

1 **Pharmacoproteomic characterisation of human colon and rectal cancer**

2 Martin Frejino^{1,2}, Riccardo Zenezini Chiozzi^{2,3}, Mathias Wilhelm², Heiner Koch^{2,4,5}, Runsheng Zheng², Susan
3 Klaeger^{2,4,5}, Benjamin Ruprecht^{2,6}, Chen Meng², Karl Kramer², Anna Jarzab², Stephanie Heinzlmeir^{2,4,5}, Elaine
4 Johnstone¹, Enric Domingo^{1,7}, David Kerr⁸, Moritz Jesinghaus⁹, Julia Slotta-Huspenina⁹, Wilko Weichert⁹,
5 Stefan Knapp¹⁰, Stephan M Feller^{11,12,*}, Bernhard Kuster^{2,4,6,13,*}

6 ¹Department of Oncology, University of Oxford, Oxford OX3 7DQ, United Kingdom

7 ²Chair of Proteomics and Bioanalytics, Technical University of Munich, Freising 85354, Germany

8 ³Chemistry Department, University of Rome, La Sapienza, Rome 00185, Italy

9 ⁴German Cancer Consortium (DKTK), Munich 85354, Germany

10 ⁵German Cancer Research Center (DKFZ), Heidelberg 69120, Germany

11 ⁶Center for Integrated Protein Science (CIPSM), Munich 81377, Germany

12 ⁷Wellcome Trust Centre for Human Genetics (WTCHG), University of Oxford, Oxford OX3 7BN, United
13 Kingdom

14 ⁸Nuffield Division of Clinical Laboratory Sciences (NDCLS), University of Oxford, Oxford OX3 9DU, United
15 Kingdom

16 ⁹Institute of Pathology, Technical University of Munich, Munich 81675, Germany

17 ¹⁰Institute of Pharmaceutical Chemistry, Goethe University, Frankfurt am Main 60439, Germany

18 ¹¹Weatherall Institute of Molecular Medicine, University of Oxford, Oxford OX3 9DS, United Kingdom

19 ¹²Institute of Molecular Medicine, Martin-Luther-University, Halle 06120, Germany

20 ¹³Bavarian Biomolecular Mass Spectrometry Center (BayBioMS), Freising 85354, Germany

21 *Corresponding author

22 Correspondence

23 kuster@tum.de

24 stephan.feller@uk-halle.de

25 Final character count including spaces: 76,048

26	Appendix Supplementary Methods	Error! Bookmark not defined.
27	Quantification of LC-MS/MS raw data	3
28	giBAQ quantification.....	3
29	LFQ quantification.....	4
30	Post-processing of quantitative LC-MS/MS data	5
31	Full Proteomes	5
32	Kinobeads	5
33	Reanalysis of public mRNA datasets.....	6
34	Pre-processing.....	6
35	Combination of transcriptomics datasets.....	10
36	Multi-omics data integration strategy	12
37	Computation of protein/mRNA fold-changes.....	12
38	Integration of transcriptomics and proteomics	12
39	mRNA-guided missing value imputation	13
40	Minimum-guided missing value imputation.....	14
41	Combination of the CRC65 and CPTAC datasets.....	16
42	Identification & characterisation of subtypes	16
43	Identification of subtypes.....	16
44	Differential expression analysis.....	17
45	Functional annotation enrichment & clustering.....	18
46	Reanalysis of public dose-response datasets.....	19
47	Pre-processing.....	19
48	Dose-response models & parameter extraction	20
49	Modelling of Drug Sensitivity	21
50	Cox proportional hazards regression	24
51	Annotation information.....	25
52	Appendix Supplementary References	26
53		

