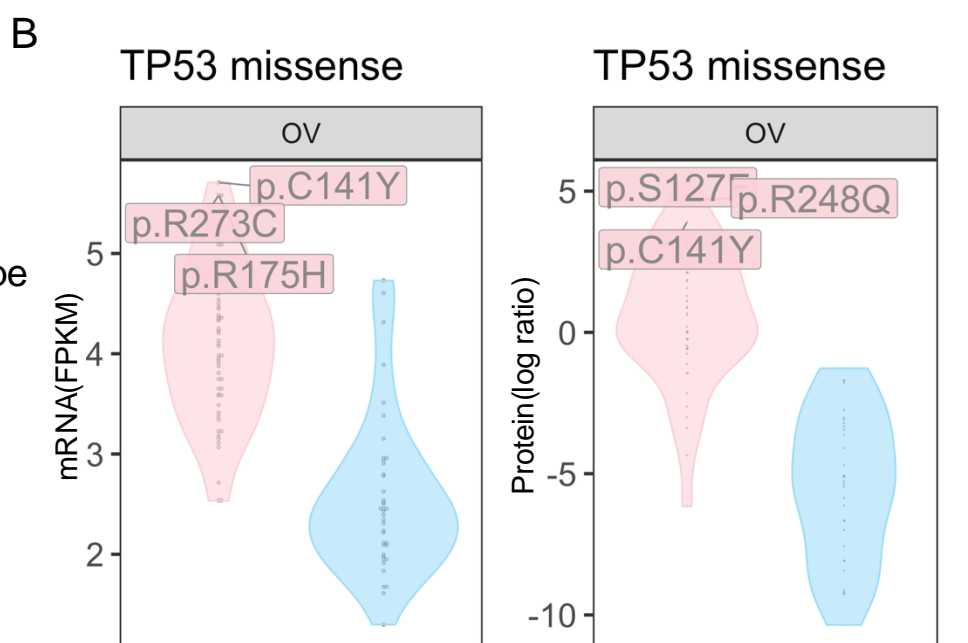
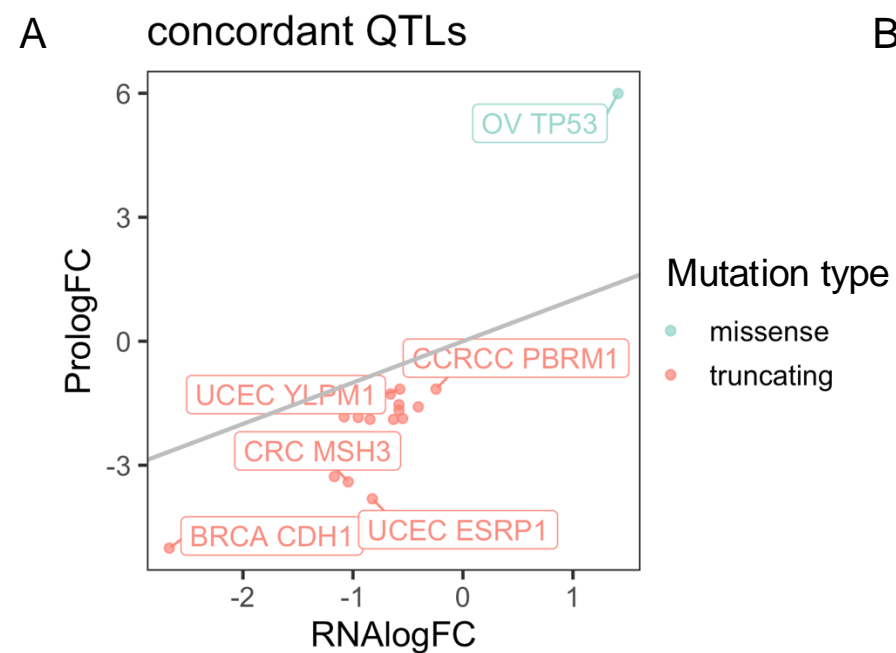
B

