

1 **Genomic determinants of protein abundance**

2 **variation in colorectal cancer cells**

3 Theodoros I. Roumeliotis^{1,*,#}, Steven Paul Williams^{1,*}, Emanuel Gonçalves²,
4 Fatemeh Zamanzad Ghavidel², Nanne Aben⁴, Magali Michaut⁴, Michael Schubert²,
5 James C. Wright¹, Mi Yang³, Clara Alsinet¹, Rodrigo Dienstmann^{6,7}, Justin Guinney⁶,
6 Pedro Beltrao², Alvis Brazma², Oliver Stegle², David J. Adams¹, Lodewyk
7 Wessels^{4,5,†}, Julio Saez-Rodriguez^{2,3,†}, Ultan McDermott^{1,†}, Jyoti S. Choudhary^{1,†,#}

8 ¹Wellcome Trust Sanger Institute, Wellcome Trust Genome Campus, Hinxton, Cambridgeshire CB10 1SA, UK.

9 ²European Molecular Biology Laboratory, European Bioinformatics Institute, Wellcome Trust Genome Campus, Cambridge
10 CB10 1SA, UK.

11 ³RWTH Aachen University, Faculty of Medicine, Joint Research Center for Computational Biomedicine, Aachen, Germany.

12 ⁴Computational Cancer Biology, The Netherlands Cancer Institute, Amsterdam, The Netherlands.

13 ⁵Faculty of EEMCS, Delft University of Technology, Delft, The Netherlands.

14 ⁶Computational Oncology, Sage Bionetworks, Fred Hutchinson Cancer Research Center, Seattle, USA.

15 ⁷Oncology Data Science, Vall d'Hebron Institute of Oncology, Barcelona, Spain.

16

17 ***Co-first author**

18 **†Co-senior author**

19 **#Corresponding author**

20 `jc4@sanger.ac.uk` (J.S.C.),

21 `tr6@sanger.ac.uk` (T.I.R.) (lead contact)

22

23 **Summary**

24 Assessing the extent to which genomic alterations compromise the integrity of the
25 proteome is fundamental in identifying the mechanisms that shape cancer
26 heterogeneity. We have used isobaric labelling and tribrid mass spectrometry to
27 characterize the proteomic landscapes of 50 colorectal cancer cell lines and to
28 decipher the relationships between genomic and proteomic variation. The robust
29 quantification of 12,000 proteins and 27,000 phosphopeptides revealed how protein
30 symbiosis translates to a co-variome which is subjected to a hierarchical order and
31 exposes the collateral effects of somatic mutations on protein complexes. Targeted
32 depletion of key chromatin modifiers confirmed the transmission of variation and the
33 directionality as characteristics of protein interactions. Protein level variation was
34 leveraged to build drug response predictive models towards a better understanding
35 of pharmacoproteomic interactions in colorectal cancer. Overall, we provide a deep
36 integrative view of the molecular structure underlying the variation of colorectal
37 cancer cells.

38

39 **Keywords**

40 Proteomics, Isobaric labelling, protein complexes, co-variation networks, mutations,
41 colorectal cancer, cell lines, drug response

42

43 **Highlights**

- 44 • The cancer cell functional “co-variome” is a strong attribute of the proteome.
- 45 • Mutations can have a direct impact on protein levels of chromatin modifiers.
- 46 • Transmission of genomic variation is a characteristic of protein interactions.
- 47 • Pharmacoproteomic models are strong predictors of response to DNA
- 48 damaging agents.

49 **Abbreviations**

50 COREAD, Colorectal Adenocarcinoma

51 IMAC, Immobilized Metal ion Affinity Chromatography

52 ROC, Receiver Operating Characteristic

53 AUC, Area Under the Curve

54 WGCNA, Weighted Correlation Network Analysis

55 CNA, Copy Number Alteration

56 SOM, Self-Organizing Map

57 QTL, Quantitative Trait Loci

58 MSI, Microsatellite Instability

59 CPS, Colorectal Proteomic Subtypes

60 **Introduction**

61 Tumours exhibit a high degree of molecular and cellular heterogeneity due to the
62 impact of genomic aberrations on protein networks underlying physiological cellular
63 activities. Modern mass spectrometry based proteomic technologies have now the
64 capacity to perform highly reliable analytical measurements of proteins in large sizes
65 of subjects and analytes providing a powerful tool in the quest for regulatory
66 associations between genomic features, gene expression patterns, protein networks
67 and phenotypic traits (Mertins et al., 2016; Zhang et al., 2014; Zhang et al., 2016).
68 However, understanding how genomic variation leads to variable proteomic
69 landscapes and distinct cellular phenotypes remains challenging due to the
70 enormous diversity in the biological characteristics of proteins. Studying protein co-
71 regulation holds the promise to overcome the challenges associated with molecular
72 complexity and is now gaining ground in the study of molecular networks as it can
73 efficiently predict gene functions (Stefely et al., 2016; Wang et al., 2016). Colorectal
74 cancer cell lines are widely used as a model that approximates cancer behaviour in a
75 variety of cellular and biochemical assays however their proteome based
76 characteristics and the genomic factors underlying protein variation remain largely
77 unexplored.

78 Here we leveraged the accurate quantification of a 12k total proteome obtained by
79 the application of isobaric labelling and tribrid mass spectrometry analysis on a panel
80 of 50 colorectal cancer cell lines, first to build *de novo* proteome-wide
81 representations of biological functions inferred by protein co-variation, highly
82 predictive of protein complexes and interactions, and second to rationalize the

83 impact of genomic variation in the context of the cancer cell co-variation protein
84 interaction network. We selected to study the colorectal cancer cell lines panel as it
85 has been extensively characterised by whole exome sequencing, gene expression,
86 copy number and methylation arrays, and the frequency of molecular alterations is
87 similar to that seen in clinical cohorts (Iorio et al., 2016). The cancer cell functional
88 “co-variome” appeared to be a strong attribute of the proteome, revealing the
89 interdependencies of protein complexes and assigning putative functions to
90 uncharacterized gene products. Additionally, our proteomics data were
91 complemented by protein phosphorylation measurements encompassing a total of
92 27,000 phosphorylated sequences which demonstrated that co-variation of
93 phosphorylation can also highlight known and novel biology. We assessed the direct
94 effects of mutations on protein abundances and we integrated these effects with the
95 cancer cell co-variome to uncover protein network vulnerabilities by the identification
96 of possible collateral effects on protein complexes. Proteomic analysis of human iPS
97 cells engineered with gene knockouts of key chromatin modifiers confirmed that
98 genomic variation can be transferred from directly affected proteins to tightly co-
99 regulated distant gene products through protein interactions. A significant number of
100 drug response predictive models were uniquely attributed to protein level variation
101 leading to a better understanding of pharmacoproteomic interactions in colorectal
102 cancer. Our results constitute a comprehensive in-depth resource elucidating the
103 molecular organization of colorectal cancer cells widely used in cancer research.

104

105 **Results**

106 **Proteome and phosphoproteome coverage**

107 To assess the extent of variation in protein and phosphorylation abundances within a
108 panel of 50 colorectal cancer cell lines (COREAD) we utilized isobaric peptide
109 labelling (TMT-10plex) and MS3 quantification (**Figure S1A**). We obtained relative
110 quantification between the different cell lines in a log₂ scale for 12,306 unique
111 proteins and 27,423 non-redundant phosphopeptides at FDR<1% (9,489 and 12,061
112 in at least half of the samples respectively) (**Figure S1B**) (**Table S1**, **Table S2**). The
113 average correlation between biological replicates was significantly higher than that of
114 non-replicates (Pearson's $r=0.74$ and -0.02 respectively, $p\text{-value}=5.1e-69$) which was
115 also observed in the inter-laboratory comparison with previously published TMT data
116 for six colorectal cancer cell lines (McAlister et al., 2014) (**Figure S1F**) confirming
117 that subtle differences between the cell lines can be detected using our proteomics
118 approach. Similar levels of global protein variation were observed across the cell
119 lines with an average standard deviation of 1.7-fold denoting highly variable
120 proteomes (**Figure S1C**). Correlation analysis between mRNA (publicly available
121 data) and protein relative abundances across the cell lines indicated a significant, yet
122 moderate concordance of the two molecular levels with an average Pearson
123 correlation $r=0.57$ for matched cell line mRNA and protein data, and $r=-0.015$ for
124 unrelated samples (Welch t-test, $p\text{-value} < 2.2e-16$) (**Figure S1G**). Overall, highly
125 variable mRNAs tend to correspond to highly variable proteins (Spearman's $r=0.62$)
126 although with a wide distribution (**Figure S1H**). Notably, several genes including

127 *TP53* displayed high variation at the protein level despite the low variation at the
128 mRNA level, suggesting a significant contribution of post-transcriptional regulatory
129 mechanisms to their total protein levels.

130 To identify genomic variants at the protein level we also searched our MS/MS
131 spectra against a customized protein database containing all amino acid
132 substitutions encoded by 77k missense mutations (Iorio et al., 2016). We identified
133 769 unique variant peptides mapped to 558 proteins (**Table S3**) (**Figure S1D**)
134 including several mutated cell differentiation markers (CD44, CD46, ITGAV, ITGB1,
135 ITGB4 and TFRC) that can be useful in targeted immunoaffinity based applications
136 for discrimination of cancer cells as well as a range of mutated protein complexes
137 (e.g spliceosomal, MCM, eIF3, CCT, Ksr1, Emerin). Two characteristic examples of
138 mutant peptides for KRAS and CTNNB1 colorectal cancer genes are depicted in
139 **Figure S1E**. Although the coverage in single amino acid variants is limited to
140 proteins with medium to higher expression, it is concordant with the overall variant
141 frequency and distinctly higher in hypermutated lines. Our proteogenomic search
142 reveals the unique proteotypes of the COREAD cell lines.

143

144 **The subunits of protein complexes maintain tight stoichiometry of total** 145 **abundance post-transcriptionally**

146 As a means of simplifying the complexity of protein abundance variation we
147 examined whether protein co-variation patterns detected across the cell lines could
148 help aggregate the thousands of single protein measurements into a much smaller

149 number of biologically meaningful clusters. The Spearman's correlation coefficients
150 between proteins with known physical interactions in protein complexes catalogued
151 in the CORUM database (Ruepp et al., 2010) was bimodal-like and clearly shifted to
152 positive values with mean 0.41. The respective distribution of all pairwise protein-to-
153 protein comparisons displayed a normal distribution with mean 0.09 (**Figure S2A, left**
154 **panel**). Specifically, 388 partially overlapping CORUM complexes, representing the
155 most significantly correlating set, showed a greater than 0.4 median correlation
156 between their components (**Table S4**). In contrast, the distribution of Spearman's
157 coefficients between CORUM pairs based on mRNA co-variation profiles was only
158 marginally shifted towards higher correlations (**Figure S2A, right panel**). This
159 indicates that the subunits of the complexes are tightly regulated post-
160 transcriptionally. Indeed, comparative Receiver Operating Characteristic (ROC)
161 curves showed that our proteomics data outperformed mRNA data in recapitulating
162 protein complexes and STRING (Szklarczyk et al., 2015) interactions (CORUM ROC
163 AUC: 0.77 vs 0.56, and STRING ROC AUC: 0.75 vs 0.58; for proteomics and gene
164 expression respectively) (**Figures 1A and 1B**). The ability to also recapitulate any
165 type of STRING interaction indicates that protein co-regulation also encompasses
166 functional relationships beyond structural physical interactions.

167

168 **Weighted correlation networks reveal the interdependencies of protein** 169 **complexes and biological processes**

170 We have shown above that the co-regulation of protein abundance is a strong
171 predictor of physical and functional associations. We therefore conducted systematic
172 genome-wide analysis of the colorectal cancer cell protein-protein correlation
173 network. To this end, we performed a weighted correlation network analysis
174 (WGCNA) (Langfelder and Horvath, 2008) using 8,469 proteins quantified in at least
175 80% of the cell lines. A total of 203 modules of co-regulated proteins ranging in size
176 from 3 to 1,065 proteins (median = 9) were detected. A comprehensive description of
177 the modules was devised based on enrichment analysis (**Table S5**) and the basic
178 structure of the colorectal cancer network is depicted in **Figure 1C**. We found that
179 approximately 60% of the modules displayed overrepresentation of protein
180 complexes (**Figure 1D**) and that the largest modules were associated with RNA
181 processing, plasma membrane, cytosolic ribosome, cell cycle, mitochondrial
182 translation, mitochondrial respiratory chain, immune response and small molecule
183 metabolic process (FDR<0.01). The full WGCNA network with weights greater than
184 0.02 is provided in **Table S6**.

185 To identify regulators of protein modules not explained by physical protein
186 interactions we examined whether enriched transcription factors from ENCODE and
187 ChEA databases in a given module were indeed co-expressed at the protein level
188 along with their transcriptional targets. We found that the “small molecule metabolic
189 process” module was enriched for the transcription factors HNF4A and CDX2 with 66
190 and 22 transcriptional targets respectively (Benj. Hoch. FDR=0.00027 and

191 FDR=0.00012 respectively). HNF4A (Hepatocyte nuclear factor 4-alpha) is an
192 important regulator of metabolism, cell junctions and the differentiation of intestinal
193 epithelial cells (Garrison et al., 2006) and has been previously associated with
194 colorectal cancer proteomic subtypes in human tumours analyzed by the CPTAC
195 consortium (Zhang et al., 2014). The “plasma membrane” module included 60
196 transcriptional targets of KLF5 (Benj. Hoch. FDR=0.00174) which itself was a
197 member of this module. KLF5 is predominantly expressed in the proliferating cells of
198 the crypt and appears to play a growth regulatory role in the intestine (Bateman et
199 al., 2004). Interestingly, KLF5 was significantly correlated with HNF4A
200 (Pearson=0.79, Benj. Hoch. FDR=4.17E-10) providing a potential link between cell
201 communication and HNF4A regulated metabolic functions. Moreover, the “plasma
202 membrane” module was enriched for an epithelial-mesenchymal transition (EMT)
203 gene set and was characterized by the anti-correlation between the epithelial marker
204 CDH1 and the mesenchymal marker Vimentin (VIM). STAT1 and STAT2 are master
205 regulators of the “immune response” module with a total of 45 targets including
206 several interferon-induced proteins (Benj. Hoch. FDR=5.06E-17 and FDR=2.05E-21
207 respectively). The protein correlation modules clearly serve as unique attributes by
208 which upstream regulatory events can be identified at the protein level. This
209 approach leverages a greater number of proteomic features and extends our
210 knowledge about cancer associated regulators beyond the use of profiles for single
211 transcription factors.

212 To better understand the interdependencies of protein complexes and biological
213 processes of the colorectal cancer cells in a global way we plotted the module-to-

214 module relationships as a correlation network. The nodes denote significant terms
215 from each module and the edges represent pairwise correlations between the
216 eigengenes (first principal component, **Table S7**) of the modules (Pearson>0.46,
217 Benj. Hoch. FDR<0.01) (**Figure S2B**). The topology of this network revealed a
218 strong coordination between the cytoplasmic ribosome and the RNA processing
219 complexes that were further linked with protein processing complexes as well as with
220 cell division and DNA repair protein complexes. Interestingly, the MCM complex, in
221 contrast to the GINS complex, was better correlated with the cytoplasmic ribosome
222 rather than the cell division processes coinciding with the proposed mechanisms
223 which couple DNA and protein syntheses (Berthon et al., 2009). The 26S
224 proteasome complex was directly associated with the RNA spliceosome generating
225 the hypothesis of another un-expected interplay between RNA processing and
226 protein turnover. In fact, the mechanistic connections between transcription and the
227 Ubiquitin-Proteasome system have been previously discussed (Muratani and
228 Tansey, 2003). The mitochondrial functions formed a distinct cluster comprised of
229 mitochondrial translation and cellular respiration complexes. The HNF4A-metabolic
230 related module is linked with the mitochondrial cluster as many of the HNF4A targets
231 are involved in mitochondrial processes. Protein processing and trafficking
232 complexes were grouped together and were associated with actin related complexes
233 as well as with proteins involved in antigen presentation. The immune response
234 signature was tightly correlated with focal adhesion, caveola proteins, septin
235 complex and key proteins of the Hippo signalling pathway such as SAV1, STK3 and
236 YAP1 revealing novel associations. Highlighting the modules with high mean mRNA-

237 to-protein correlations on the network confirmed that the HNF4A, CDX2, KLF5,
238 STAT1 and STAT2 modules were strongly driven by transcription, whereas nodes
239 representing protein complexes corresponded poorly to mRNA levels (**Figure S2B**).
240 The correlation of RNA to protein levels also appears to be modestly influenced by
241 protein class (**Figure S3A**). Interestingly, proteins characterized by degradation that
242 is best explained by a two-state model with two different degradation rates (non-
243 exponentially degraded, NED) present significantly lower mRNA-to-protein
244 correlations compared to their exponentially degraded (ED) counterparts (McShane
245 et al., 2016) that are explained by one-state model (**Figure S3B**).

246 Examination of unannotated modules revealed correlation between tRNA
247 methyltransferases and eukaryotic translation elongation factors (**Figure S3C**). This
248 suggests that the maintenance of total stoichiometry is also a feature of interactions
249 between protein complexes. We also mapped 66 uncharacterized proteins to 30
250 modules that allowed them to be functionally annotated based on their associated
251 module function. Some candidates for further investigation include: cell cycle
252 (C1orf112, C14orf80, C12orf45, C9orf78, C10orf12, C1orf52), tRNA processing
253 (C18orf21), mitochondrial translation (C6orf203, C2orf47), chaperonin-containing T-
254 complex (C12orf29), N-terminal protein amino acid acetylation (C8orf59), protein N-
255 linked glycosylation via asparagine (C20orf24), Oxidative phosphorylation (C8orf82),
256 ribonuclease H2 complex (C8orf76) and BLOC-1 complex (C10orf32). Taken
257 together, our protein correlations reveal a higher order of cellular functions in a well-
258 organized structure shaped by the compartmental interactions between protein

259 complexes and clearly divided into transcriptionally and post-transcriptionally
260 regulated sectors.