54 **Appendix Supplementary Methods**

55 **Quantification of LC-MS/MS raw data**

56 **giBAQ quantification**

57 For full proteomes, we used unique and razor peptides for quantification. MaxQuant was set to calculate iBAQ
58 intensities (Schwanhaussner et al, 2011), however we instead used a modified version of the classical iBAQ
59 approach termed giBAQ to estimate absolute quantities of gene groups as opposed to protein groups, excluding
60 reverse and contaminant hits from the proteinGroups.txt output of MaxQuant. Unless otherwise stated, the
61 CPTAC and CRC65 datasets were processed separately from one another. In order to ensure comparability with
62 transcriptomics data, we assigned every protein group to a single gene group using the following heuristic with
63 four mapping groups (MGs): First, we retrieved a mapping of Uniprot identifiers in the “Majority protein IDs”
64 column to HGNC symbols from Ensembl using biomaRt v2.25.2, collapsing isoforms to their parent Uniprot
65 identifier. This resulted in two mapping groups, one for the Swiss-Prot entries (MG1) and one for the TrEMBL
66 entries (MG2) contained in Uniprot, which were subsequently filtered to remove Swiss-Prot and TrEMBL
67 identifiers mapping to more than one HGNC symbol. For each protein group, we then selected the first Swiss-
68 Prot entry in the “Majority protein IDs” column to represent this protein group and subsequently assigned it the
69 corresponding HGNC symbol from MG1. Afterwards, protein groups left without an HGNC symbol were
70 represented by the first TrEMBL entry in the “Majority protein IDs” column and subsequently assigned the
71 corresponding HGNC symbol from MG2. For protein groups still left without an HGNC symbol, we then used
72 the R package Uniprot.ws v2.10.2 to retrieve a mapping of Uniprot identifiers in the “Majority protein IDs”
73 column to HGNC symbols from Uniprot (MG3), collapsing isoforms to their parent Uniprot identifier. Protein
74 groups left without an HGNC symbol were then represented by the first Uniprot entry in the “Majority protein
75 IDs” column and subsequently assigned the first corresponding HGNC symbol from MG3. Afterwards, all
76 protein groups still left without an HGNC symbol were assigned the first HGNC symbol from the “Gene
77 names” column of the proteinGroups.txt output file (MG4). This heuristic guaranteed that the majority of
78 HGNC-symbol-assignments were done using biomaRt (99% for both the CPTAC and CRC65 datasets), a
79 system which—as opposed to Uniprot.ws or MaxQuant—could also be used to map HGNC symbols to the various
80 identifiers encountered in the transcriptomics datasets described below, ensuring maximum comparability of
81 transcriptomics and proteomics data. We favoured MG1 over MG2 during assignment of HGNC symbols, since
82 it was based on Swiss-Prot entries, which are manually curated. Selecting the first Swiss-Prot (MG1), TrEMBL

83 (MG2) or Uniprot (MG3) entry in the “Majority protein IDs” column to represent a given protein group during
84 HGNC-symbol-assignment in their respective mapping group further made sure that these assignments were
85 always based on the identifier corresponding to the protein with the highest number peptides in the respective
86 protein group. It is worth noting that 96% and 97% of our heuristic HGNC-symbol assignments for the CPTAC
87 patient and CRC65 cell line datasets were identical to the corresponding values in the “Gene names” column of
88 the proteinGroups.txt output file, even though it only served as the source for 1% and 0.1% of the HGNC-
89 symbols, respectively.

90 Next, we extracted the number of theoretical peptides per protein group used by MaxQuant for the calculation of
91 iBAQ intensities as the ratio of raw intensities to iBAQ intensities, rounded to the nearest integer. As expected,
92 these values were constant for each protein group across all samples in the two datasets and were equal to the
93 values determined by *in silico* digestion of the fasta file of UniprotKB used in the database search by MaxQuant
94 with the Protein Digestion Simulator v2.2.5679. Protein groups with only missing intensity values across one of
95 the two datasets or no assigned HGNC-symbol were excluded from the respective dataset. For each gene group
96 as defined by a common HGNC-symbol, we then summed up the raw intensities of protein groups in this gene
97 group in a sample-wise manner and divided them by the sum of the number of theoretical peptides for these
98 proteins, which satisfied the standard iBAQ criteria adapted to our MaxQuant search parameters. These were
99 fully tryptic peptides between 7 and 30 amino acids without missed cleavages, since the minimum peptide
100 length in MaxQuant was set to 7 instead of 6. For a gene group only containing proteins which do not share any
101 theoretical peptides, this quantity termed giBAQ value will be equal to the iBAQ value obtained for a
102 hypothetical protein constructed from them, thereby normalising to the length of the portion of the gene
103 accessible by mass spectrometry, rather than the individual protein lengths.