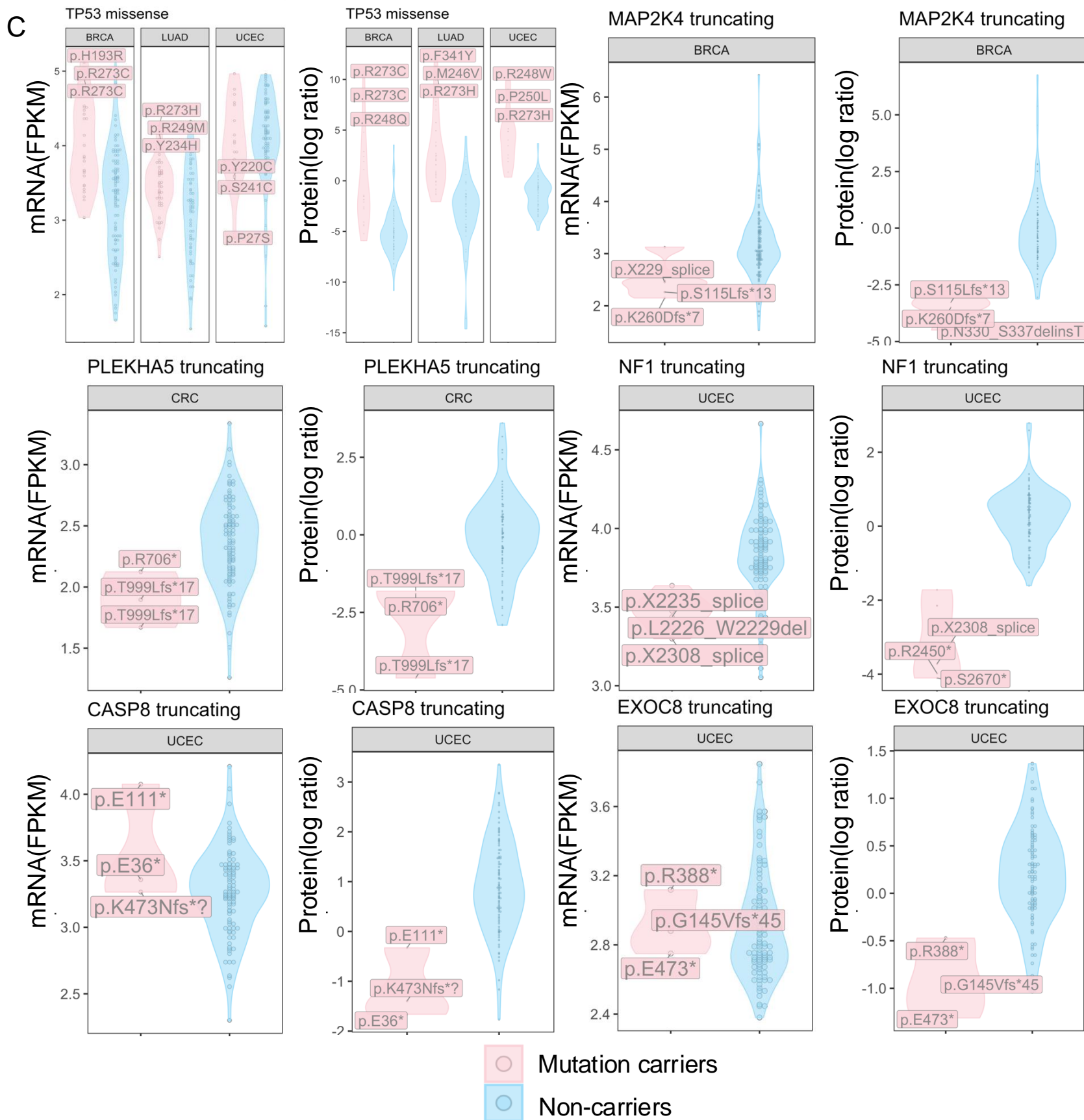
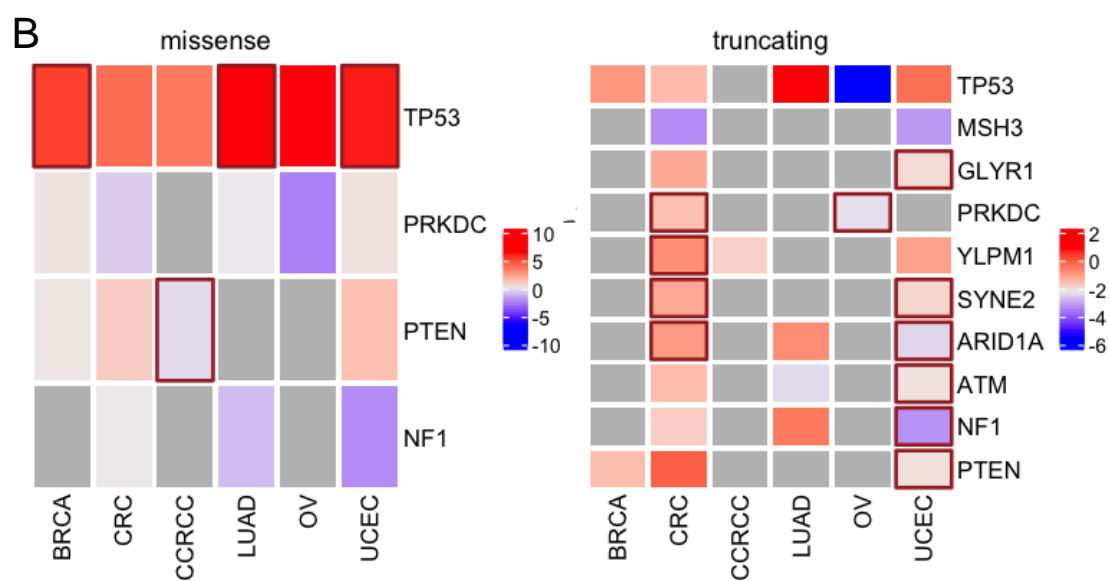
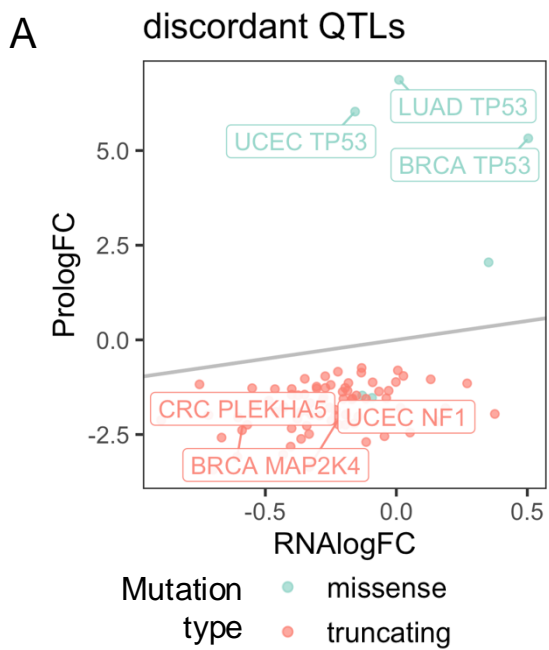
Cancer Type	Breast Cancer	Clear Cell Renal Cell Carcinoma	Colorectal Cancer	Lung Adenocarcinoma	Ovarian Cancer	Uterine Corpus Endometrial Carcinoma
Abbreviation	<b>BRCA</b>	<b>CCRCC</b>	<b>CRC</b>	<b>LUAD</b>	<b>OV</b>	<b>UCEC</b>
Data Source	Krug et al. 2020 (PMID: 33212010)	Clark et al. 2019 (PMID: 31675502)	Vasaikar et al. 2019 (PMID: 31031003)	Gillette et al. 2020 (PMID: 32649874)	McDermott et al. 2020 (PMID: 32529193)	Dou et al. 2020 (PMID: 32059776)
Sample Size (Tumors/Normals)	T: 115 N: 18	T: 110 N: 84	T: 95 N: 100	T: 109 N: 102	T: 84 N: 19	T: 97 N: 20
Female %	100%	25.2%	57.4%	34.6%	100%	100%
Average Onset (yr)	60.4	60.6	65.2	62.7	59.1	63.7
Tumor Stage	1: 3% 2: 60.4% 3: 26.9% NA: 9.7%	1: 41.8% 2: 15.5% 3: 34.5% 4: 8.2%	1: 10.2% 2: 40.6% 3: 41.1% 4: 8.1%	1: 53.5% 2: 27.5% 3: 18.5% 4: 0.5%	1: 1% 2: 1% 3: 72.8% 4: 15.5% NA: 9.7%	1: 76.1% 2: 6.8% 3: 14.5% 4: 2.6%