261

262 ***De novo* prediction of phosphorylation networks reveals novel** 263 **functional relationships**

264 The scale of global phosphorylation survey accomplished here offers the opportunity
265 for the *de novo* prediction of kinase-substrate associations inferred by co-changing
266 phosphorylation patterns that involve kinases (Ochoa et al., 2016; Petsalaki et al.,
267 2015). Phosphorylated proteins are highly enriched for spliceosomal and cell cycle
268 functions and cover a range of cancer related pathways (**Figure S3D**). Notably, for
269 about 450 partially overlapping CORUM complexes more than 60% of their subunits
270 were found to be phosphorylated. To detect differential phosphorylation we
271 regressed protein abundances from the respective phosphorylation profiles as the
272 two levels of information are strongly correlated (**Figure S3E**). Pairwise correlation
273 analysis among 213 variable phosphopeptides belonging to 144 kinases, and the
274 787 most variable phosphopeptides from other protein types (**Table S8**) revealed a
275 strong enrichment of nucleosome assembly proteins (mainly histones) (FDR=5.09E-
276 08) correlating with three kinases, namely CSNK1A1, CDK13 and VRK3 (**Figure**
277 **S4A**). CSNK1A1 is a casein kinase that participates in Wnt signalling where it is
278 essential for β -Catenin phosphorylation and degradation (Liu et al., 2002). CSNK1A1
279 has also been implicated in the segregation of chromosomes during mitosis and may
280 be cell cycle-regulated (Brockman et al., 1992). VRK3 has recently been shown to

281 be an active kinase as well as a signalling scaffold in cells, with a specific role in
282 DNA replication and chromatin dissociation during interphase (Park et al., 2015).
283 Interestingly, VRK3 phosphorylation status was strongly correlated with the
284 phosphorylation of 7 histones of which H2AFX presented the strongest correlation in
285 spite of poor correlation of their respective protein abundance profiles (**Figure S4B,**
286 **top panel**). Strongly maintained co-phosphorylation was also observed between
287 RAF1, MAPK1, MAPK3 and RPS6KA3 of the MAPK pathway (**Figure S4C, left**
288 **panel**) as well as between CDK1 and CDK7 of the cell cycle pathway (**Figure S4C,**
289 **right panel**). The correlation plots of MAPK1 and MAPK3 phosphorylation and total
290 protein are depicted in **Figure S4B, bottom panel**. Overall, these examples
291 demonstrate that functional relationships are encrypted in the patterns of co-
292 regulated phosphorylation events.

293

294 **Protein abundance and phosphorylation variation are associated with** 295 **genomic alterations**

296 Assessing the impact of non-synonymous protein coding variants and copy number
297 alterations on protein abundance is fundamental in understanding the link between
298 cancer genotypes and the dysregulated biological processes. To investigate this, we
299 first examined whether driver mutations in any of the 17 colorectal cancer driver
300 genes (Iorio et al., 2016) with at least 5 occurrences across the cell lines could alter
301 the levels of their protein products. Strikingly, for 7 such genes (*PTEN*, *B2M*, *CD58*,
302 *PIK3R1*, *ARID1A*, *BMP2* and *MSH6*) we found that driver mutations had a
303 significant negative impact on the respective protein abundances, in line with their

304 function as tumour suppressors, whereas missense mutations in *TP53* were
305 associated with elevated protein levels (ANOVA test, permutation-based FDR<0.05)
306 (**Figure 2A**). Protein abundance variation of APC could not be systematically
307 attributed to genomic variants although it is possible that extreme protein changes in
308 individual cell lines could be the result of specific mutations. For example, the low
309 relative abundance of APC protein in 5 out of 7 cell lines with $\log_2\text{Ratio}<-0.8$ could be
310 explained by the presence of frameshift and nonsense mutations in these particular
311 cell lines. Distinctly, for the majority of driver mutations in oncogenes, there was no
312 clear relationship between the presence of mutations and protein expression. From
313 these observations we conclude, that mutations in canonical tumour suppressor
314 genes predicted to cause premature stop codons and ultimately nonsense-mediated
315 decay of transcript were significantly associated with decreased protein abundance
316 compared to driver mutations in oncogenes (Log_2 mean protein abundance -0.65 vs
317 $+0.45$ respectively) (**Figure 2B**), suggesting that these have distinct regulation or
318 gain new function upon mutation.

319 We extended our analysis to globally assess the effect of mutations on the protein
320 abundances. For 5,498 genes harbouring any type of non-synonymous protein
321 coding variants in at least three cell lines, 626 proteins exhibited lower (N=566) or
322 higher (N=60) abundances in the mutated versus the wild-type cell lines at ANOVA
323 $p\text{-value}<0.05$ (all 77 proteins with FDR < 0.1 were associated with decreased levels
324 of expression) (**Figure 3A**). This high confidence subset was enriched for
325 phosphatidylinositol signalling proteins (KEGG FDR=0.0307: PTEN, PIK3R1,
326 PIK3C2A, MTM1 and PIK3C2B) and included 5 tumour suppressors (MLH1, MSH2,

327 NF1, PIK3R1, PTEN). Restricting the analysis to frameshift mutations only (the
328 second most frequent mutation type), we found that 136 of the 389 genes presented
329 lower abundances in the mutated cell lines with ANOVA p-value<0.05 of which 121
330 passed the 10% FDR cut-off (**Figure 3B**). Notably, the significantly affected proteins
331 were strongly enriched for chromatin modification proteins (FDR=2.66E-10, N=23)
332 and included 10 oncogenes and 4 tumour suppressors. The STRING network of the
333 most significant hits is depicted in **Figure 3E**. A less pronounced impact of frameshift
334 mutations was found at the mRNA level where only 15% of the 349 genes (with both
335 mRNA and protein data) exhibited altered mRNAs abundances in the mutated
336 samples at ANOVA p-value<0.05, only 19 of these were below the 10% FDR (**Figure**
337 **3C**). The overlap between the different analyses is depicted in **Figure 3D**.
338 Considering all proteins negatively affected by mutations we found
339 overrepresentation of proteins with certain domains (e.g. helicase) as well as
340 enrichment of certain classes of enzymes such as kinases, transferases and
341 hydrolases (**Figure 3F**) highlighting the protein classes that are subject to protein
342 abundance reduction upon structural changes. Notably, 59 out of the 677 genes
343 affected by genomic variants at p-value<0.05 are currently catalogued in the
344 COSMIC Census list of genes for which mutations have been causally implicated in
345 cancer.

346 We also explored the effect of 20 recurrent copy number alterations (CNAs) using
347 binary-type data on the protein abundances of 212 falling within these intervals.
348 Amplified genes tended to display increased protein levels whereas gene losses had
349 an overall negative impact on protein abundances although with several exceptions

350 **(Figure 3G, top panel)**. The 51 significant genes with ANOVA p-value <0.05 (31
351 genes at $FDR < 0.1$) were mapped to 13 genomic loci. The 13q33.2 amplification
352 encompasses the highest number of affected proteins **(Figure 3G, bar plot)**. Losses
353 in 18q21.2, 5q21.1 and 17p12 loci are associated with reduced protein levels of
354 three important colorectal cancer drivers, SMAD4, APC and MAP2K4 respectively
355 ($FDR < 0.1$). Increased levels of CDX2 and HNF4A were modestly associated with
356 13q12.13 and 20q13.12 amplifications ($p\text{-value} < 0.1$, $FDR < 30\%$). Global correlation
357 using normalized log₂ copy number ratios obtained from the Cancer Cell Line
358 Encyclopaedia showed median CNA to mRNA and protein correlations 0.35 and
359 0.23 respectively **(Figure 3H)**. To summarize the possible levels of regulation we
360 trained a Self-Organizing Map (SOM) using the Pearson correlations coefficients
361 between CNA and mRNA, CNA and protein, mRNA and protein (three vectors) for
362 each protein, which indicated three main regulatory routes: good concordance
363 between the three levels **(Figure 3I, left side cluster)**, CNAs corresponding to
364 mRNAs but buffered at the protein level **(Figure 3I, top-middle cluster)** and
365 proteins well corresponding to mRNA irrespective of CNAs **(Figure 3I, bottom right**
366 **cluster)**. Taken together, our results show that copy number alterations more often
367 affect the mRNA levels than the protein levels, which needs to be taken under
368 consideration when gene expression data are used as a proxy of the protein levels
369 for the identification of actionable pathways. A summary of all proteins affected by
370 mutations and recurrent CNAs is in **Table S9**.

371 Next we assessed the direct impact of mutations on net protein phosphorylation. We
372 found 72 differentially phosphorylated proteins in the mutated cell lines (Welch's t-

373 test, FDR<10%) with both positive and negative effects (**Figure S5A**). The SRC
374 kinase and the RUNX1 transcription factor were among the top over-phosphorylated
375 proteins while APC was among the top hypo-phosphorylated proteins. We then
376 focused on eight colorectal cancer genes (*APC*, *TP53*, *KRAS*, *BRAF*, *PIK3CA*,
377 *PTEN*, *RNF43* and *PIK3R1*) to individually assess the extended effects of driver
378 mutations on the phosphorylation status of the colorectal cancer pathway. We found
379 that *APC* mutations were associated with decreased phosphorylation of APC and
380 increased phosphorylation of AXIN1, and that *PTEN* mutations were related to
381 increased TP53 phosphorylation at 10% FDR (ANOVA test) (**Figure S5B**). Mutations
382 in *PIK3CA* were associated with increased inactivating phosphorylation of BAD,
383 while *BRAF* V600E mutants exhibited increased AKT1 and decreased ARAF
384 phosphorylation (**Figure S5B**). As expected, for all the associations the respective
385 total protein levels were undifferentiated (**Figure S5C**). These observations indicate
386 a sophisticated level of cross talk between cancer genes.

387 Overall, we show that not all driver mutations have the same effect on protein
388 abundance. We identify key mutations that significantly impact abundance levels of
389 proteins, which converge in certain protein classes. We conclude that for only a
390 small portion of the proteome the variation in abundance can be directly explained by
391 mutations and that driver mutations also alter the phosphorylation status of colorectal
392 cancer proteins.

393

394 **The consequences of genomic alterations extend to protein complexes**

395 As tightly controlled maintenance of protein abundance appears to be pivotal for
396 protein complexes and interactions, we hypothesize that genomic variation can be
397 transferred from directly affected genes to distant gene protein products through
398 protein interactions there by explaining another layer of protein variation. We
399 retrieved strongly co-regulated interactors of the affected proteins and constructed
400 mutation-vulnerable protein networks, comprised of 1,108 total protein nodes
401 (**Figure 4A**) encompassing at least 25 protein complexes. One characteristic
402 example was the BAF complex characterized by disruption of ARID1A protein
403 abundance. Driver mutations in ARID1A were also significantly associated with
404 decreased levels of the respective module (p-value = 0.01707) (**Figure 4B**)
405 indicating a central role of ARID1A in the regulation of the profile of the complex. The
406 WGCNA networks also revealed a correlation between the BAF complex and the
407 functionally related PBAF complex containing the ARID2 and PBRM1 proteins which
408 were however mapped to different modules (**Figure 4C**). We noticed that the
409 PBRM1 sub-network displayed unusually poor overlap with STRING interactions,
410 and the correlations were attributed to the effects of co-occurring mutations on the
411 protein abundances, specifically in the hypermutated HT-115 cell line (**Figure S6A**).
412 This indicates that in addition to functional relationships, protein co-regulation can
413 also classify the effect of co-occurring genomic variants. These events are infrequent
414 and observed on small modules, which lack functional connection between their
415 components. To confirm whether the down-regulation of ARID1A, ARID2 and

416 PBRM1 can indeed affect the abundance levels of their interactors we performed
417 proteomic analysis on the respective CRISPR-Cas9 knockout (KO) clones derived
418 from human iPS cells (**Table S10**). Down-regulation of ARID1A coincided with
419 diminished levels of 8 partners in the predicted interactome that were closer to the
420 core of the network and were known components of the BAF complex whilst more
421 distant interactors were not affected (**Figure 4D**). This provides an indication that the
422 topology of the correlation network can predict the relative strengths of interactions.
423 Reduced levels of ARID2 resulted in the down-regulation of all three direct
424 interactors (BRD7, PHF10 and SCAF11) and the significant loss of PBRM1 protein.
425 Four components of the BAF complex were also weakly compromised in the ARID2
426 KO reflecting the overlap between the BAF and PBAF complexes. On the other
427 hand, loss of PBRM1 had no effect on ARID2 or any of its interactors demonstrating
428 that collateral effects transmitted through protein interactions can be directional.
429 Distinctly, loss of PBRM1 had no impact on the abundance of the constituents of its
430 module confirming that the co-variation here is due to co-occurring genomic variation
431 rather than direct interactions. Pathway enrichment analysis on the changing
432 proteins detected in the KO cell lines, revealed the differential regulation of a number
433 of biological processes reflecting the modulation of a wide range of target genes
434 (**Figure 4E**). Notably, down-regulation of ARID2 specifically activated the MAPK
435 pathway, actin cytoskeleton, ubiquitin proteolysis and immune related signalling
436 pathways that were not affected by ARID1A and PBRM1 depletion.

437 The above examples confirm that the loss of a subunit in a protein complex can
438 diminish the protein abundance of its partners but not always to the same degree

439 and with subsequent changes in the total stoichiometry. We devised linear models to
440 detect specific mutations that cause severe deviations from strongly correlating
441 protein profiles across the cell lines, thus significantly compromising the
442 maintenance of stoichiometry between the interacting partners. We detected 50 such
443 mutations (p-value<0.05); mostly frameshift or nonsense alterations but also several
444 single amino acid substitutions (**Table S11**). Two examples are provided in **Figure**
445 **S6B**. The outlier points in the correlation plots of PIK3R1-PIK3CB and SEC31A-
446 SEC13 involved in the PI3K pathway and protein processing in endoplasmic
447 reticulum respectively could be explained by a truncating mutation in PIK3R1 and a
448 missense mutation in SEC31A that significantly disrupted the total abundance
449 stoichiometry impairing their co-functionality in the associated biological processes.
450 This analysis highlights a subset of specific mutations with the highest impact on
451 protein abundance and reveals their cell line-specific consequences on protein
452 interactions. Overall, our findings indicate that an additional layer of protein variation
453 can be potentially explained by the collateral effects of mutations on tightly co-
454 regulated partners.

455

456 **Protein quantitative trait loci analysis of colorectal cancer drivers**

457 We performed Quantitative Trait Loci (QTL) analysis to systematically interrogate the
458 distant effects of colorectal cancer driver genomic alterations on protein abundance
459 (pQTL) and gene expression (eQTL). We identified 86 proteins and 196 mRNAs with
460 at least one pQTL (**Table S12**) and eQTL respectively at 10% FDR (**Figure 5A**,

461 **Figure S6C**). To assess the replication rates between independently tested QTL for
462 each phenotype pair we also performed the mapping using 6,456 commonly
463 quantified genes and we found that 64% of the pQTLs (N=74) validated as eQTLs
464 and 54% of the eQTLs (N=86) validated as pQTLs (**Figure 5B**). Ranking the pQTLs
465 (FDR<30%) by the number of associations showed that mutations in *BMP2R2*,
466 *RNF43* and *ARID1A*, as well as CNAs of regions 18q22.1, 13q12.13, 16q23.1,
467 9p21.3, 13q33.2 and 18q21.2 loci accounted for 62% of the total variant-protein pairs
468 (**Figure 5C**). The above-mentioned genomic events were also among the top 10
469 eQTL hotspots (**Figure S6D**). High frequency hotspots in chromosomes 13, 16 and
470 18 associated with CNAs have been previously identified in colorectal cancer tissues
471 (Zhang et al., 2014). Enrichment analysis of the gene sets associated with each
472 pQTLs showed overrepresentation of 12 distinct protein complexes and 36 partially
473 redundant GO terms (Fisher's test, Benj. Hoch. FDR<0.1). Interestingly, increased
474 levels of the mediator complex were associated with *FBXW7* mutations (**Figure S6E**,
475 **first panel**), an ubiquitin ligase that targets *MED13/13L* for degradation (Davis et al.,
476 2013) and *TP53* mutant cell lines were associated with up-regulation of cell division
477 related proteins (**Figure S6E, second panel**). Examination of the pQTL for other
478 functional relationships showed that driver mutations in *RNF43*, an E3 ubiquitin-
479 protein ligase that negatively regulates the Wnt signaling pathway, were positively
480 associated with *APC* protein abundance (**Figure S6E, third panel**) and that *BMP2R2*
481 mutations were negatively correlated with *TGFBR2* protein levels (**Figure S6E**,
482 **fourth panel**), both being members of the TGF-beta superfamily (Massague, 2012).
483 Our data clearly demonstrate that a large portion of genomic variation affecting

484 mRNA levels is not always transferred to the proteome. We see that distant protein
485 changes attributed to variation in cancer driver genes can be regulated directly at the
486 protein level and are not conspicuous at the mRNA level, with indication of further
487 causal effects including enzyme substrate relationships.

488

489 **Protein complexes associated with microsatellite instability**

490 Loss of DNA mismatch repair activity is responsible for the microsatellite instability
491 (MSI) observed in 15% of all colorectal cancers. MSI tumours are associated with
492 better prognosis and differential response to chemotherapy (Boland and Goel, 2010).
493 An improved understanding of the effect of MSI on cellular processes thus has the
494 potential to explain some of these clinical features. We detected 10 differentially
495 regulated modules between the MSI-high and MSI-low cell lines (Welch's t-test,
496 permutation based FDR<0.05). This encompasses 172 proteins (**Figure 6A**) that
497 include a subset of 33 genes previously attributed to MSI events in colorectal tumors
498 (Kim et al., 2013) such as MSH6, MSH3, PMS2, BAX and RAD50. The STRING
499 interactions between the MSI associated proteins are depicted in **Figure 6B**,
500 substantiating the functional relationships between this set of proteins. We
501 additionally identified epigenetic dysregulation characterized by reduced levels of
502 histone methyltransferase (KMT2D, PAXIP1, NCOA6, SETD1B, KMT2C, KMT2B
503 and KDM6A). We also detect histone deacetylation (HDAC3 and NCOR1) protein
504 complexes associated with these events. This network indicates the suppression of
505 two members of the INO80 chromatin remodelling complex (INO80D and ACTL6A)

506 in MSI-high cells and can explain the down-regulation of the Arp2/3 protein complex
507 by protein interactions with ACTL6A. Other distinct complexes we have identified
508 negatively affected by MSI were the exocyst complex, which is implicated in
509 targeting secretory vesicles to specific docking sites on the plasma membrane
510 (Heider and Munson, 2012) and the SKI complex (SKIV2L, TTC37 and WDR61),
511 which is involved in exosome-mediated RNA decay (Wang et al., 2005). Overall, this
512 alludes to multiple epigenetic mechanisms playing a role in the MSI pathogenesis,
513 and also suggests a new role for exocyst in this phenotype. The MSI up-regulated
514 modules showed over-representation of proteins from the loci 8p21 and 18q21
515 including SMAD4. Although high SMAD4 levels have been previously associated
516 with MSI and better prognosis in colon cancer (Isaksson-Mettavainio et al., 2012),
517 our data suggest that the observed differences stem from a mutually exclusive
518 SMAD4 copy number alteration (loss) with MSI.