104 Lfq quantification

105 For Kinobeads experiments, we used unique and razor peptides for quantification. MaxQuant was set to
106 calculate Lfq intensities with an Lfq minimum ratio count of two. Fast Lfq was used to determine
107 normalisation factors, with the minimum number of neighbours set to three, the average number of neighbours
108 set to six and the “Stabilize large Lfq ratios” option enabled. We did not require MS2 spectra for pair-wise
109 peptide intensity comparisons during the calculation of Lfq intensities. Gene names associated with protein
110 groups were remapped as described above for full proteome data. For each gene group and sample, we then

111 summed up the LFQ intensities of protein groups in this gene group to generate a gene-level quantification,
112 which was used in all subsequent analyses.

113 Post-processing of quantitative LC-MS/MS data

114 Full Proteomes

115 The CRC65 and CPTAC datasets were processed separately from one another. For the CPTAC dataset, we
116 restricted our analysis to CPTAC sample IDs not removed in the original publication (Zhang et al, 2014) due to
117 duplication. Following common practice, giBAQ values were log₂-transformed and median-centred. Since the
118 CRC65 cell line panel contained two pairs of concordant cell lines (HDC-54/HDC-57 and LS 180/LS 174T;
119 Bracht et al, 2010; Klijn et al, 2015; Medico et al, 2015), we replaced missing values in HDC-54 and LS 180
120 with measured values from HDC-57 and LS 174T, respectively. HDC-57 and LS 174T were subsequently
121 excluded from the analysis. It is worth noting that we did not just calculate the arithmetic mean of these cell line
122 pairs, since this would prevent minimum-guided missing value imputation on the peptide level with raw-file-
123 specific backgrounds, as described in its respective section.

124 Kinobeads

125 Since we observed batch effects between the three different biological replicates, we decided to log₂-transform
126 and median-centre them separately from one another, followed by batch-effect removal using ComBat (Johnson
127 et al, 2007). Briefly, we first restricted our dataset to gene groups consistently measured between the different
128 replicates and across the dataset, selecting gene groups with at least two out of three reported LFQ values per
129 cell line in at least 23 cell lines. Next, we removed all gene groups with one or fewer reported LFQ values per
130 replicate in order to eliminate noise and enable subsequent parametric batch adjustment using ComBat,
131 implemented in the sva R package v3.18.0, with MSI status as a covariate as described by Guinney and
132 colleagues (Guinney et al, 2015). We used the median across all three replicates to summarise the LFQ values
133 for each cell line, taking the median across all six replicates for the two pairs of concordant cell lines described
134 above; HDC-54 was selected to represent HDC-54/HDC-57 and LS 180 was selected to represent
135 LS 180/LS 174T.

136 Reanalysis of public mRNA datasets

137 Pre-processing

138 We downloaded 10 mRNA datasets (GSE36133 - Barretina et al, 2012; E-MTAB-783 - Garnett et al, 2012; E-
139 MTAB-2706 - Klijn et al, 2015; GSE28567 - Loboda et al, 2011; GSE59857 - Medico et al, 2015;
140 Supplementary Table S6 - Mouradov et al, 2014; platform codes "IlluminaGA_RNASeq" +
141 "AgilentG4502A_07_3" - TCGA Network, 2012; GSE8332 - Wagner et al, 2007; GSE24795 - Wilding et al,
142 2010, respectively) and processed them as follows, automatically obtaining the required chip definition files
143 (.cdf) from Bioconductor for microarray datasets, if not stated otherwise. The same cell lines have inconsistent
144 names in the different transcriptomics datasets, which is why the first pre-processing step always involved the
145 mapping of all cell line names to a set of consistent sample names (see also Table EV6A).