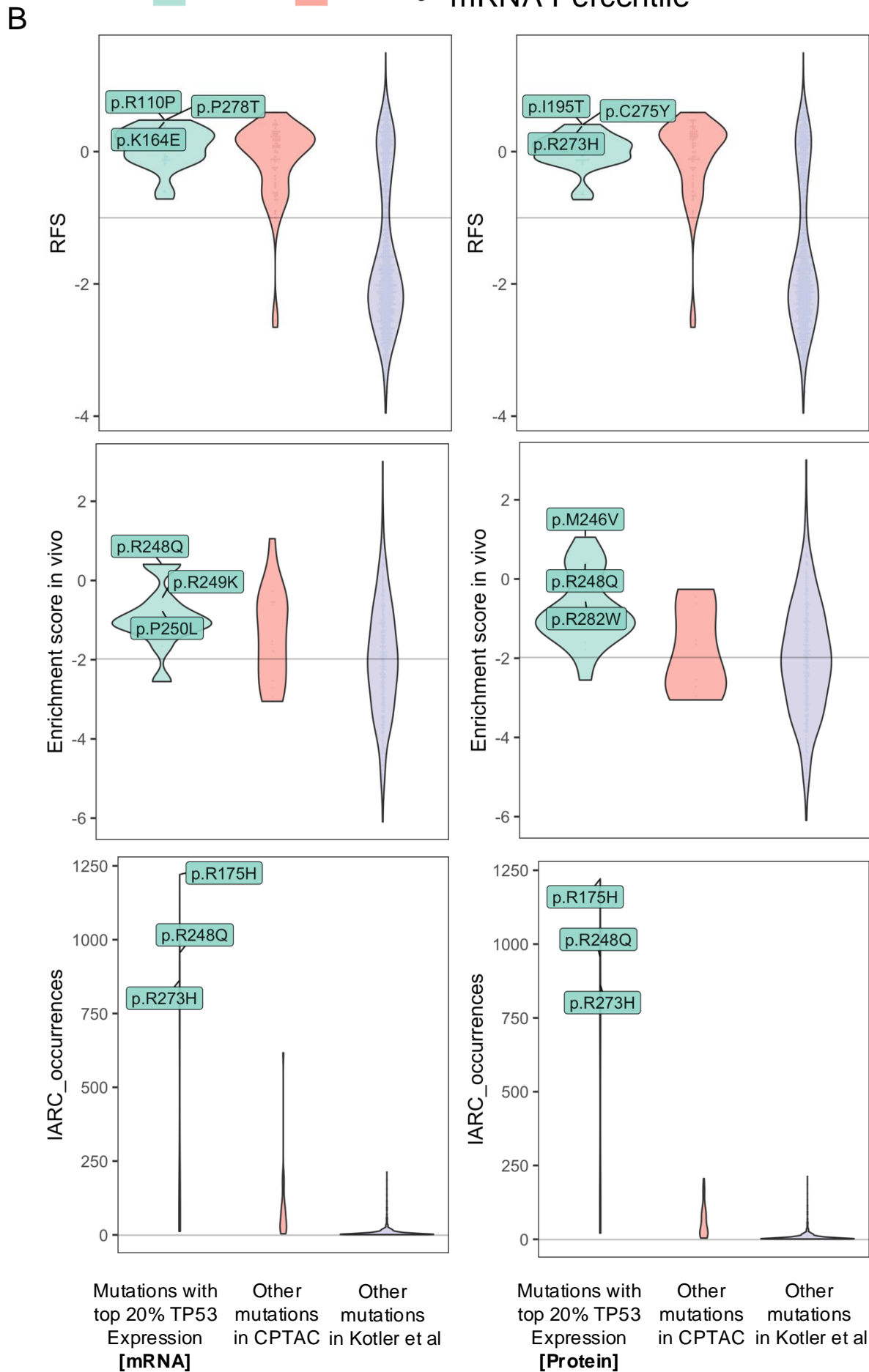