519

520 **Proteomic subtypes of colorectal cancer cell lines**

521 To explore whether our deep proteomes recapitulate tissue level subtypes of
522 colorectal cancer and to provide insight into the cellular and molecular heterogeneity
523 of the colorectal cancer cell lines, we performed unsupervised clustering based on
524 the quantitative profiles of 7,330 proteins without missing values by class discovery
525 using the ConsensusClusterPlus method (Wilkerson and Hayes, 2010). Optimal
526 separation by k-means clustering was reached using 7 colorectal proteomic
527 subtypes (CPS) (**Figure S7A and Figure 7A**).

528 Our proteomic clusters overlapped very well with previously published tissue
529 subtypes (annotations from Medico et al., **Figure S7B**) (Medico et al., 2015),
530 especially with the classification described by De Sousa E Melo et al. (De Sousa E
531 Melo et al., 2013). Previous classifiers have commonly subdivided samples along the
532 lines of 'Epithelial' (lower crypt and crypt top), 'MSI-H' and 'Stem-like', with varying
533 descriptions (Guinney et al., 2015). In contrast, our high depth proteomic dataset not
534 only captured the commonly identified classification features, but provides increased
535 resolution to further subdivide these groups. The identification of unique proteomic
536 features pointing to key cellular functions, gives insight into the molecular basis of
537 these subtypes, and provides clarity as to the differences between them (**Figure 7A**).
538 A detailed description of the unique proteomic features of our COREAD classification
539 is provided in **Table S13**.

540 Cell lines with a canonical epithelial phenotype (previously classified as CCS1 by De
541 Sousa E Melo et al., 2013) clustered together, but are now divided into 3 subtypes
542 (CPS1, CPS2, CPS3). These subtypes all displayed high expression of HNF4A,
543 indicating a more differentiated state. While subtypes CPS1 and CPS3 contain
544 Transit Amplifying cell phenotypes (Sadanandam et al., 2013), CPS2 is largely
545 characterised by a Goblet cell signature (**Figure S7B**). CPS2 is also enriched in lines
546 that are hypermutated, and while some are MSI-H, the MSI-negative/hypermutated
547 lines (HT115, HCC2998, HT55) (Medico et al., 2015) all cluster in this group (**Figure**
548 **S7B**). Transit Amplifying subtype CPS3 can be distinguished from CPS1 by lower
549 expression of cell cycle proteins (eg. CCND1), and low histone phosphorylation
550 (possibly mediated by VRK3), as well as higher activation of PPAR signalling and

551 amino-acid metabolism pathways. CPS3 also contains lines (DIFI, NCI-H508) that
552 are most sensitive to the anti-EGFR antibody Cetuximab (Medico et al., 2015).
553 Further, this group correlates with a crypt top description 'Subtype A' (Budinska et
554 al., 2013) while subtypes CPS1 and CPS2 are associated with the lower crypt
555 'Subtype B' (Budinska et al., 2013), (**Figure S7B**).

556 The CPS4 subtype is the canonical MSI-H cluster, with a strong correlation with the
557 CCS2 cluster identified by De Sousa E Melo et al.. These lines have also been
558 commonly associated with a less differentiated state by other classifiers, and this is
559 reinforced by our dataset; subtype CPS4 has low levels of the HNF4A-CDX2
560 module, rendering this group clearly distinguishable from CPS2 (**Figure 7A**). The
561 separation into two distinct MSI-H/Hypermethylated classifications was also observed
562 by Guinney et al., (2015), and may have implications for patient therapy and
563 prognosis. Significantly, CPS4 displays low expression of ABC transporters, which
564 may contribute to the better response rates seen in MSI-H patients (Popat et al.,
565 2005).

566 The origin of CPS5 is less well defined, as it expresses intermediate levels of
567 HNF4A. However, it is characterized by moderate down-regulation of cell cycle,
568 ribosome and spliceosome modules, and displays low levels of MLH1 and AKT3
569 proteins (**Figure 7A**). Interestingly, the phosphorylation landscape of CPS5 exhibits
570 a global low phosphorylation, particularly in microtubule cytoskeleton and adherens
571 junction proteins.

572 Lastly, we capture the commonly observed colorectal ‘Stem-like’ subgroup, which is
573 represented in subtypes CPS6 and CPS7 (**Figure 7A, S7B**). Both subtypes exhibit
574 stem-like expression profiles, with very low levels of HNF4A and CDX1 transcription
575 factors (Chan et al., 2009; Garrison et al., 2006; Jones et al., 2015). Cells in both
576 subtypes commonly exhibit loss of 9p21.3 including *CDKN2A* and *CDKN2B*, while
577 this is rarely seen in other subtypes. Interestingly, while CPS6 displays activation of
578 the Hippo signalling pathway and loss of 18q21.2 (*SMAD4*), CPS7 has a
579 mesenchymal profile, with low expression of CDH1 and AXIN2, and high Vimentin.
580 The overall strong suppression of cell adhesion and tight junction components may
581 be influenced by low expression of KLF5 (Zhang et al., 2013).

582

583 **Pharmacoproteomic models are strong predictors of response to DNA** 584 **damaging agents**

585 Although a number of recent studies have investigated the power of different
586 combinations of molecular data to predict drug response in colorectal cancer cell
587 lines, these have been limited to using genomic (mutations and copy number),
588 transcriptomic and methylation datasets (Iorio et al., 2016). We have shown above
589 that the DNA and gene expression variations are not perfectly mirrored in the protein
590 measurements. As such one might expect to gain predictive power for some
591 phenotypic associations when also using the protein abundance changes. To date
592 there has not been a comprehensive analysis of the effect on the predictive power
593 from the addition of proteomics datasets in colorectal cancer. All of the colorectal cell

594 lines included in this study have been extensively characterised by sensitivity data
595 (IC50 values) for 265 compounds (lorio et al., 2016). These include clinical drugs (n
596 = 48), drugs currently in clinical development (n = 76), and experimental compounds
597 (n = 141).

598 We built Elastic Net models that use as input features genomic (mutations and copy
599 number gains/losses), methylation (CpG islands in gene promoters), gene
600 expression and proteomic datasets. We were able to generate predictive models
601 where the Pearson correlation between predicted and observed IC50 >0.4 in 91 of
602 the 265 compounds (**Table S14**). Importantly, using the proteomics data enabled the
603 construction of more predictive models than with any other feature type (**Figure 7B,**
604 **top panel**). Examples of the proteomics predictive models for etoposide are shown
605 in **Figure S7C**. Response to most drugs was often specifically predicted by one data
606 type, with very little overlap (**Figure 7B bottom panel**). Interestingly, when response
607 to a drug was predicted by both gene expression and proteomics, the protein-RNA
608 correlation for genes associated with response tended to be higher (Mann-Whitney U
609 test, p-value: 0.006) (**Figure S7D**).

610 Within the proteomics-based signatures found to be predictive for drug response, we
611 frequently observed the drug efflux transporters ABCB1 and ABCB11 (8 and 7 out of
612 43 respectively, 9 unique) (**Table S14**). In all models containing these proteins,
613 elevated expression of the drug transporter was associated with drug resistance.
614 Interestingly, ABCB1 and ABCB11 are very tightly co-regulated (Pearson's $r=0.94$,
615 $FDR= 5.99E-21$), suggesting a novel interaction. Notably, protein measurements of

616 these transporters correlated more strongly with response to these drugs than the
617 respective mRNA measurements (mean Pearson's $r=0.68$ and $r=0.43$ respectively,
618 Wilcoxon test $p\text{-value}=0.0005$). This suggests that the protein expression levels of
619 drug efflux pumps play a key role in determining drug response, and while predictive
620 genomic biomarkers may still be discovered, the importance of proteomic
621 associations with response should not be under-estimated.

622 To detect whether any specific classes of drug might be best predicted by each of
623 the 4 molecular features, we initially classified each of the 91 agents into 21 classes
624 depending on target class and biological activity (**Figure S7E**), and subsequently
625 further reduced the dimensionality of the data by then classifying into 4 groups
626 (termed 'kinase', 'DNA', 'chromatin' and 'other'). There was a significant enrichment
627 for proteomics within the predictive models for the 'DNA' group of agents, which
628 includes many chemotherapy agents and mitotic poisons (Mann-Whitney U test,
629 nominal $p\text{-value}$: 0.01, FDR-corrected $p\text{-value}$: 0.0486) (**Figure 7C**). In contrast, 47%
630 of the models specifically predicted using genomics features were for kinase
631 inhibitors. This suggests that while the response to targeted kinase inhibitors can be
632 modulated by point mutation (eg. BRAF mutations predict response to BRAF
633 inhibitors), the response to broader DNA damaging agents is dependent on the
634 expression levels of key proteomic subsets. Interestingly, NBEAL1 and PARD3B
635 proteins that were found to be down-regulated by mutations were also among the top
636 10 proteomics predictive models, were absent from the gene expression models,
637 and 6 out of the 11 associated drugs were from the "DNA" category. Although the
638 mechanism by which these proteins may affect drug sensitivity are unclear, it is

639 known that PARD3B is involved in asymmetrical cell division and cell polarization
640 processes (Williams et al., 2014) and that both genes are located in chr2q33 and are
641 associated with Amyotrophic Lateral Sclerosis 2 which suggests functional
642 similarities. These examples also highlight the value of proteomics in better
643 understanding the consequences of non-driver genomic alterations in drug sensitivity
644 through the proteome.

645 **Discussion**

646 Our analysis of colorectal cancer cells using in-depth proteomics has yielded several
647 significant insights into both fundamental molecular cell biology, and the molecular
648 heterogeneity of colorectal cancer subtypes. Beyond static measurements of protein
649 abundances, the quality of our dataset enabled the construction of a reference
650 proteomic co-variation map with topological features reflecting the dynamic interplay
651 between protein complexes and biological processes in colorectal cancer cells.
652 Notably, identification of protein complexes and network topologies in such a global
653 scale would require the analysis of hundreds of protein pull-downs and thousands
654 hours of analysis (Hein et al., 2015) thus our approach can serve as a time effective
655 screening tool for the study of protein networks. Another novel aspect that emerged
656 from our analysis is the maintenance of co-regulation at the level of net protein
657 phosphorylation. This seems to be more pronounced in signalling pathways where
658 the protein abundances are insufficient to indicate functional associations.
659 Analogous study of co-regulation between different types of protein modifications
660 could also enable the identification of modification cross-talk (Beltrao et al., 2013).

661 We show that the subunits of protein complexes tend to tightly maintain their total
662 abundance stoichiometry post-transcriptionally which forms the basis for the better
663 understanding of the higher order organization of the proteome. The primary level of
664 co-regulation between proteins allows for prediction of human gene functions and
665 the secondary assembly of the co-variome reveals the interdependencies of protein
666 complexes and biological processes and uncovers possible pathway interplays.
667 Importantly, our data can be used in combination with genetic interaction screens
668 (Costanzo et al., 2016) to explore whether gene essentiality meets the protein co-
669 regulation principles. Our catalogue of 210,000 weighted interactions can help the
670 selection of protein hubs representing the best predictors of interactomes in pull-
671 down assays. Moreover, the identification of proteins with outlier profiles from the
672 conserved profile of their known interactors, within a given complex, can point to
673 their pleiotropic roles in the associated processes.

674 The simplification of the complex proteomic landscapes enables a more direct
675 alignment of genomic features with cellular functions and delineates how genomic
676 variation is received by protein networks and how this is disseminated throughout the
677 proteome. This framework also proved very efficient to identify upstream regulatory
678 events that link transcription factors to their transcriptional targets at the protein level
679 and explained the components of the co-variome not strictly shaped by physical
680 protein interactions. To a smaller degree the module-based analysis was predictive
681 of co-occurring genomic variants exposing paradigms of simple cause-and-effect
682 proteogenomic features of the cell lines.

683 We show that mutations largely affect protein abundances directly at the protein level
684 with a higher pressure on the chromatin modification protein class. Targeted
685 depletion of key chromatin modifiers by CRISPR/cas9 followed by proteomic
686 analysis confirmed that the effects of genomic variation on distant gene products,
687 physically related with the directly affected proteins, can be explained by the
688 mechanisms that define protein co-variation. The latter is supported by the
689 observation that the severity of the distant effects is well predicted by the co-variation
690 consistency and the topological features of the correlation networks. Additionally, this
691 analysis indicated that directionality can be another characteristic of such
692 interactions.

693 We provide evidence that colorectal cancer subtypes derived from tissue level gene
694 expression datasets are largely reproduced at the proteome level which further
695 resolves the main subtypes into groups that reflect a possible cell type of origin and
696 the underlying differences in genomic alterations. This robust functional
697 characterization of the COREAD cell lines can be a useful resource to guide cell line
698 selection in targeted cellular and biochemical experimental designs where cell line
699 specific biological features can bias the results. Importantly, proteomics analysis
700 highlighted that protein variation better predicts responses to drugs that interfere with
701 cell cycle and DNA replication and that the expression of key protein components
702 such as ABC transporters is critical to predicting drug response in colorectal cancer.
703 While further work is required to establish these as validated biomarkers of patient
704 response in clinical trials, numerous studies have noted the role of these channels in
705 aiding drug efflux (Chen et al., 2016). Overall, our Elastic Net models suggest that

706 expression of a single protein alone may not be sufficient to predict drug resistance,
707 and consideration of a panel of markers may be required. This study demonstrates
708 that proteomics is the technology of choice for functional systems biology and
709 provides a valuable resource for the study of regulatory variation in cancer cells.

710 **Author Contributions**

711 Conceptualization, J.S.C. & U.M.; Methodology, T.I.R., S.W.; Mass Spectrometry
712 T.I.R.; Proteomics Data Analysis, T.I.R., E.G., M.S., S.W., J.C.W., P.B., J.S-R.; QTL
713 Analysis, F.Z.G., E.G., A.B., J.S-R.; O.S.; Drug Data Analysis, N.A., M.M., M.S.,
714 M.Y., J.S-R.; S.W., T.I.R., L.W., U.M.; CRISPR Lines, C.A., D.J.A.; Writing – Original
715 Draft, T.I.R., S.W., L.W., U.M. J.S.C.; Writing - Review & Editing, All.

716 **Acknowledgments**

717 This work was funded by the Wellcome Trust (086375 and 102696). SPW is funded
718 by the ERC Synergy Project CombatCancer. We would like to thank members of the
719 Cancer Genome Project for helpful discussions and Sarah A. Teichmann for
720 discussions and general suggestions about the manuscript.

721

722

723

724

725

726

727

728

729 **References**

730

- 731 Bateman, N.W., Tan, D.F., Pestell, R.G., Black, J.D., and Black, A.R. (2004). Intestinal tumor
732 progression is associated with altered function of KLF5. *Journal of Biological Chemistry* 279,
733 12093-12101.
- 734 Beltrao, P., Bork, P., Krogan, N.J., and van Noort, V. (2013). Evolution and functional cross-
735 talk of protein post-translational modifications. *Molecular systems biology* 9, 714.
- 736 Berthon, J., Fujikane, R., and Forterre, P. (2009). When DNA replication and protein
737 synthesis come together. *Trends in biochemical sciences* 34, 429-434.
- 738 Boland, C.R., and Goel, A. (2010). Microsatellite Instability in Colorectal Cancer.
739 *Gastroenterology* 138, 2073-U2087.
- 740 Brockman, J.L., Gross, S.D., Sussman, M.R., and Anderson, R.A. (1992). Cell cycle-
741 dependent localization of casein kinase I to mitotic spindles. *Proceedings of the National*
742 *Academy of Sciences of the United States of America* 89, 9454-9458.
- 743 Budinska, E., Popovici, V., Tejpar, S., D'Ario, G., Lapique, N., Sikora, K.O., Di Narzo, A.F.,
744 Yan, P., Hodgson, J.G., Weinrich, S., *et al.* (2013). Gene expression patterns unveil a new
745 level of molecular heterogeneity in colorectal cancer. *The Journal of pathology* 231, 63-76.
- 746 Chan, C.W.M., Wong, N.A., Liu, Y., Bicknell, D., Turley, H., Hollins, L., Miller, C.J., Wilding,
747 J.L., and Bodmer, W.F. (2009). Gastrointestinal differentiation marker Cytokeratin 20 is
748 regulated by homeobox gene CDX1. *Proceedings of the National Academy of Sciences of*
749 *the United States of America* 106, 1936-1941.
- 750 Chen, Z., Shi, T., Zhang, L., Zhu, P., Deng, M., Huang, C., Hu, T., Jiang, L., and Li, J.
751 (2016). Mammalian drug efflux transporters of the ATP binding cassette (ABC) family in
752 multidrug resistance: A review of the past decade. *Cancer letters* 370, 153-164.
- 753 Costanzo, M., VanderSluis, B., Koch, E.N., Baryshnikova, A., Pons, C., Tan, G., Wang, W.,
754 Usaj, M., Hanchard, J., Lee, S.D., *et al.* (2016). A global genetic interaction network maps a
755 wiring diagram of cellular function. *Science* 353.
- 756 De Sousa E Melo, F., Wang, X., Jansen, M., Fessler, E., Trinh, A., de Rooij, L.P.M.H., de
757 Jong, J.H., de Boer, O.J., van Leersum, R., Bijlsma, M.F., *et al.* (2013). Poor-prognosis
758 colon cancer is defined by a molecularly distinct subtype and develops from serrated
759 precursor lesions. *Nature medicine* 19, 614-618.
- 760 Garrison, W.D., Battle, M.A., Yang, C.H., Kaestner, K.H., Sladek, F.M., and Duncan, S.A.
761 (2006). Hepatocyte nuclear factor 4 alpha is essential for embryonic development of the
762 mouse colon. *Gastroenterology* 130, 1207-1220.
- 763 Guinney, J., Dienstmann, R., Wang, X., de Reynies, A., Schlicker, A., Soneson, C., Marisa,
764 L., Roepman, P., Nyamundanda, G., Angelino, P., *et al.* (2015). The consensus molecular
765 subtypes of colorectal cancer. *Nature medicine* 21, 1350-1356.
- 766 Heider, M.R., and Munson, M. (2012). Exorcising the exocyst complex. *Traffic* 13, 898-907.