146 GSE36133

147 We used GEOquery v2.35.6 to download annotation information associated with GSE36133 and obtained raw
148 .CEL files, as well as the "Expression arrays samples info file" from the CCLE portal of the Broad Institute
149 (<http://portals.broadinstitute.org/ccle/data/browseData?conversationPropagation=begin>, 31 December 2015).
150 The dataset was restricted to cell lines from the large intestine, as well as C32 and Colo 741, which were derived
151 from skin according to the "Expression arrays samples info file". The `preproPara()` function from the `affyPara`
152 package v1.29.0 was then used with background correction through RMA, quantile normalisation and probe
153 summarisation through median polish of perfect match (PM) probes to calculate the classical RMA expression
154 measure. Probes with only missing values across the dataset were excluded. Next, we retrieved a mapping of
155 AffyIDs to HGNC symbols from Ensembl using `biomaRt` and filtered it to remove AffyIDs, which mapped to
156 more than one HGNC symbol. This mapping was then used to assign each AffyID to its corresponding HGNC
157 symbol. AffyIDs we were not able to map this way were excluded from the dataset. This was one of two
158 datasets for which we were not able to first map primary IDs (in this case AffyIDs) to GeneIDs with the default
159 annotation data downloaded from GEO or related sources. For consistency reasons, GeneIDs served as the basis
160 for HGNC symbol assignment during the analysis of other transcriptomics datasets whenever possible. For
161 GSE36133, this resulted in an expression matrix of 36,877 AffyIDs across 63 cell lines.

162 E-MTAB-783

163 We obtained raw .CEL files, as well as the corresponding `eSet.r` and `adf.txt` files from ArrayExpress
164 (<https://www.ebi.ac.uk/arrayexpress/experiments/E-MTAB-783/files/>, 31 December 2015) and downloaded the

165 “gdsc_manova_input_w5.csv” file from the GDSC portal in order to be able to map cell lines to tissues of origin
166 (<ftp://ftp.sanger.ac.uk/pub4/cancerrxgene/releases/release-5.0/>, 12 March 2016). The dataset was restricted to
167 cell lines from the large intestine, as well as C32 and SW626, which were derived from skin and ovary,
168 respectively, according to the tissue labels in the “gdsc_manova_input_w5.csv” file. The rest of the pre-
169 processing was similar to GSE36133, with the only difference being that we first mapped AffyIDs to Ensembl
170 GeneIDs using the adf.txt file and afterwards assigned each GeneID to an HGNC symbol using biomaRt as
171 described above. This resulted in an expression matrix of 20,106 AffyIDs across 37 cell lines.

172 E-MTAB-2706

173 We obtained Supplementary Table 1 from the original publication by Klijn et al., as well as RPKM data for
174 coding genes from the website accompanying the publication (<http://research-pub.gene.com/KlijnEtAl2014/>, 06
175 October 2015). Supplementary Table 1 was used to assign cell lines to their tissues of origin and the dataset was
176 subsequently restricted to colorectal cancer cell lines as well as C32, which was derived from skin. RPKM
177 values equal to zero were treated as missing values and GeneIDs only containing missing values were removed
178 from the dataset. Afterwards, we quantile normalised and log₂-transformed the data, followed by assigning each
179 GeneID to an HGNC symbol using biomaRt as described above. This resulted in an expression matrix of 20,965
180 GeneIDs across 56 cell lines.

181 GSE28567

182 We obtained annotation information for GSE28567 through GEOquery and downloaded the appropriate custom
183 chip definition file (.cdf) and the corresponding raw .CEL files from NCBI GEO (Edgar et al, 2002; 06
184 November 2015). The R package makecdfenv v1.46.0 was used to construct an annotation package from the
185 .cdf file, which was subsequently used to relate probes to locations on the microarray. The rest of the pre-
186 processing was similar to GSE36133, with the only difference being that we first mapped the custom probe IDs
187 to Ensembl GeneIDs using the annotation information from GEO and afterwards assigned each GeneID to an
188 HGNC symbol using biomaRt as described above. This resulted in an expression matrix of 33,749 probe IDs
189 across 67 cell lines.