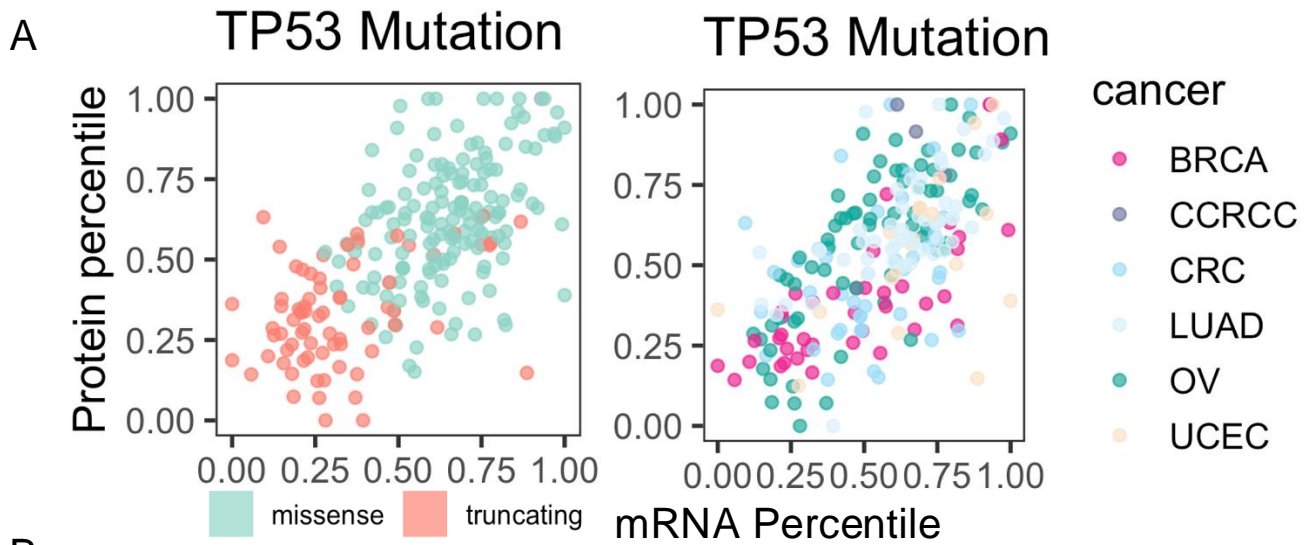