767 Hein, M.Y., Hubner, N.C., Poser, I., Cox, J., Nagaraj, N., Toyoda, Y., Gak, I.A., Weisswange,
768 I., Mansfeld, J., Buchholz, F., *et al.* (2015). A human interactome in three quantitative
769 dimensions organized by stoichiometries and abundances. *Cell* 163, 712-723.
770 Iorio, F., Knijnenburg, T.A., Vis, D.J., Bignell, G.R., Menden, M.P., Schubert, M., Aben, N.,
771 Goncalves, E., Barthorpe, S., Lightfoot, H., *et al.* (2016). A Landscape of Pharmacogenomic
772 Interactions in Cancer. *Cell* 166, 740-754.
773 Isaksson-Mettavainio, M., Palmqvist, R., Dahlin, A.M., Van Guelpen, B., Rutegard, J., Oberg,
774 A., and Henriksson, M.L. (2012). High SMAD4 levels appear in microsatellite instability and
775 hypermethylated colon cancers, and indicate a better prognosis. *International journal of*
776 *cancer Journal international du cancer* 131, 779-788.
777 Jones, M.F., Hara, T., Francis, P., Li, X.L., Bilke, S., Zhu, Y.L., Pineda, M., Subramanian, M.,
778 Bodmer, W.F., and Lal, A. (2015). The CDX1-microRNA-215 axis regulates colorectal
779 cancer stem cell differentiation. *Proceedings of the National Academy of Sciences of the*
780 *United States of America* 112, E1550-E1558.
781 Kim, T.M., Laird, P.W., and Park, P.J. (2013). The Landscape of Microsatellite Instability in
782 Colorectal and Endometrial Cancer Genomes. *Cell* 155, 858-868.
783 Langfelder, P., and Horvath, S. (2008). WGCNA: an R package for weighted correlation
784 network analysis. *BMC bioinformatics* 9, 559.
785 Liu, C., Li, Y., Semenov, M., Han, C., Baeg, G.H., Tan, Y., Zhang, Z., Lin, X., and He, X.
786 (2002). Control of beta-catenin phosphorylation/degradation by a dual-kinase mechanism.
787 *Cell* 108, 837-847.
788 Massague, J. (2012). TGF beta signalling in context. *Nat Rev Mol Cell Bio* 13, 616-630.
789 McAlister, G.C., Nusinow, D.P., Jedrychowski, M.P., Wuhr, M., Huttlin, E.L., Erickson, B.K.,
790 Rad, R., Haas, W., and Gygi, S.P. (2014). MultiNotch MS3 Enables Accurate, Sensitive, and
791 Multiplexed Detection of Differential Expression across Cancer Cell Line Proteomes. *Anal*
792 *Chem* 86, 7150-7158.
793 McShane, E., Sin, C., Zauber, H., Wells, J.N., Donnelly, N., Wang, X., Hou, J., Chen, W.,
794 Storchova, Z., Marsh, J.A., *et al.* (2016). Kinetic Analysis of Protein Stability Reveals Age-
795 Dependent Degradation. *Cell*.
796 Medico, E., Russo, M., Picco, G., Cancelliere, C., Valtorta, E., Corti, G., Buscarino, M.,
797 Isella, C., Lamba, S., Martinoglio, B., *et al.* (2015). The molecular landscape of colorectal
798 cancer cell lines unveils clinically actionable kinase targets. *Nature communications* 6, 7002.
799 Mertins, P., Mani, D.R., Ruggles, K.V., Gillette, M.A., Clauser, K.R., Wang, P., Wang, X.,
800 Qiao, J.W., Cao, S., Petralia, F., *et al.* (2016). Proteogenomics connects somatic mutations
801 to signalling in breast cancer. *Nature* 534, 55-62.
802 Muratani, M., and Tansey, W.R. (2003). How the ubiquitin-proteasome system controls
803 transcription. *Nat Rev Mol Cell Bio* 4, 192-201.
804 Ochoa, D., Jonikas, M., Lawrence, R.T., El Debs, B., Selkrig, J., Typas, A., Villen, J., Santos,
805 S.D., and Beltrao, P. (2016). An atlas of human kinase regulation. *Molecular systems biology*
806 12, 888.
807 Park, C.H., Ryu, H.G., Kim, S.H., Lee, D., Song, H., and Kim, K.T. (2015). Presumed
808 pseudokinase VRK3 functions as a BAF kinase. *Bba-Mol Cell Res* 1853, 1738-1748.
809 Petsalaki, E., Helbig, A.O., Gopal, A., Pasculescu, A., Roth, F.P., and Pawson, T. (2015).
810 SELPHI: correlation-based identification of kinase-associated networks from global phospho-
811 proteomics data sets. *Nucleic acids research* 43, W276-282.
812 Popat, S., Hubner, R., and Houlston, R.S. (2005). Systematic review of microsatellite
813 instability and colorectal cancer prognosis. *J Clin Oncol* 23, 609-618.
814 Ruepp, A., Waegele, B., Lechner, M., Brauner, B., Dunger-Kaltenbach, I., Fobo, G.,
815 Frishman, G., Montrone, C., and Mewes, H.W. (2010). CORUM: the comprehensive
816 resource of mammalian protein complexes--2009. *Nucleic acids research* 38, D497-501.

817 Sadanandam, A., Lyssiotis, C.A., Homicsko, K., Collisson, E.A., Gibb, W.J., Wullschleger,
818 S., Ostos, L.C.G., Lannon, W.A., Grotzinger, C., Del Rio, M., *et al.* (2013). A colorectal
819 cancer classification system that associates cellular phenotype and responses to therapy.
820 *Nature medicine* *19*, 619-625.

821 Stefely, J.A., Kwiecien, N.W., Freiburger, E.C., Richards, A.L., Jochem, A., Rush, M.J.,
822 Ulbrich, A., Robinson, K.P., Hutchins, P.D., Veling, M.T., *et al.* (2016). Mitochondrial protein
823 functions elucidated by multi-omic mass spectrometry profiling. *Nature biotechnology* *34*,
824 1191-1197.

825 Szklarczyk, D., Franceschini, A., Wyder, S., Forslund, K., Heller, D., Huerta-Cepas, J.,
826 Simonovic, M., Roth, A., Santos, A., Tsafou, K.P., *et al.* (2015). STRING v10: protein-protein
827 interaction networks, integrated over the tree of life. *Nucleic acids research* *43*, D447-452.

828 Wang, J., Ma, Z., Carr, S.A., Mertins, P., Zhang, H., Zhang, Z., Chan, D.W., Ellis, M.J.,
829 Townsend, R.R., Smith, R.D., *et al.* (2016). Proteome profiling outperforms transcriptome
830 profiling for co-expression based gene function prediction. *Mol Cell Proteomics*.

831 Wang, L.N., Lewis, M.S., and Johnson, A.W. (2005). Domain interactions within the Ski2/3/8
832 complex and between the Ski complex and Ski7p. *Rna* a Publication of the Rna Society *11*,
833 1291-1302.

834 Wilkerson, M.D., and Hayes, D.N. (2010). ConsensusClusterPlus: a class discovery tool with
835 confidence assessments and item tracking. *Bioinformatics* *26*, 1572-1573.

836 Williams, S.E., Ratliff, L.A., Postiglione, M.P., Knoblich, J.A., and Fuchs, E. (2014). Par3-
837 mInsc and G alpha(i3) cooperate to promote oriented epidermal cell divisions through LGN.
838 *Nat Cell Biol* *16*, 758-U231.

839 Zhang, B., Wang, J., Wang, X., Zhu, J., Liu, Q., Shi, Z., Chambers, M.C., Zimmerman, L.J.,
840 Shaddox, K.F., Kim, S., *et al.* (2014). Proteogenomic characterization of human colon and
841 rectal cancer. *Nature* *513*, 382-387.

842 Zhang, B., Zhang, Z., Xia, S., Xing, C., Ci, X., Li, X., Zhao, R., Tian, S., Ma, G., Zhu, Z., *et al.*
843 (2013). KLF5 activates microRNA 200 transcription to maintain epithelial characteristics and
844 prevent induced epithelial-mesenchymal transition in epithelial cells. *Mol Cell Biol* *33*, 4919-
845 4935.

846 Zhang, H., Liu, T., Zhang, Z., Payne, S.H., Zhang, B., McDermott, J.E., Zhou, J.Y., Petyuk,
847 V.A., Chen, L., Ray, D., *et al.* (2016). Integrated Proteogenomic Characterization of Human
848 High-Grade Serous Ovarian Cancer. *Cell* *166*, 755-765.

849

850

851

852

853

854

855 **Figure Legends**

856 **Figure 1**

857 Protein co-variation networks in colorectal cancer cell lines. A) Receiver Operating
858 Characteristic (ROC) curve illustrating the performance of proteomics and
859 microarrays based correlations to predict known interactions from the CORUM
860 protein complexes database and from B) the STRING database. C) The basic
861 structure of the full WGCNA network. Protein modules are color-coded and
862 representative enriched terms are used for the annotation of the network. D) Protein
863 abundance correlation networks derived from WGCNA analysis for enriched
864 CORUM complexes.

865

866 **Figure 2**

867 The effect of colorectal cancer driver mutations on protein abundances. A)
868 Association of driver mutations in colorectal cancer genes with the respective protein
869 abundance levels (ANOVA test, permutation based FDR<5%). The cell lines are
870 ranked by protein abundance and the bar on the top indicates the presence of
871 mutations with brown mark. B) Volcano plot summarizing the effect of loss of
872 function (LoF) and missense driver mutations on the respective protein abundances.

873

874 **Figure 3**

875 The global effects of genomic alterations on protein and mRNA abundances. A)
876 Volcano plot summarizing the effect of all protein coding non-synonymous variants
877 on the respective protein abundances (ANOVA test). B) Volcano plot summarizing
878 the effect of all protein coding frameshift mutations on the respective protein
879 abundances (ANOVA test). C) Volcano plot summarizing the effect of all protein
880 coding frameshift mutations with both mRNA and protein measurements on the
881 respective mRNA abundances (ANOVA test). D) Venn diagrams displaying the

882 overlap between proteins affected by all types of mutations and proteins affected by
883 frameshift only mutations at different confidence levels (top panel) and the overlap
884 between proteins and mRNAs affected by frameshift mutations at $p\text{-value} < 0.05$
885 (bottom panel). E) STRING network of the proteins down-regulated by frameshift
886 mutations (permutation based $FDR < 0.1$). F) Overrepresentation of Pfam domains
887 (left panel) and PANTHER protein classes (right panel) for proteins negatively
888 affected by mutations. G) Volcano plot summarizing the effect of recurrent copy
889 number alterations on the protein abundances of the contained genes (binary data,
890 ANOVA test). Red and blue points highlight genes with amplifications and losses
891 respectively. Enlarged points highlight genes at $FDR < 10\%$. The bar plot (bottom
892 panel) illustrates the number of affected proteins per genomic locus. Red and blue
893 bars indicate amplifications and losses respectively. H) Distributions of the Pearson's
894 correlation coefficients for CNA to mRNA (red) and CNA to protein (green)
895 correlations considering all genes across 38 cell lines with normalized \log_2 copy
896 number values (source: Cancer Cell Line Encyclopedia). I) Self-Organizing Map
897 trained on the Pearson correlation between CNA, mRNA and protein levels per gene
898 across the cell lines. The fan plot within each neuron displays the magnitude of each
899 one of the three vectors. Three main regulatory routes can be distinguished: good
900 concordance between the three levels (left side cluster), CNAs corresponding to
901 mRNAs but buffered at the protein level (top-middle cluster) and proteins well
902 corresponding to mRNA irrespective of the presence of CNAs (bottom right cluster).

903

904

905

906

907

908

909 **Figure 4**

910 The consequences of mutations on protein complexes. A) Correlations networks
911 filtered for known STRING interactions of proteins affected by mutations at p-
912 value<0.05. The font size is proportional to the $-\log_{10}(\text{p-value})$ and the font colour
913 displays the effect of the mutations on protein abundances (blue=negative,
914 red=positive). CORUM interactions are highlighted as green thick edges and
915 representative protein complexes are labelled. B) Boxplots illustrating the association
916 of ARID1A driver mutations with lower levels of the ARID1A complex. C) Protein
917 abundance correlation network of the ARID1A, ARID2 and PBRM1 modules. Green
918 edges denote known STRING interactions and the edge thickness is increasing
919 proportionally to the interaction weight. The node colour displays the mRNA-to-
920 protein Pearson correlation and the size of the nodes shows the protein variation. D)
921 Heatmap summarizing the protein abundance \log_2 fold-change values in the
922 knockout clones compared to the WT clones for the proteins in the ARID1A, ARID2
923 and PBRM1 modules. E) KEGG pathway and CORUM enrichment analysis for the
924 proteomic analysis results of ARID1A, ARID2 and PBRM1 CRISPR knockouts in
925 human iPS cells.

926

927 **Figure 5**

928 Proteome-wide quantitative trait loci (pQTL) analysis of cancer driver genomic
929 alterations. A) Identification of cis and trans pQTLs in colorectal cancer cell lines
930 considering colorectal cancer driver variants. The p-value and genomic coordinates
931 for the most confident non-redundant protein-variant association tests are depicted in
932 the Manhattan plot. B) Replication rates between independently tested QTL for each
933 phenotype pair using common sets of genes and variants (N=6,456 genes). C)
934 Representation of pQTLs as 2D plot of variants (y-axis) and associated genes (x-
935 axis). Associations with $q < 0.3$ are shown as dots coloured by the beta value (red:
936 positive association, blue: negative association) while the size is increasing with the

937 confidence of the association. Cumulative plot of the number of associations per
938 variant is shown below the 2D matrix.

939

940 **Figure 6**

941 Protein complexes associated with MSI. A) Heatmap of MSI-high associated proteins
942 (Welch t-test, permutation based FDR<0.05). Columns represent proteins sorted
943 horizontally based on color-coded modules and rows correspond to cell lines. The
944 modules are labelled by significantly enriched terms on the right panel. B) STRING
945 network of interconnected MSI-high associated proteins. The nodes are color-coded
946 by module and distinct complexes are highlighted. Nodes with black outline have
947 been previously found with MSI events in colorectal cancer by Kim et al.

948

949 **Figure 7**

950 Proteomics subtypes of colorectal cancer cell lines and drug associations. A) Cell
951 lines are represented as columns, horizontally ordered by seven color-coded
952 proteomics consensus clusters and aligned with microsatellite instability (MSI),
953 published colorectal cancer subtypes by DeMelo classification, HNF4A protein
954 abundance, cancer driver genomic alterations, differentially regulated proteins,
955 selected enriched KEGG pathways, differentially regulated colorectal cancer proteins
956 and differentially regulated phosphopeptides. The heatmap of the differentially
957 regulated proteome was divided into 50 color-coded clusters. Enriched terms for
958 each cluster are shown on the left. B) The number of drugs for which predictive
959 models (i.e. models where the Pearson correlation between predicted and observed
960 IC50s exceeds $r > 0.4$) could be fitted is stratified per data type (top panel).
961 Predictive models for more drugs are found by the use of proteomics data. A
962 heatmap indicating for each drug and each data type whether a predictive model
963 could be fitted (bottom panel). Drugs for which no predictive model could be fitted
964 using any data type were omitted. Most drugs were specifically predicted by one
965 data type. Examples of predictive models for each of the four data types are

966 highlighted. C) The number of drugs where response was specifically predicted by
967 one molecular data type stratified by each of the four molecular data types and by
968 four classes defined based on drug target and biological activity.

969

970

971 **Supplementary figures legends**

972

973 **Figure S1**

974 Proteome and phosphoproteome coverage. A) Workflow for quantitative global
975 proteome and phosphoproteome analysis. 50 colorectal cancer cell lines (COREAD)
976 were analysed using TMT-10plex in seven multiplex sets. The SW48 cell line was
977 used as the reference sample in each set. Biological replicates of MDST8 cell line
978 were included in two different sets and the 7th set corresponds to a biological
979 replicate of the 6th set. These were used to evaluate the normalization and the batch
980 effect correction methods. B) Number of protein groups (left panel) and unique
981 phosphopeptides (right panel) identified per multiplex set are depicted as blue bars
982 and cumulative number of identifications shown as red lines. C) Box plots of
983 normalized \log_2 Ratio values per cell line. D) Heatmap of the TMT scaled S/N values
984 of the identified mutant peptides shown as rows. The columns represent the
985 COREAD cell lines sorted from the least mutated (left) to the most mutated (right)
986 cell line. E) Two example identification spectra of mutant peptides from KRAS and
987 CTNNB1 along with the TMT quantification profiles. F) Box plots of Pearson
988 correlation coefficients between un-related samples (All), paired samples from the
989 7th replicate batch, the MDST8 replicate samples and inter-laboratory replicates for
990 six cell lines from McAlister et al. G) Box plots of the mRNA to protein Pearson
991 correlation between unrelated and paired cell lines. H) Scatter plot of protein
992 variation versus mRNA variation expressed as median-normalized standard
993 deviation (SD) across the cell lines.

994

995 **Figure S2**

996 Global distributions of gene-to-gene correlations. A) Distributions of Spearman's
997 correlation coefficients between protein-protein pairs (left panel) and mRNA-mRNA
998 pairs (right panel) for all pairs and for pairs with known relationships in the CORUM
999 database. B) Correlation network of the WGCNA modules using the eigengene
1000 profiles (Pearson>0.46, Benj. Hoch. FDR<0.01). Nodes represent WGCNA modules
1001 labelled with enriched terms (GO-Slim, KEGG, CORUM, GSEA, ChEA, ENCODE
1002 and Pfam) and are color-coded by ReactomeFI clusters. The size of the nodes is
1003 proportional to the number of proteins in the module. Transcriptionally controlled
1004 processes are highlighted with orange font and the processes with enriched
1005 transcription factors are outlined. Black thick edges highlight examples of
1006 associations between biological processes or protein complexes.

1007

1008 **Figure S3**

1009 Transcriptome-to-proteome correlation per protein class, protein modules correlation
1010 plot and phosphorylated pathways. A) Gene-level mRNA-to-protein Pearson
1011 correlations ranked by lowest to highest value. PANTHER protein classes with
1012 negative or positive enrichment relatively to the mean of all mRNA-to-protein
1013 correlations are displayed. B) Box plots illustrating the mRNA-to-protein correlation
1014 for proteins characterized as exponentially degraded (ED) and non-exponentially
1015 degraded (NED) by McShane et. al. 2016. C) Scatter plot illustrating the correlation
1016 between tRNA methyltransferases and eukaryotic translation elongation factors. D)
1017 Enriched KEGG pathways by DAVID analysis of all quantified phosphoproteins. E)
1018 The distributions of Pearson coefficients for randomized and matched pairs of
1019 phosphopeptide abundances versus protein abundances.

1020

1021

1022

1023 **Figure S4**

1024 *De novo* prediction of phosphorylation networks. A) Correlation network of 213
1025 variable phosphopeptides (with protein abundance regressed out) belonging to 144
1026 non-redundant kinases and the 787 most variable phosphopeptides of all other types
1027 of proteins (significant Pearson correlations displayed in the network were filtered for
1028 Benjamini-Hochberg adjusted p -value <0.05). The nodes are enlarged proportionally
1029 to the number of direct edges and are color-coded based on ReactomeFI clustering.
1030 Protein kinases are highlighted with bold font. B) Scatter plots of two significantly
1031 correlating phosphoprotein pairs (VRK3-H2AFX and MAPK1-MAPK3) for which the
1032 respective protein levels displayed insignificant correlation. C) Snapshots of the
1033 MAPK and cell cycle KEGG pathways highlighting (pink) significantly correlating
1034 phosphorylations. The Pearson correlation for each association is shown below the
1035 pathways.