190 GSE59857

191 We used GEOquery to obtain a normalised expression set object and annotation information associated with
192 GSE59857 and downloaded non-normalised expression data from NCBI GEO (04 November 2015), as well as
193 further annotation information for Illumina expression arrays from Bioconductor. The R packages limma

194 v3.26.2 and beadarray v2.19.1 were used for background correction, normalisation and transformation of raw
195 data. We first generated an ExpressionSetIllumina from the expression set object we downloaded from NCBI
196 GEO and replaced the normalised expression data with the non-normalised data. Subsequently, we performed
197 background correction using the backgroundCorrect function from limma, estimating the parameters of the
198 normal-exponential (normexp) convolution model adapted for microarrays other than Affymetrix arrays by
199 Ritchie et al (2007) using the maximum-likelihood estimator developed by Silver et al (2009) with an offset of
200 10 in order to ensure that expression values were positive after background correction. Afterwards, the corrected
201 expression data were quantile normalised using the normaliseIllumina function from the beadarray package,
202 followed by log₂-transformation of the dataset. This method was used in favour of the neqc method described
203 by Shi and colleagues (Shi et al, 2010) after comparison of the distributions of expression values produced by
204 the two algorithms. We restricted the dataset to probes from coding regions, excluding negative control and bad
205 quality probes, but keeping probes associated with GeneIDs contained in the single-sample predictor for the
206 CMS subtypes published by Guinney et al (2015) regardless of probe quality annotation. Ensembl GeneIDs
207 were mapped to HGNC symbols using biomaRt as described above. This resulted in an expression matrix of
208 27,293 probe IDs across 155 cell lines.

209 Supplementary Table S6 from Mouradov et al.

210 We downloaded (05 August 2014) Supplementary Table S6 from the original publication by Mouradov et al
211 (2014), excluded cell line RW2982 from the dataset as the only cell line not present in the CRC65 panel, treated
212 RPKM values equal to zero as missing values and restricted the dataset to RefSeq IDs containing at least one
213 measured value. Subsequently, RPKM values were quantile normalised and log₂-transformed, followed by re-
214 mapping RefSeq IDs to HGNC symbols using biomaRt as described above. This resulted in an expression
215 matrix of 29,973 RefSeq IDs across 12 cell lines.

216 Platform code "IlluminaGA_RNASeq" from TCGA

217 We used the TCGAbiolinks R package v1.0.4 to obtain and prepare level 3 data associated with COAD/READ
218 tumours and platform code "IlluminaGA_RNASeq", for which we also had proteomics data from the CPTAC
219 study on colorectal cancer (Zhang et al, 2014) and which were not excluded from this dataset due to duplication
220 as described above (16 December 2015). We removed duplicate rows from the RPKM expression matrix,
221 remapped the TCGA barcodes to the sample names described in Table EV6A, treated RPKM values equal to
222 zero as missing values and restricted the dataset to GeneIDs containing at least one measured value, as described

223 above. Afterwards, the dataset was log₂-transformed and median-centred, followed by re-mapping of GeneIDs
224 to HGNC symbols using biomaRt as described above. This resulted in an expression matrix of 19,670 GeneIDs
225 across 87 tumours.