**C**

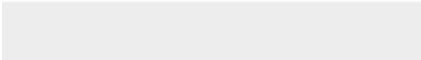








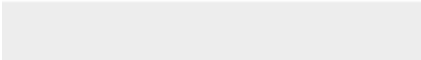



Click here to access/download  
**Supplementary Material**  
SuppFigures.pdf





Click here to access/download  
**Supplementary Material**  
TableS1.eQTLs.xlsx





Click here to access/download  
**Supplementary Material**  
TableS2.pQTLs.xlsx





Click here to access/download  
**Supplementary Material**  
TableS3.concordant\_e.pQTLs.xlsx







Click here to access/download  
**Supplementary Material**  
TableS4.significant\_spsQTLs.xlsx





Click here to access/download  
**Supplementary Material**  
TableS5.DEP\_pairedTN\_stats.xlsx

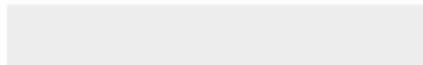




[Click here to access/download](#)

**Supplementary Material**

[TableS6.TP53\\_mutation\\_test\\_statistics.xlsx](#)





[Click here to access/download](#)

**Supplementary Material**

[TableS7.CorCoefs.ConcordantAndDiscordant.xlsx](#)





Click here to access/download  
**Supplementary Material**  
point-by-point-response.docx





Kuan-lin Huang, PhD  
Assistant Professor, Department of Genetics and Genomic Sciences  
Institute of Genomics and Multiscale Biology  
Icahn School of Medicine at Mount Sinai

1399 Park Avenue (Room 4-420C)  
Box 1498  
New York, NY 10029  
Phone: (212) 824-6134  
Email: [kuan-lin.huang@mssm.edu](mailto:kuan-lin.huang@mssm.edu)  
Web: [ComputationalOmicsLab.org](http://ComputationalOmicsLab.org)

Sep 13<sup>th</sup> 2024

Hans Zauner  
Editor, GigaScience

Dear Dr. Zauner,

We would like to express our sincere gratitude to you and the reviewers for providing valuable feedback on our manuscript titled "*Mutation Impact on mRNA Versus Protein Expression across Human Cancers*" (GIGA-D-24-00168) along with point by point response to each of the reviewer comments.

Following your suggestion and the reviewers comments, we have conducted a multitude of analyses and improvements that have significantly strengthened the manuscript. Here are the key improvements in this revised version:

- **Expanded Methods and Metrics:** We have added more detailed explanations of the statistical metrics used in our analysis. Additionally, we have discussed alternative approaches and outlined the limitations of our current methods.
- **Improved Figures:** In response to the request for clarity, we have enlarged Figure 1 and corrected the truncated legend for Figure 3, enhancing the overall presentation of the data.
- **Additional Analysis on mRNA-Protein Correlation:** We have addressed concerns regarding the low correlation between mRNA and protein expression by including additional statistical analysis. These new results highlight the variation in mRNA-protein correlations across genes and cancer types for concordant/discordant eQTL/pQTLs (**Figure S5, Table S7**), along with an expanded discussion on how these findings stress the need to consider protein-specific impacts of mutations.

Additionally, we have carefully edited the manuscript, where you can find a tracked changes version attached. We believe the revised manuscript adequately addresses all the reviewers' concerns and hope you will find it to be satisfactory for publication.



Kuan-lin Huang, PhD  
Assistant Professor, Department of Genetics and Genomic Sciences  
Institute of Genomics and Multiscale Biology  
Icahn School of Medicine at Mount Sinai

1399 Park Avenue (Room 4-420C)  
Box 1498  
New York, NY 10029  
Phone: (212) 824-6134  
Email: [kuan-lin.huang@mssm.edu](mailto:kuan-lin.huang@mssm.edu)  
Web: [ComputationalOmicsLab.org](http://ComputationalOmicsLab.org)

Sincerely and on behalf of the team,

Kuan-lin Huang, Ph.D.  
Assistant Professor of Genetics and Genomic Sciences & Artificial Intelligence and Human Health  
Icahn School of Medicine at Mount Sinai  
New York, NY 10029