1036

1037 **Figure S5**

1038 The effects of mutations on protein phosphorylation. A) Volcano plot summarizing
1039 the direct effects of mutations on protein phosphorylation. B) Box plots illustrating the
1040 differential phosphorylation between mutated and wild-type cells considering
1041 colorectal cancer driver mutations. C) The respective un-differentiated protein
1042 abundances.

1043

1044

1045

1046

1047

1048 **Figure S6**

1049 Identification of mutations that cause loss of correlation stoichiometry and expression
1050 quantitative trait loci analysis of cancer driver genomic alterations. A) Line plot
1051 displaying the abundance profile of the proteins in the PBRM1 module. The top bar
1052 indicates the total number of mutated genes within the module across the cell lines.
1053 B) Scatter plots of protein pairs in which specific mutation cause severe divergence
1054 from highly correlating profiles. C) Identification of *cis* and *trans* eQTLs in colorectal
1055 cancer cell lines considering cancer driver variants. The p-value and genomic
1056 coordinates for the most confident non-redundant mRNA-variant association tests
1057 are depicted in the Manhattan plot. D) Representation of eQTLs as 2D plot of
1058 variants (y-axis) and associated genes (x-axis). Associations with $q < 0.3$ are shown
1059 as dots coloured by the beta value (red: positive association, blue: negative
1060 association) while the size is increasing with the confidence of the association.
1061 Cumulative of the number of associations per variant is plotted below the 2D matrix.
1062 E) Selected examples of protein networks and individual proteins with pQTLs
1063 functionally associated with the cancer variants.

1064

1065 **Figure S7**

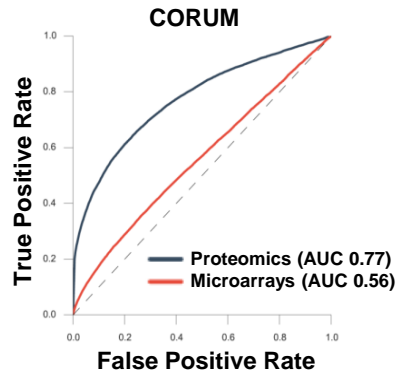
1066 Consensus clustering of colorectal cancer cell lines and drug response models for
1067 each drug target class and for each data type. A) Proteome clusters were derived
1068 based on consensus clustering. Optimal classification of the cell lines based on the
1069 proteome quantified across all cell lines ($N=7,330$) was derived by the application of
1070 the ConsensusClusterPlus R package using 1,000 resampling repetitions in the
1071 range of 2 to 10 clusters. The consensus matrices for target values $k=5,6$ and 7 are
1072 visualized (top left panel) along with the empirical cumulative distribution function
1073 (CDF) plot which indicates the k at which the distribution reaches an approximate
1074 maximum (top right panel). Cluster-consensus plot displaying the mean of all
1075 pairwise consensus values between a cluster's members at each k . Balanced mean
1076 consensus values are obtained at $k=7$. B) Overlap of the proteomics subtypes with

1077 tissue level classifications. C) Heatmap showing the proteomic signature associated
1078 with response to Etoposide. D) Gene-level mRNA-to-protein Pearson correlation for
1079 genes associated with drug response, for drugs that could be predicted by both gene
1080 expression and proteomics data (overlapping) or for drugs that could only be
1081 predicted gene expression or proteomics (non-overlapping). E) The number of drugs
1082 where response was specifically predicted by one molecular data type, stratified by
1083 each of the four molecular data types and by the 21 drug classes defined by Iorio et
1084 al. (2016).

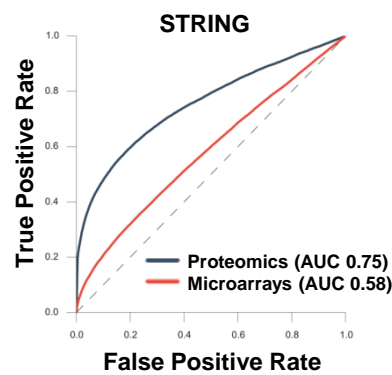
1085

Figure 1

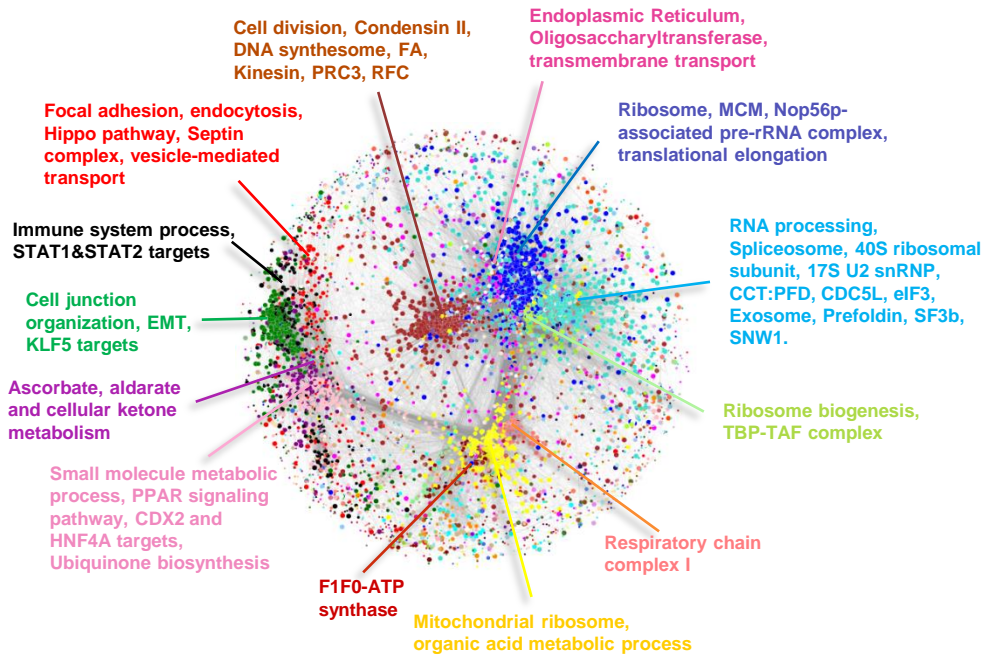
A



B



C



D

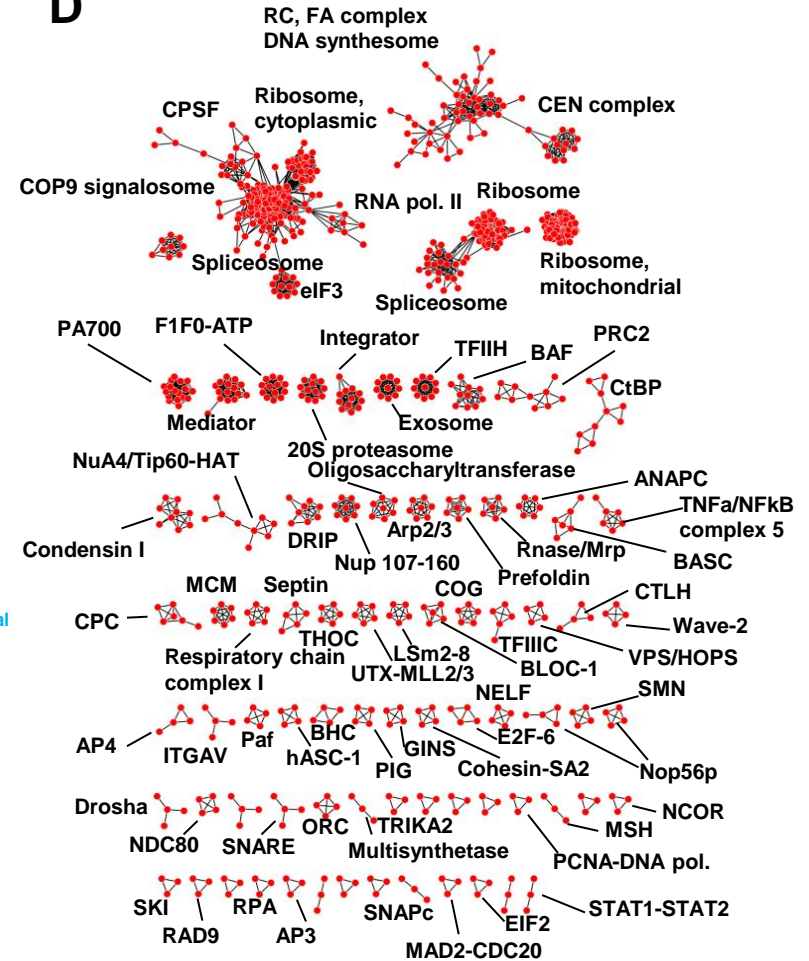
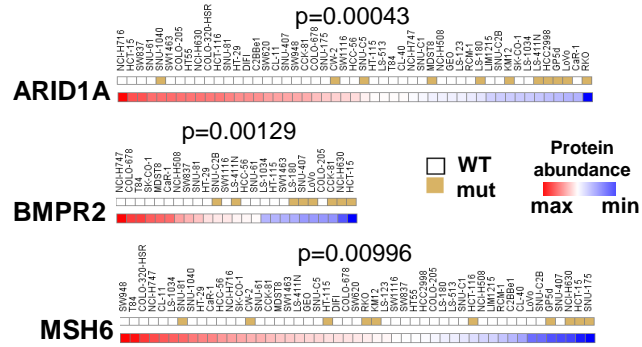
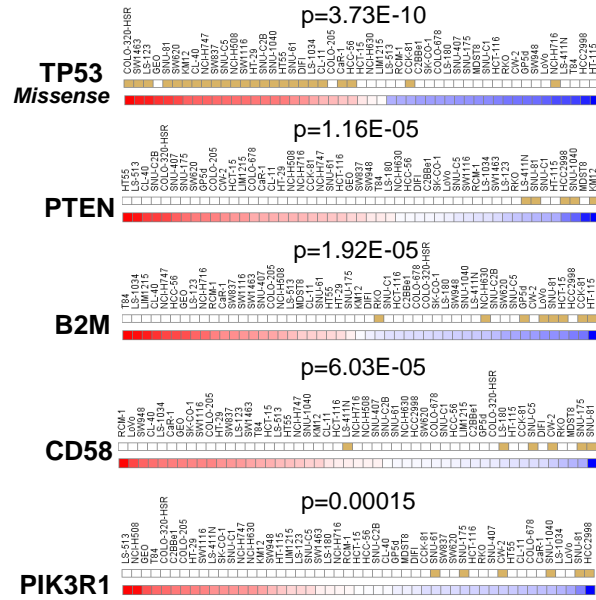