226 Platform code "AgilentG4502A_07_3" from TCGA

227 We again used the TCGAbiolinks R package v1.0.4 to obtain and prepare level 1 data and MAGE-TAB files
228 associated with COAD/READ tumours and platform code "AgilentG4502A_07_3", for which we also had
229 proteomics data from the CPTAC study on colorectal cancer (Zhang et al, 2014; 25 January 2016). The
230 corresponding platform design file (.adf files) was downloaded from the TCGA portal at [https://tcga-](https://tcga-data.nci.nih.gov/docs/publications/tcga/platformdesign.html)
231 [data.nci.nih.gov/docs/publications/tcga/platformdesign.html](https://tcga-data.nci.nih.gov/docs/publications/tcga/platformdesign.html) (26 January 2016). We used limma for the
232 processing of txt files generated from microarray scans using Agilent's Feature Extraction Software (TCGA
233 Network, 2011). Briefly, an RGList was generated, followed by annotation of the probes using the platform
234 design file in order to perform quality control through annotated MA-plots, boxplots of background intensities,
235 image plots of foreground and background intensities, density plots and the QC metrics calculated by the R
236 package arrayQualityMetrics v3.26.1. Subsequently, we performed background correction using the
237 backgroundCorrect function from limma with the method described by Edwards (2003) and an offset of 50,
238 since it produced the most stable results as judged by inspection of MA-plots following within-array loess
239 normalisation with standard parameters for several background correction methods. Afterwards, between-array
240 normalisation using the "Gquantile" method implemented in the normalizeBetweenArrays function of limma
241 was carried out under the assumption that the signal distributions of the green channels across all arrays should
242 be the same, since the green channel always contained the Universal Human Reference RNA (TCGA Network,
243 2012). We took the median of replicate probes, filtered out probes without GeneID annotation and finally
244 mapped the GeneIDs to HGNC symbols using biomaRt as described above. This resulted in an expression
245 matrix of 66,728 Agilent IDs across 73 tumours.

246 GSE8332

247 We used GEOquery to download annotation information associated with GSE8332 and also obtained raw .CEL
248 files from NCBI GEO (20 November 2015). The dataset was pre-processed similarly to GSE36133 and
249 subsequently restricted to CRC cell lines based on Figure 1 of the original publication (Wagner et al, 2007), as
250 well as C32 and Colo 741, which were derived from skin as described above. It is worth noting that the authors
251 categorised Colo 741 as a CRC cell line, which contradicts the annotations found in the other mRNA datasets

252 analysed thus far. Using the annotation information downloaded from GEO, AffyIDs were then annotated with
253 GeneIDs, which were subsequently mapped to HGNC symbols using biomaRt as described above. This resulted
254 in an expression matrix of 39,546 AffyIDs across 36 cell lines.

255 GSE24795

256 Raw data for GSE24795 were obtained (05 November 2015) and pre-processed as described for GSE8332, with
257 the only difference that no samples needed to be omitted, since all cell lines included in the study were CRC cell
258 lines. This resulted in an expression matrix of 39,546 AffyIDs across 30 cell lines.

259 Combination of transcriptomics datasets

260 We aggregated the different transcriptomics datasets described above into one expression matrix using a scheme
261 similar to the one published by Guinney and colleagues (Guinney et al, 2015). As mentioned in the original
262 publication, we expected strong batch effects between the different transcriptomics datasets and also needed to
263 select a single probe or primary ID to represent the expression of a given gene in each dataset in order to avoid
264 inconsistent mRNA quantification. But since we intended to compare the abundance of as many mRNA species
265 to as many protein species across as many samples as possible, respectively, we refrained from removing outlier
266 samples from each dataset separately and also did not restrict the list of quantifiable mRNA species to the ones
267 with the largest median absolute deviation (MAD). We also did not select a reference dataset *a priori*, but rather
268 computed similar consistency measures as the ones described by Guinney and colleagues for each dataset and
269 selected the reference dataset afterwards, thereby maximising the consistency between the different datasets as
270 described below. We started out by reducing all datasets to a common set of 10,044 reference genes based on
271 the HGNC-symbol-assignments described in the previous sections. Next, each of the reference genes was
272 temporarily represented by the primary identifier (AffyID, Probe ID, etc.) with the largest MAD in each of the
273 datasets (reference identifiers). For each dataset, we then calculated the correlation between all primary
274 identifiers and these reference identifiers, generating ten correlation matrices with a number of rows equal to the
275 number of primary identifiers in the respective dataset and 10,044 columns each; one for each reference
276 identifier. For each dataset, this correlation matrix is equivalent to the collection of correlation vectors C of
277 length $|G_{REF}|$ mentioned in the publication by Guinney and colleagues. For each of the reference genes and each
278 dataset, we then calculated the correlation of all correlation vectors associated with primary identifiers
279 quantifying the respective reference gene to all other correlation vectors from the different datasets associated
280 with primary identifiers quantifying the same reference gene. This resulted in a collection of 10,044 “gene-

281 lists”, one for each reference gene. Each of these “gene-lists” itself consisted of ten “reference-dataset-lists”, one
282 for each dataset. Each of these “reference-dataset-lists” itself contained a “dataset-list” of correlation matrices,
283 one for each dataset. These correlation matrices finally contained the correlations of all correlation vectors
284 associated with primary IDs quantifying the respective reference gene in the respective dataset against all
285 corresponding correlation vectors associated with primary IDs quantifying the same reference gene in all other
286 datasets. This enabled us to avoid having to select a reference dataset *a priori* as described by Guinney et al.,
287 thereby allowing different selection criteria to be consulted before deciding on a reference dataset. In order to
288 select a reference dataset, we used the following heuristic: For each potential reference dataset and each primary
289 ID in it, we calculated a consistency metric *comet* by first determining which primary IDs in the other datasets
290 quantifying the same reference gene show the highest “correlation of correlations” value with the respective
291 primary ID in the potential reference dataset and subsequently summing up these “correlations of correlations”
292 for each primary ID in the potential reference dataset. For each dataset and reference gene, this resulted in a
293 vector of *comet* values, one for each primary ID quantifying the respective reference gene. The primary ID with
294 the maximum *comet* value (*maxcomet* value) was selected to represent the respective reference gene in the
295 respective dataset. These primary IDs were the ones resulting in the highest consistency between the different
296 datasets if we were to select the respective dataset as a reference. We then compared the distribution of these
297 *maxcomet* values across all datasets in order to find the reference dataset resulting in the highest consistency
298 between the different datasets. Four of the ten datasets showed similarly promising distributions of *maxcomet*
299 values (E-MTAB-2706, GSE28567, GSE8332 and GSE36133) and we decided to use GSE36133 as the
300 reference dataset, since it was the biggest dataset in terms of both primary IDs and cell lines. In this reference
301 dataset, primary IDs associated with *maxcomet* values were selected to represent their corresponding reference
302 genes. For each of the other datasets, the primary ID with the maximum “correlation of correlations” value with
303 the primary ID chosen to represent a given reference gene in the reference dataset was selected to represent the
304 respective reference gene in the respective dataset. This resulted in ten datasets with 10,044 genes each, which
305 were represented by a combination of primary IDs ensuring maximum consistency between the different
306 datasets.

307 We next merged the eight datasets quantifying transcription in cell lines into one expression matrix and also
308 generated a second expression matrix containing all ten datasets. These two expression matrices were treated
309 separately from one another during adjustment of batch effects in order to maximise the number of genes in
310 each of them. First, we reduced these matrices to genes quantified in at least two samples per dataset, followed

311 by the adjustment of batch effects between the different datasets using ComBat (Johnson et al, 2007) with MSI
312 status as a covariate. Next, we restricted the expression matrix containing all ten datasets to samples also present
313 in the CPTAC patient dataset. At this stage, we had two mRNA expression matrices, which—together with the
314 giBAQ quantitation for cell lines and patients described above—served as the basis for our multi-omics data
315 integration strategy.

316 Multi-omics data integration strategy

317 Computation of protein/mRNA fold-changes

318 Wilhelm et al. showed that protein/mRNA ratios or fold-changes are reasonably stable across a number of
319 different tissues and that median ratios can be used to estimate protein abundance from mRNA abundance for a
320 given protein (Wilhelm et al, 2014). Since these fold-changes were calculated based on RPKM values as
321 determined by mRNA-Seq, we were curious as to whether this also holds true when transcript abundance is
322 estimated using microarrays. We therefore calculated the median expression for each transcript and sample in
323 the two mRNA expression matrices described above and divided the giBAQ values of the corresponding
324 proteins in the CRC65 cell line and CPTAC patient dataset by these values in a sample-wise manner. This
325 resulted in two ratio matrices, one for the CRC65 dataset and one for the CPTAC dataset, which were used to
326 produce Figures EV3A & B. The median mRNA expression values were the basis for Figure EV2B. These
327 figures showed that there are systematic differences between protein and mRNA datasets, which manifest as
328 reasonably stable protein/mRNA fold-changes within the CRC65 and CPTAC datasets. Notably, the greater
329 proteomic ‘depth’ of the CRC65 dataset resulted in systematically higher protein/mRNA fold-changes compared
330 to the CPTAC dataset. This could not have been adjusted using global total-sum normalisation, since the peptide
331 coverage differed drastically between the two datasets in a non-linear way (Fig. EV2C).