Figure 2

A



B

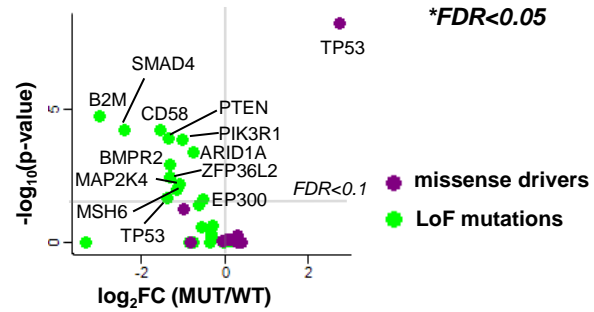


Figure 3

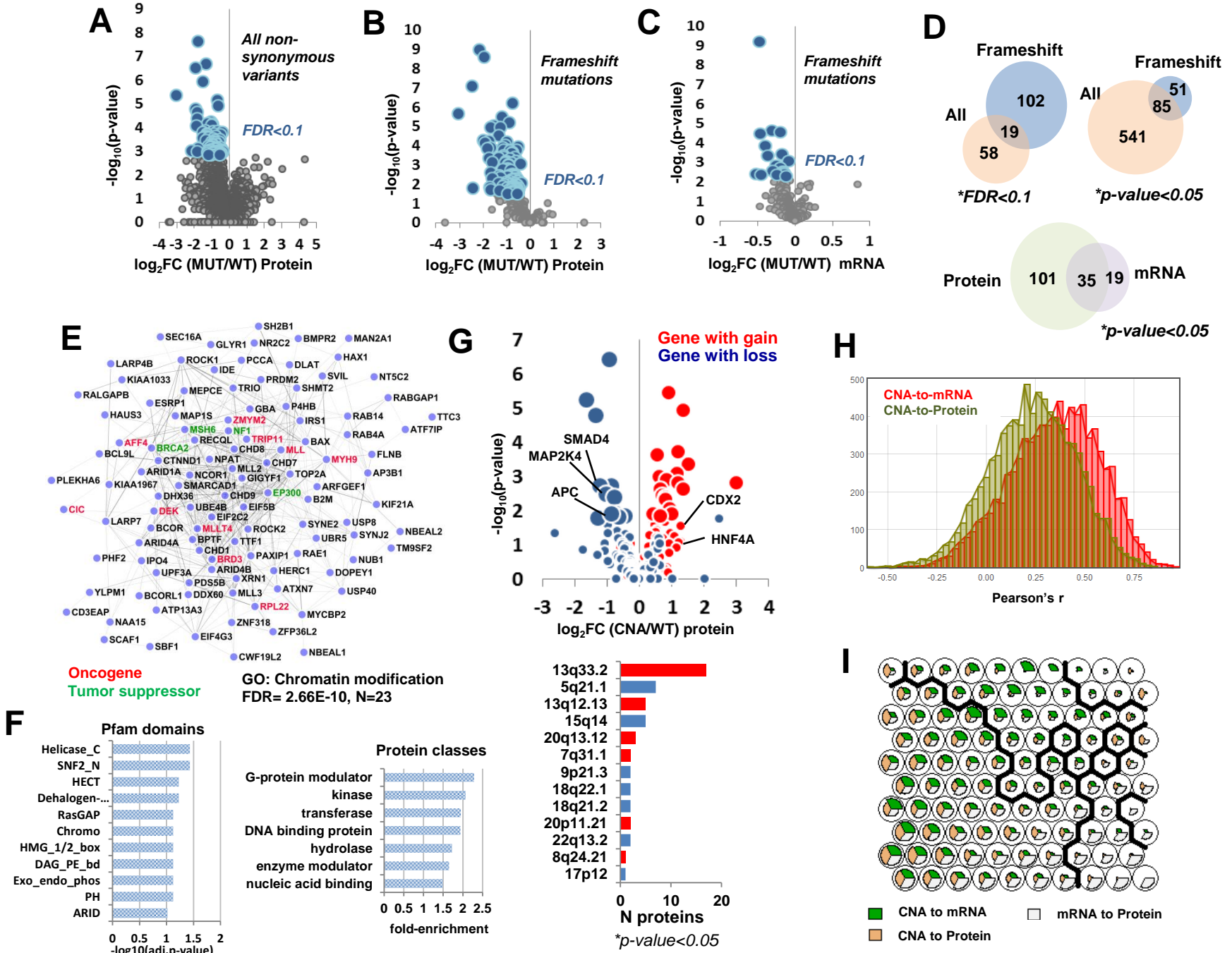


Figure 4

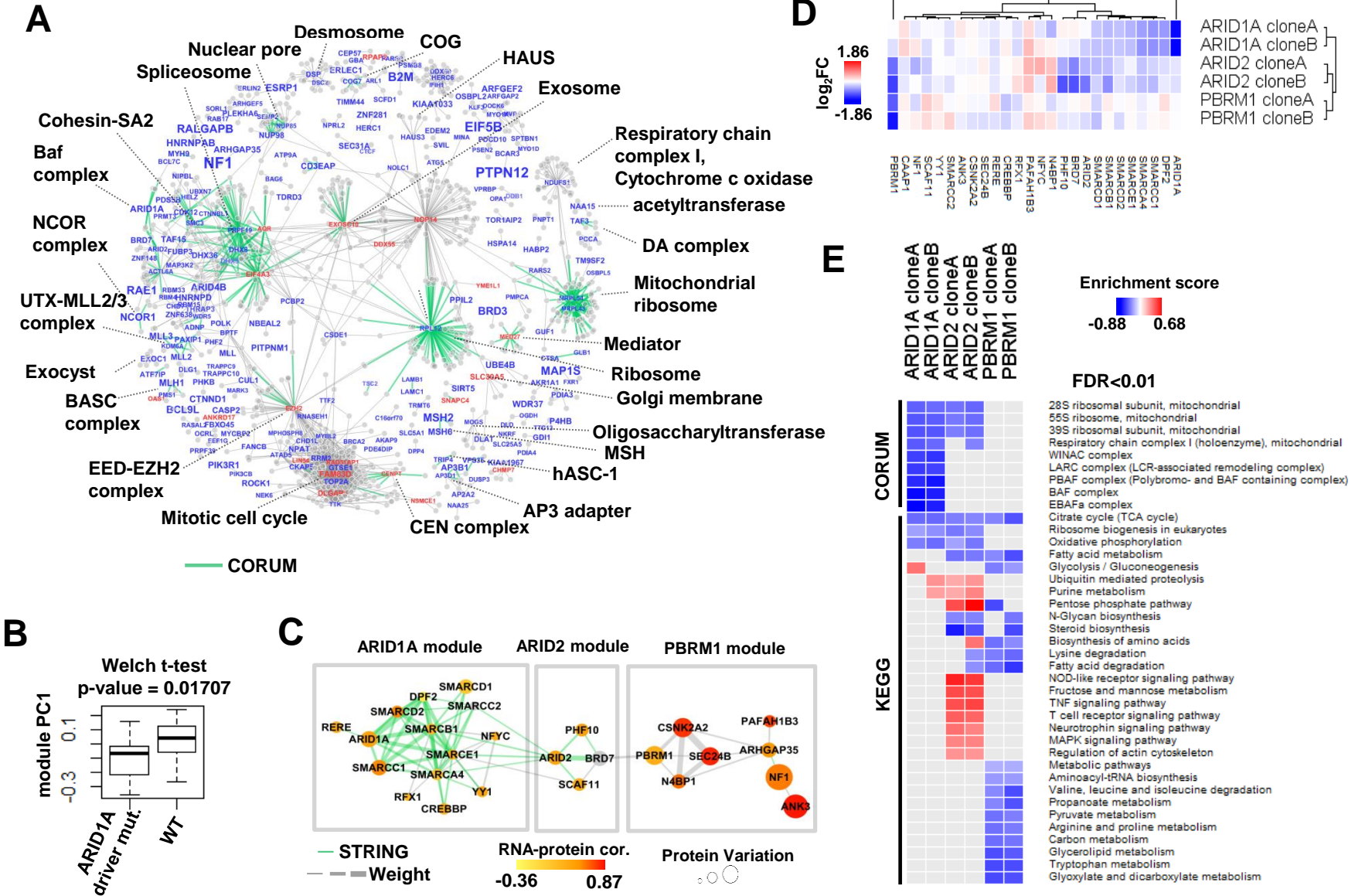


Figure 5

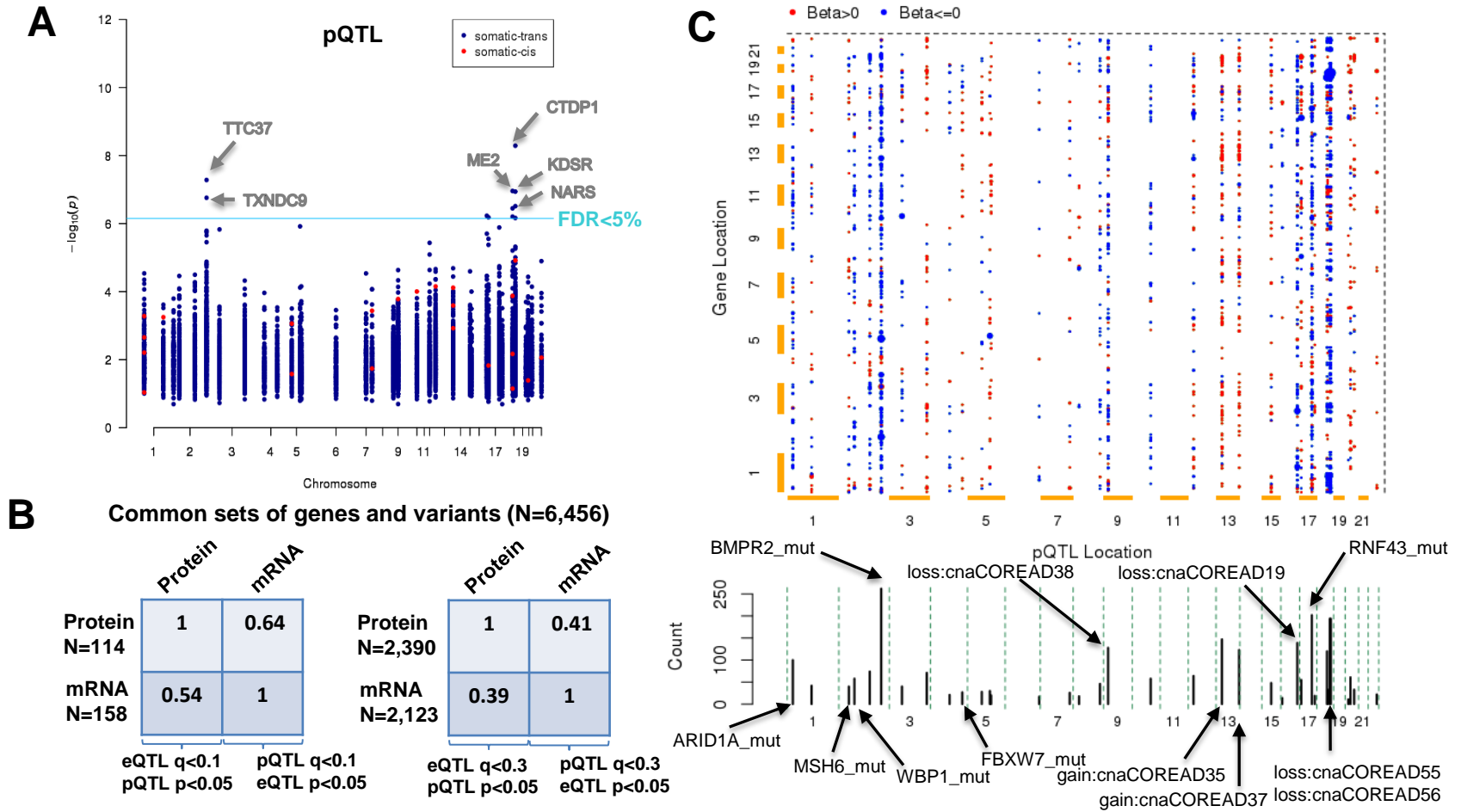
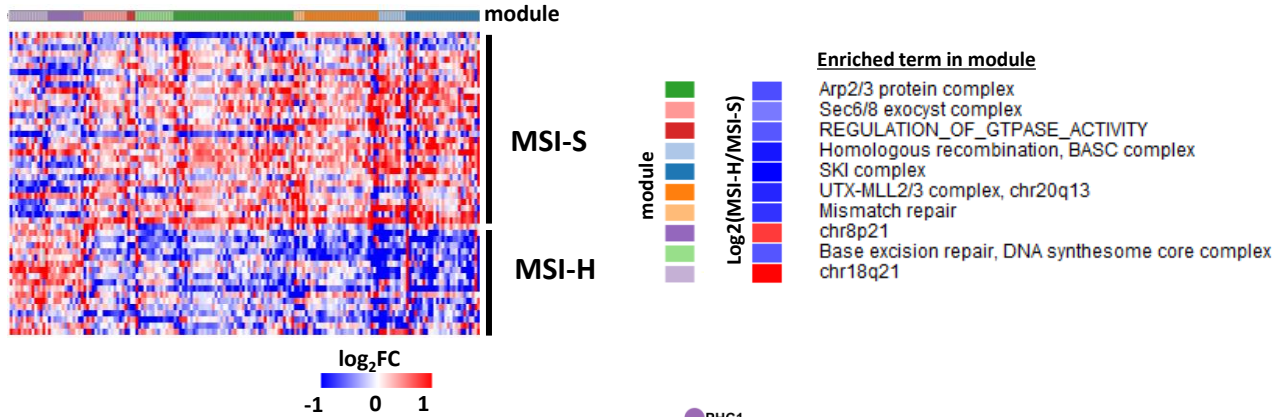
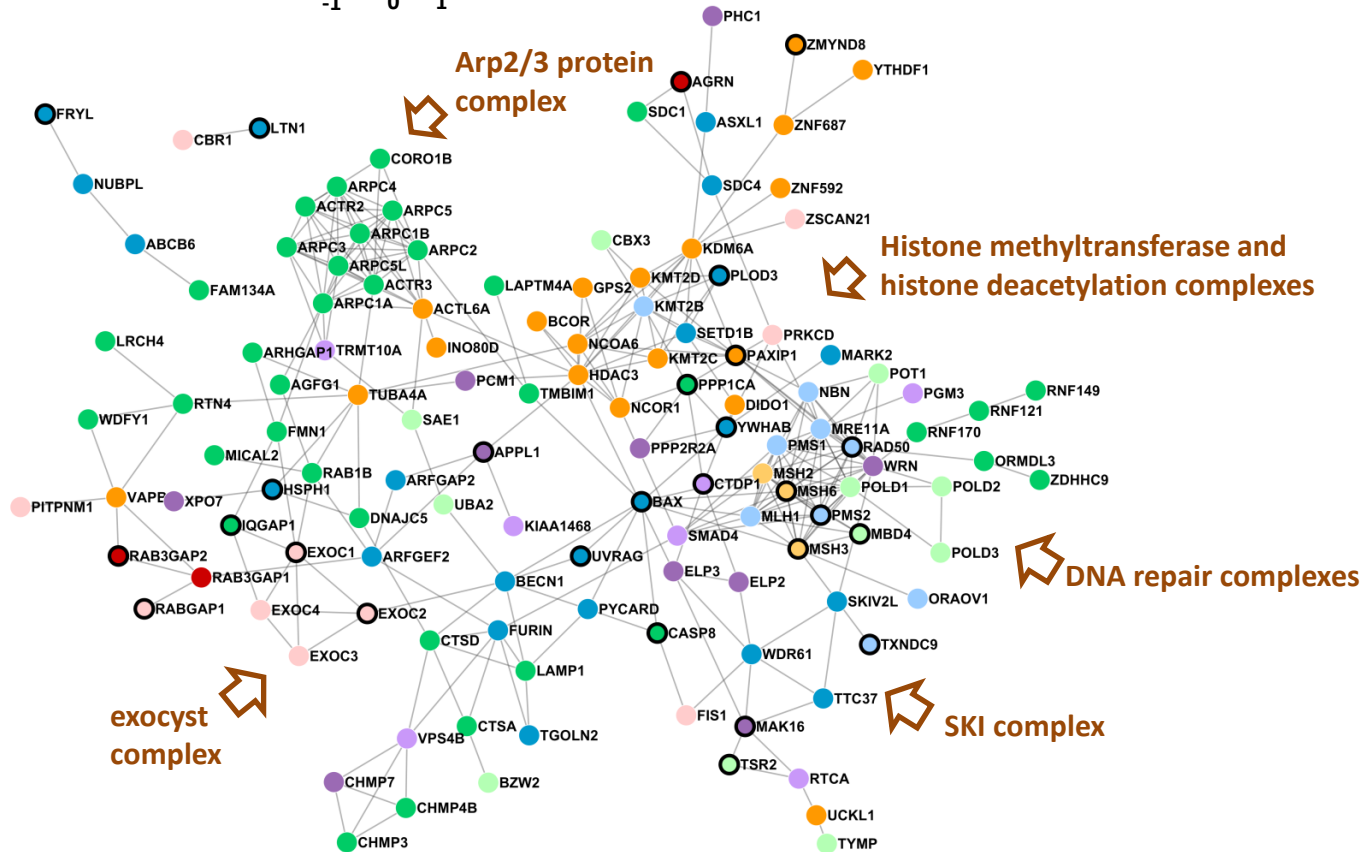


Figure 6

A



B



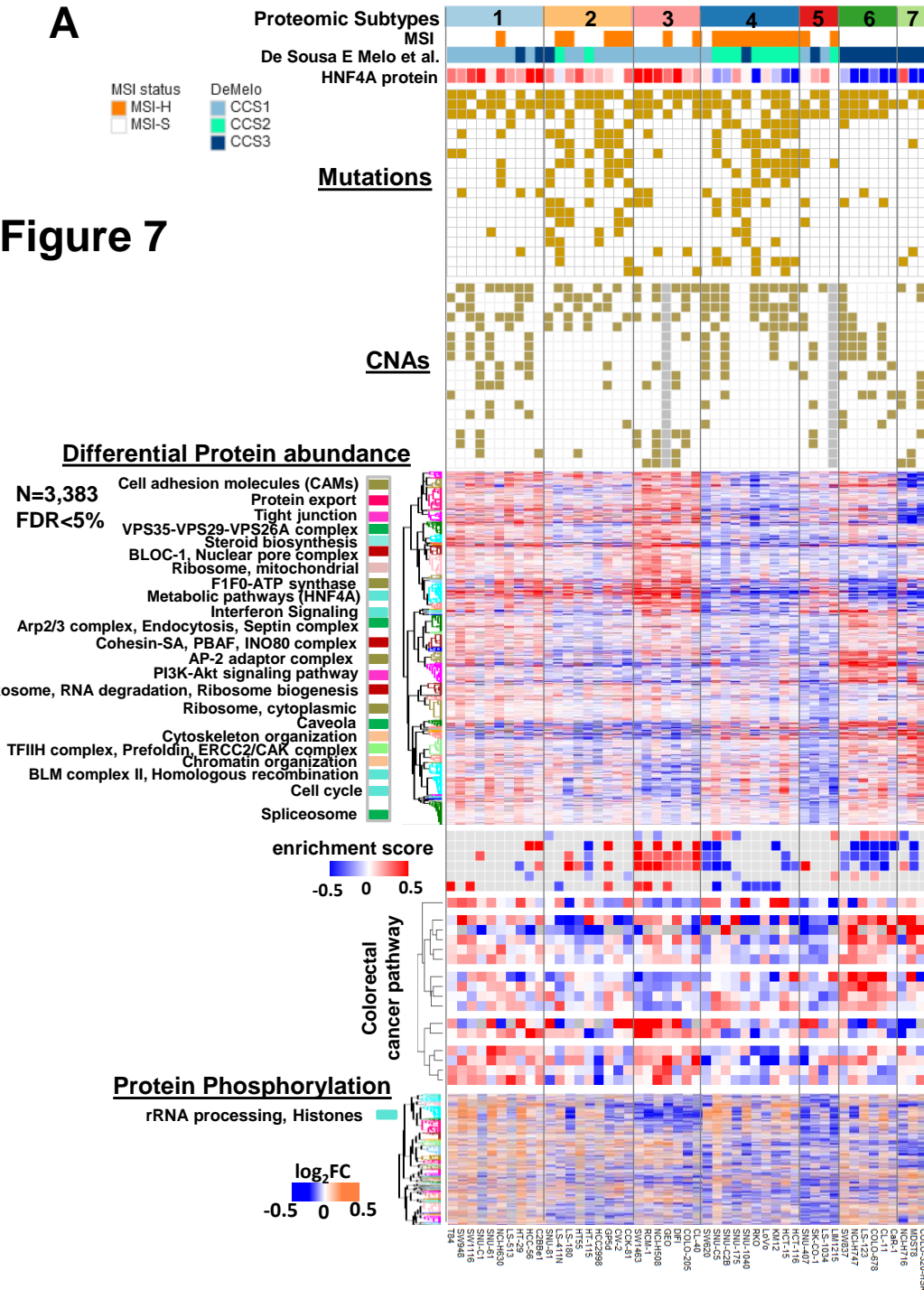
A

Figure 7

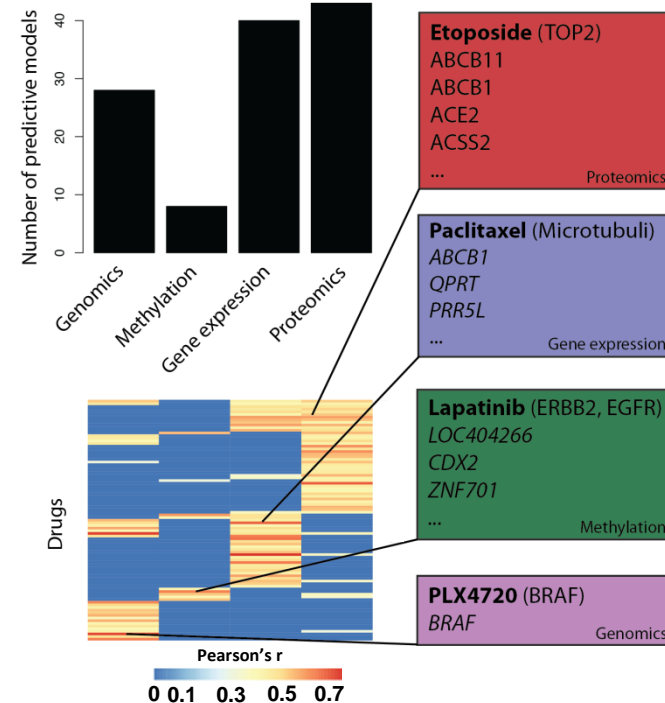
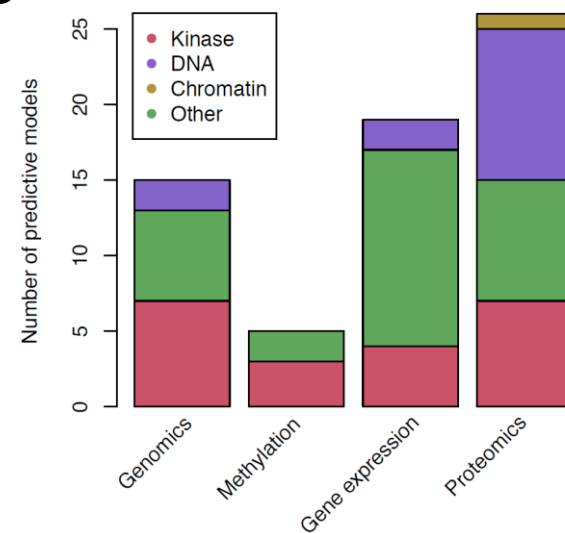
B**C**

Figure S1

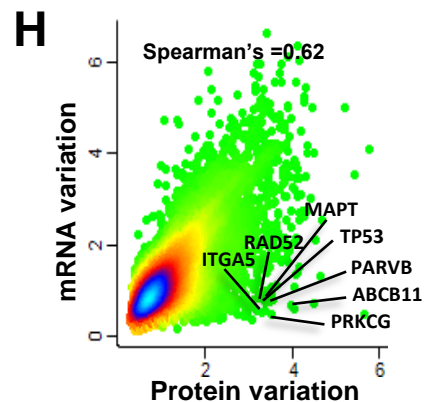
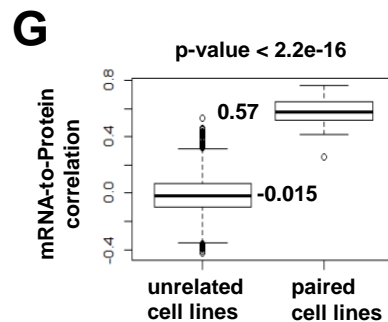
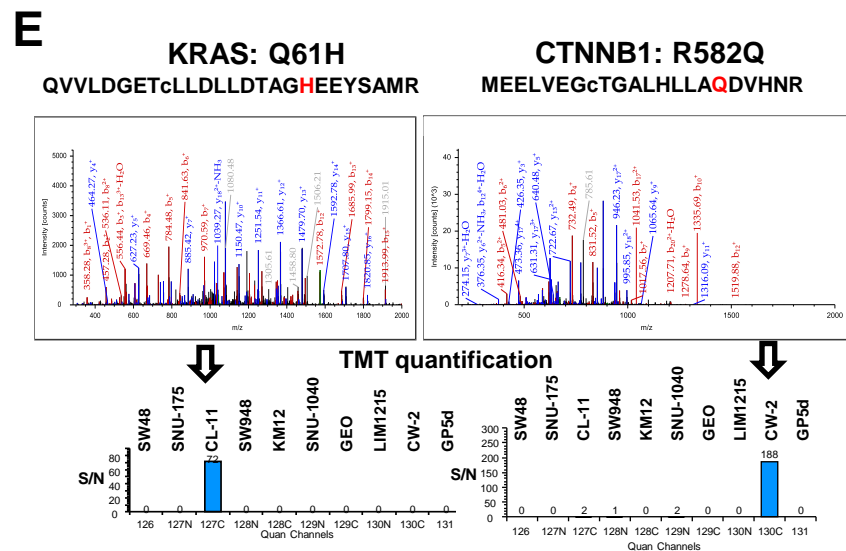
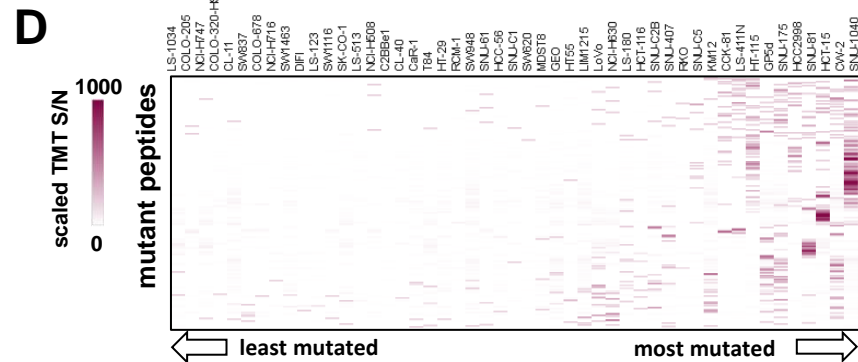
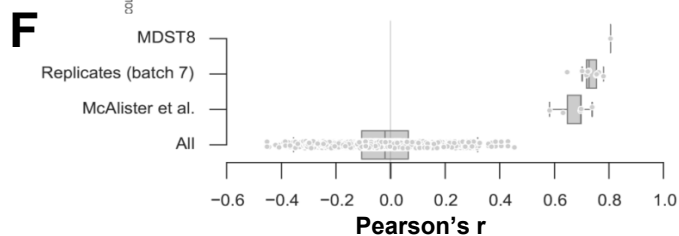
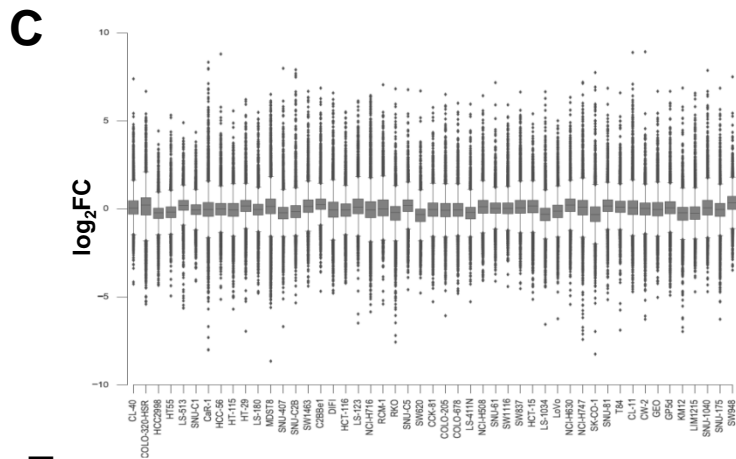
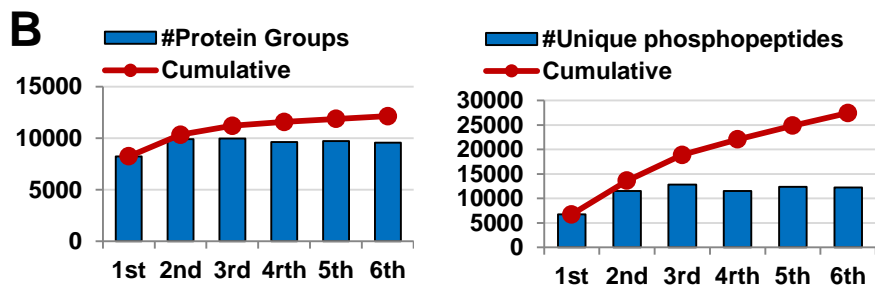
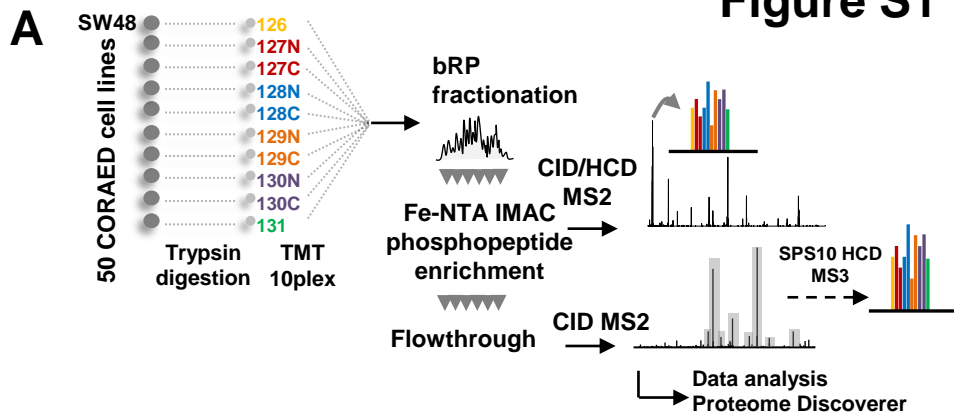
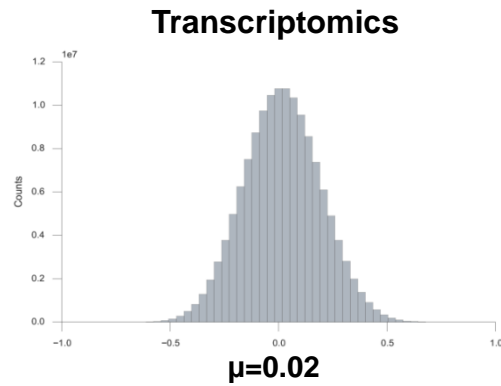
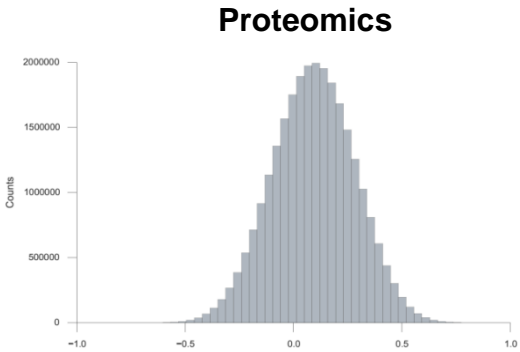


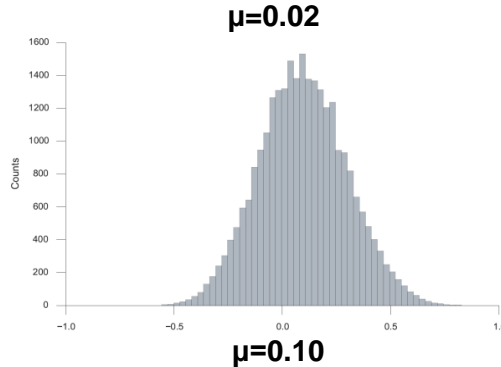
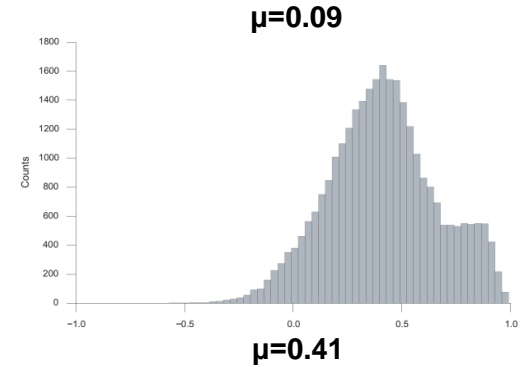
Figure S2

A

All pairs



CORUM



*Spearman's *r*

B

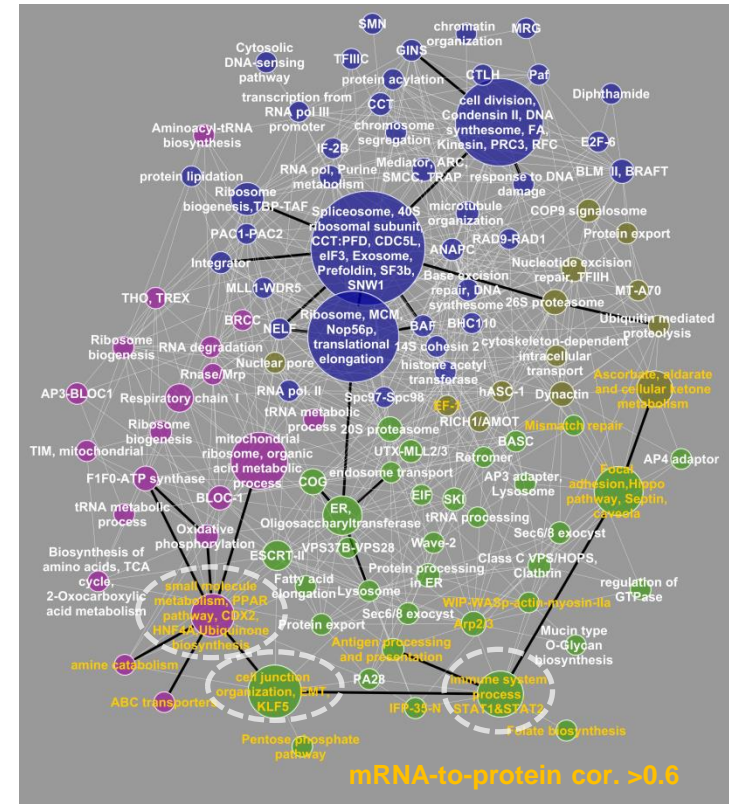


Figure S3

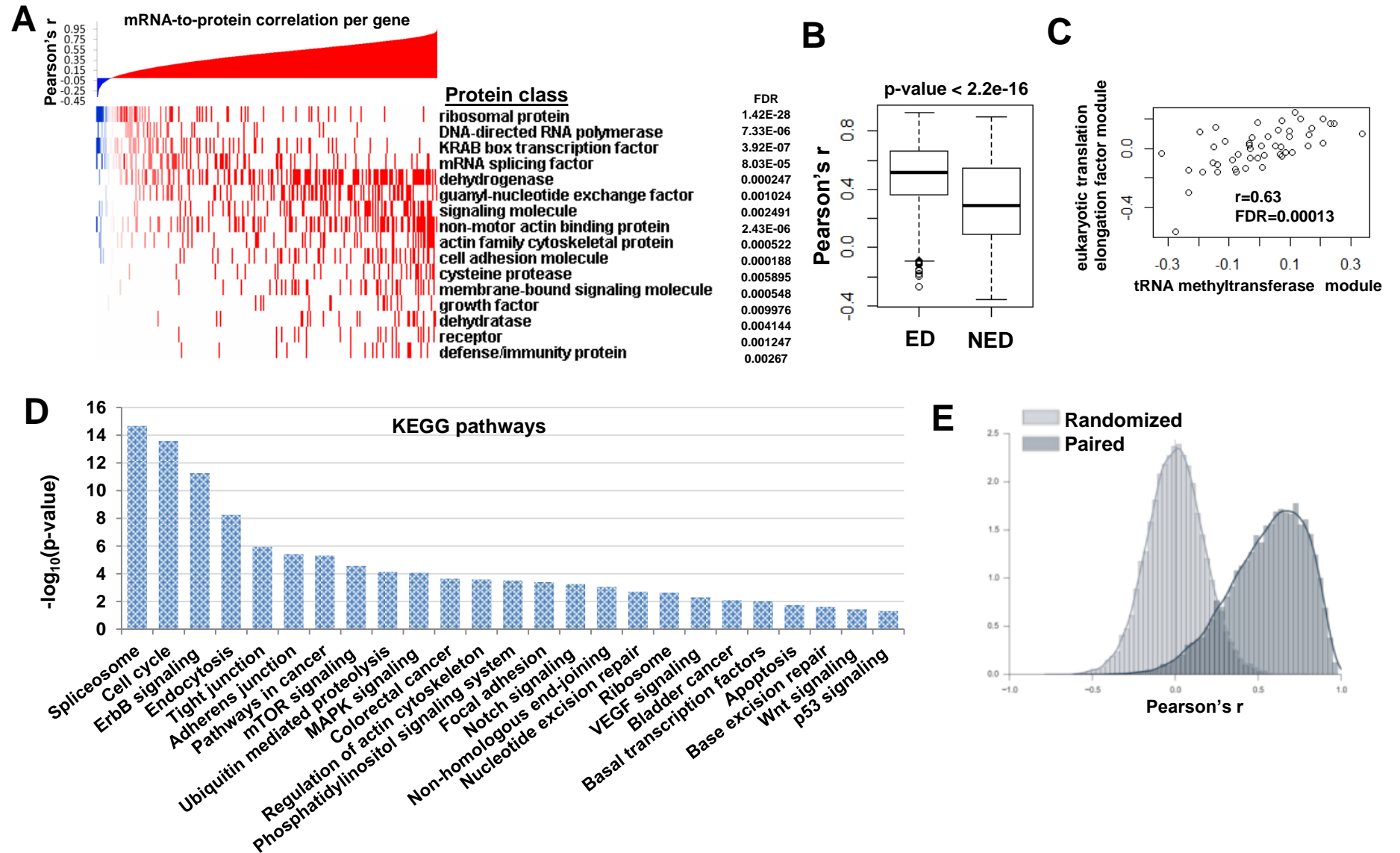
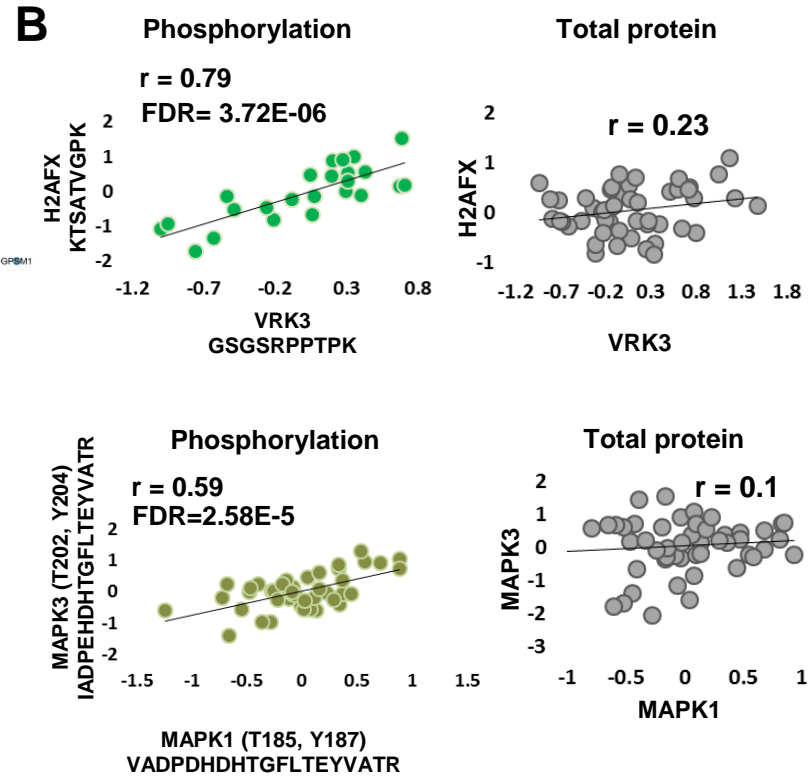
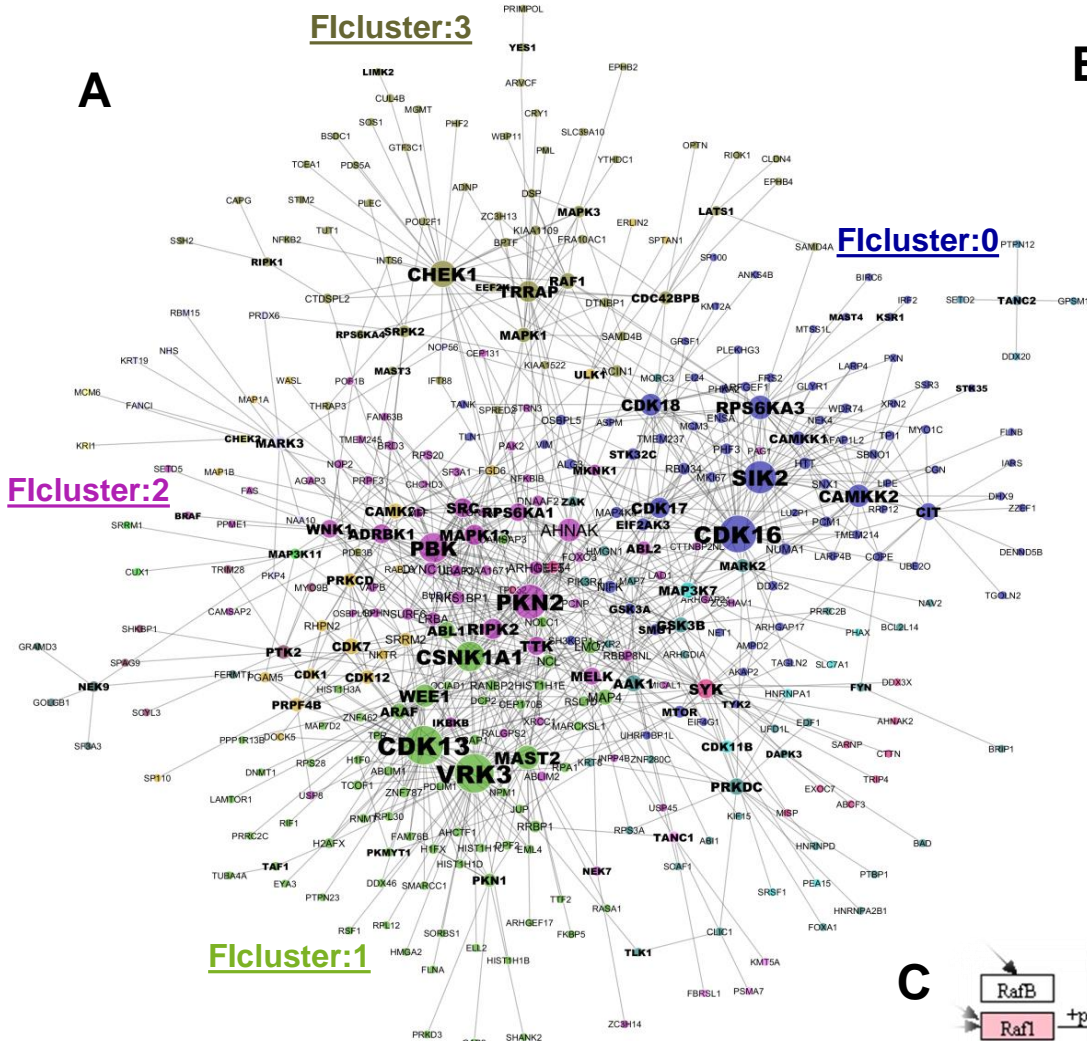


Figure S4



nucleosome assembly: $FDR=5.09E-08$
mitosis: $FDR=0.024$

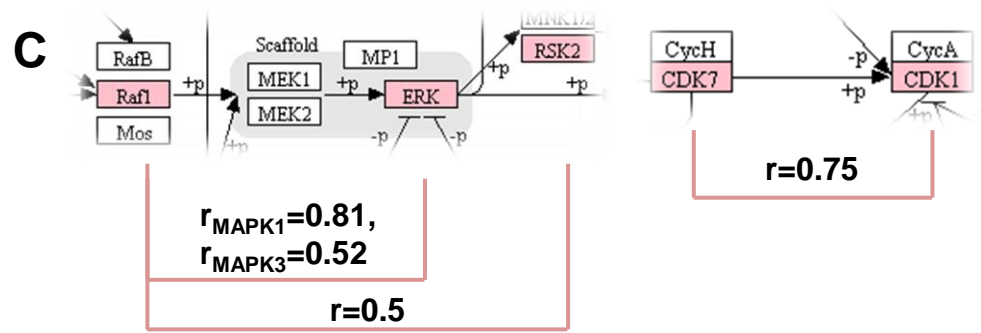


Figure S5

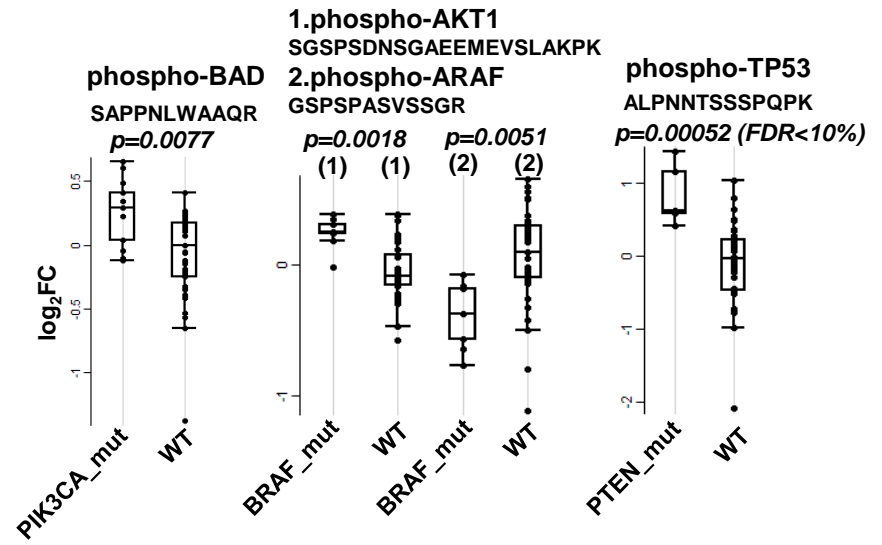
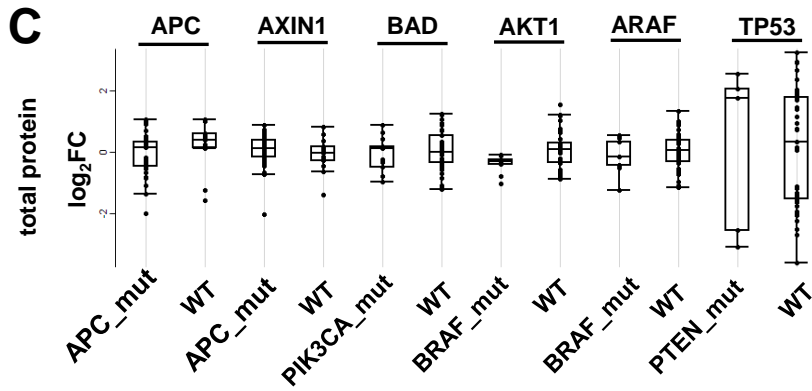
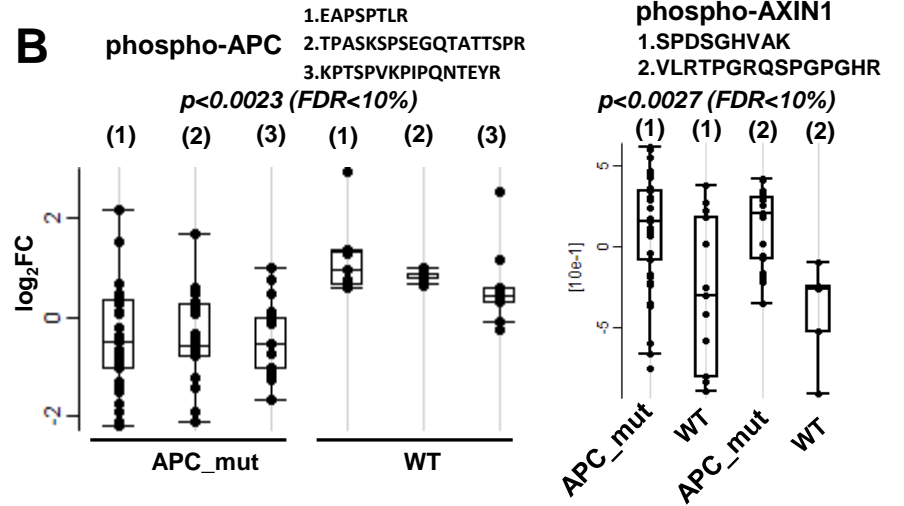
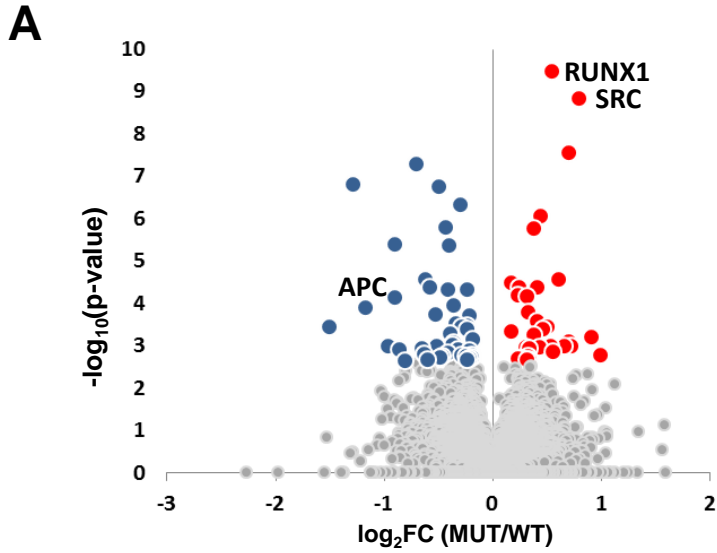
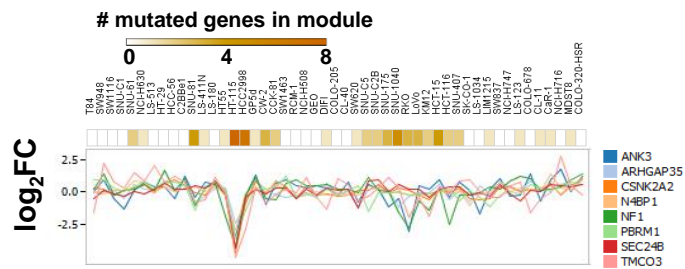
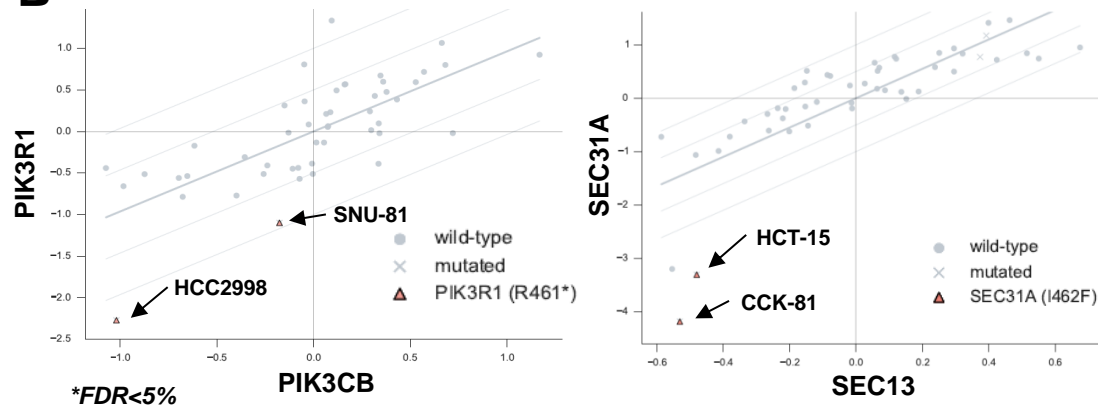


Figure S6

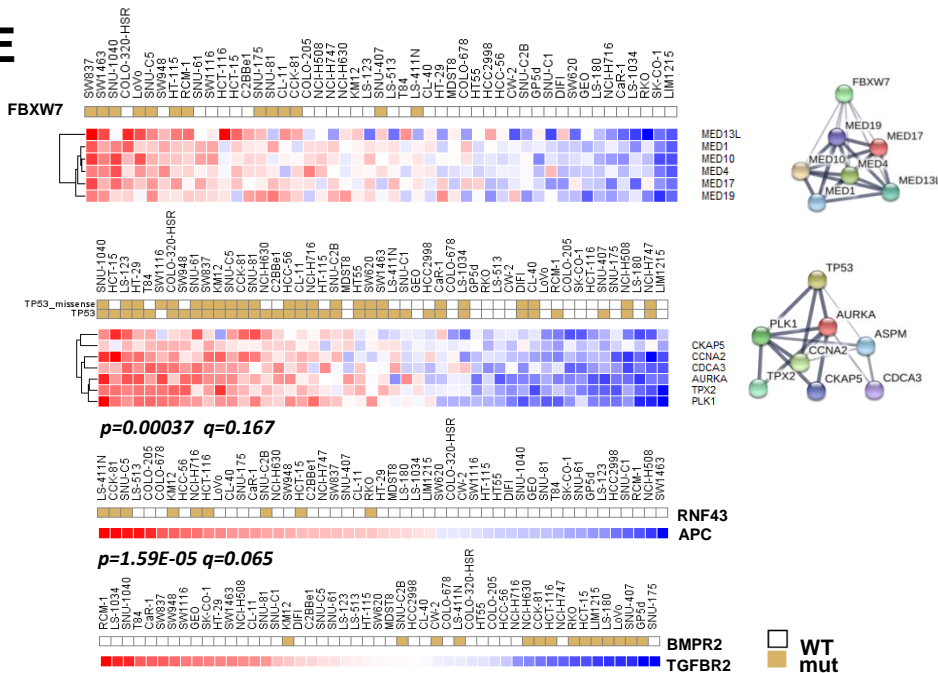
A



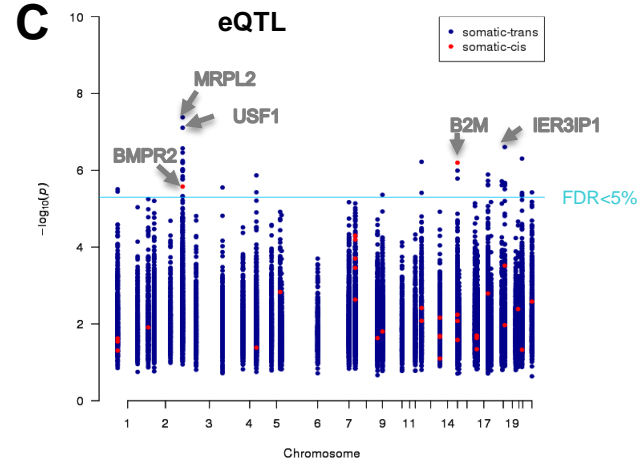
B



E



C



D

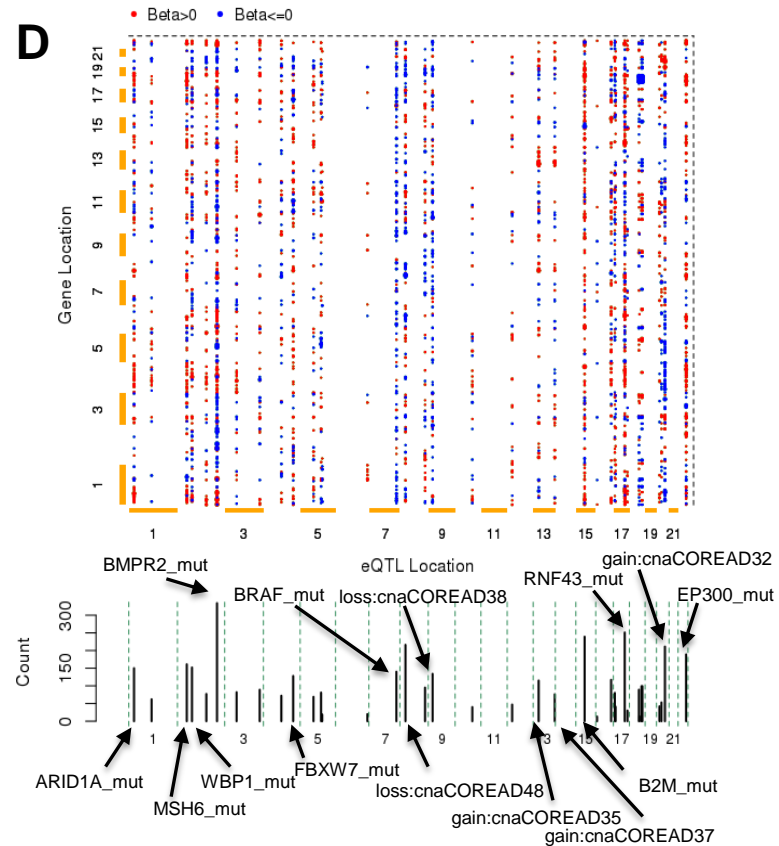
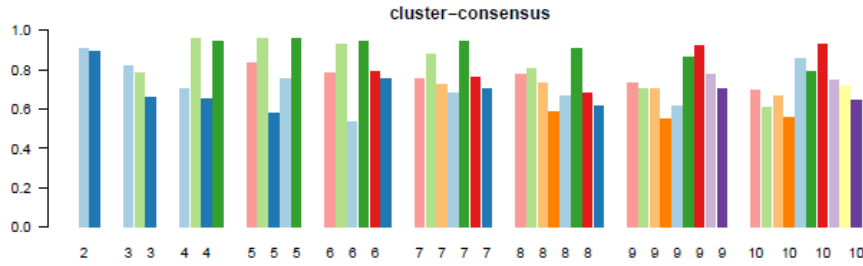
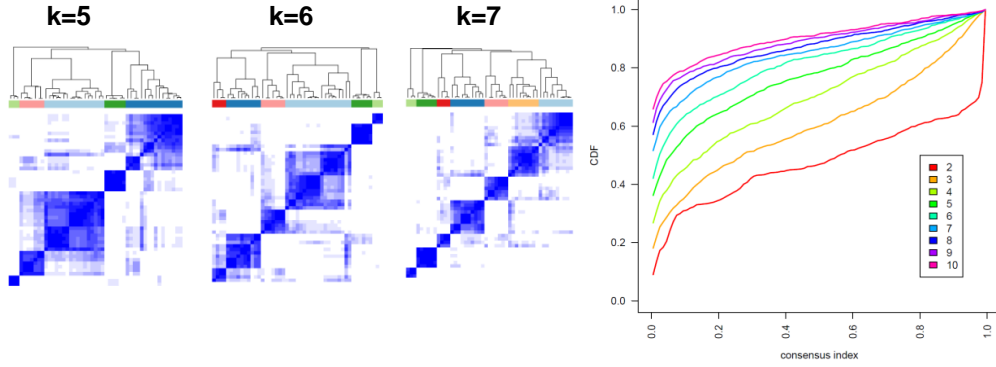
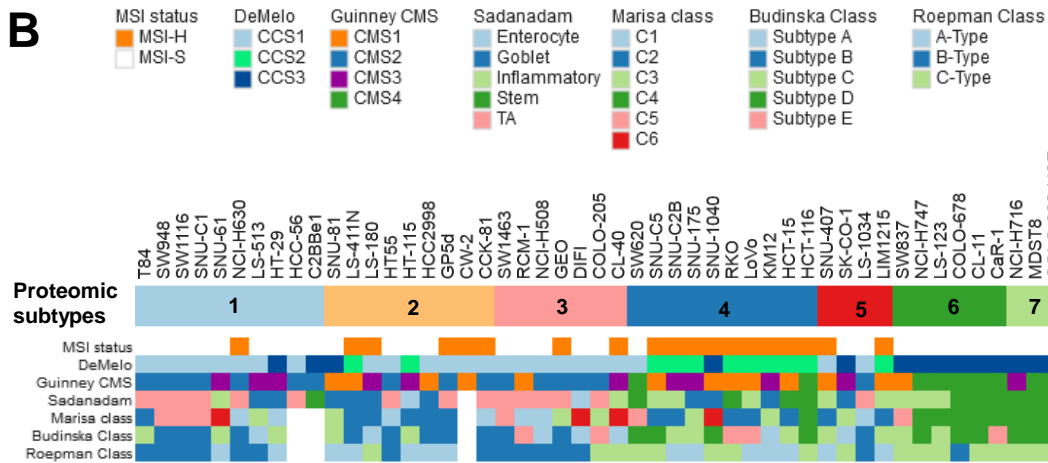
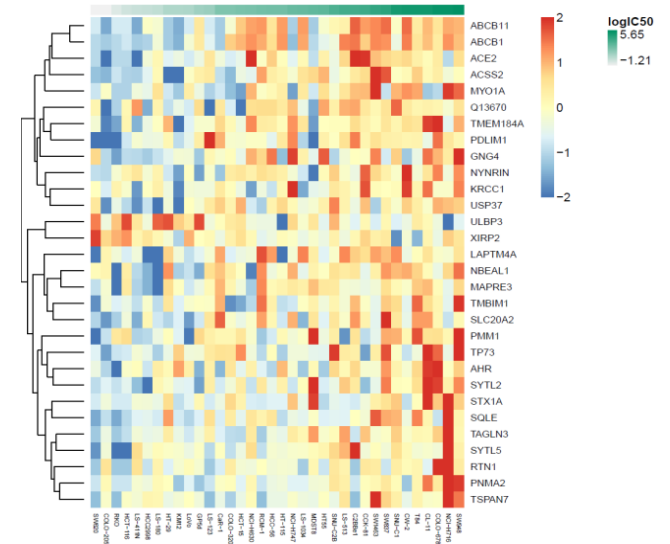
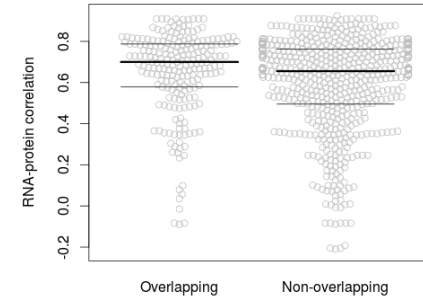


Figure S7

A

B

C

D

E