332 Integration of transcriptomics and proteomics

333 Because protein/mRNA fold-changes were reasonably stable within both datasets, we hypothesised that they
334 could be treated as systematic differences and therefore be adjusted using e.g. ComBat as described above.
335 ComBat however shifts genes in each dataset to the overall mean and pooled variance across datasets, thereby
336 altering the location and scale of each dataset. In our case, we wanted to avoid altering the giBAQ expression
337 values as they represent actual measurements of gene group abundance on the protein level. Therefore, we
338 decided to instead use MComBat (Stein et al, 2015), which allows the specification of a “gold-standard” dataset

339 towards which all other datasets are adjusted and which is not changed. MComBat takes advantage of the
340 systematic differences between protein and mRNA measurements and models them in a protein-wise fashion.
341 Stein and colleagues published the R code for MComBat alongside their manuscript at
342 <https://github.com/SteinCK/M-ComBat/blob/master/M-ComBat%20R%20Script>, however this script is setup to
343 stop if the combined expression matrix used as input contains any number of missing values. We modified this
344 script to instead work similarly to the original ComBat function implemented in the sva R package v3.18.0 with
345 respect to missing values; the modified script is available at
346 https://github.com/mfrejno/pharmacoproteomics_crc. With this modified MComBat function at hand, we set out
347 to integrate transcriptomics and proteomics data for the CRC65 cell line and CPTAC patient datasets separately
348 from one another. After merging proteomics and transcriptomics data for each dataset separately, we reduced
349 these matrices to gene groups quantified in at least two samples on both the protein and transcript level,
350 followed by the adjustment of systematic differences between proteomics and transcriptomics data using the
351 aforementioned modified MComBat function without MSI status as a covariate. The resulting adjusted
352 expression matrices were the basis for Figure EV3C, as well as for mRNA-guided missing value imputation.
353 Preserving differences between mRNA measurements of distinct cell lines, this adjustment resulted in protein
354 measurements of one cell line to cluster together with mRNA measurements of the same cell line. Before,
355 protein and mRNA measurements clustered together in their respective dataset and were not similar to each
356 other. Here, we would like to stress that we only made use of MComBat in order to be able to perform mRNA-
357 guided missing value imputation. The increase in correlation between protein and mRNA data is only due to the
358 fact that MComBat removed systematic differences between them.

359 mRNA-guided missing value imputation

360 The increase in the overall correlation between transcriptomics and proteomics data after adjusting for
361 systematic differences between them (see Fig. EV2B and EV3C) enabled calculating protein abundance from
362 adjusted mRNA abundance. In order to do so, the CRC65 and CPTAC datasets had to be processed separately
363 from one another. For each protein sample and all of its cognate mRNA samples, we assembled all pairs of
364 protein and mRNA measurements corresponding to the same gene group and subsequently modelled protein
365 abundance as a function of adjusted mRNA abundance using a single linear model. For each cognate mRNA
366 sample, we then used the corresponding linear model to calculate the protein abundance of gene groups with
367 missing values in the respective protein sample. For each gene group with missing values in the respective
368 protein sample, we then used the mean across all calculated protein abundances based on cognate mRNA

369 samples to impute the giBAQ value of the corresponding gene group. Imputed giBAQ values smaller than zero
370 were treated as missing values (see Fig. EV2D).

371 Minimum-guided missing value imputation

372 Since we were not able to impute all missing values in the proteomics datasets using mRNA-guided missing
373 value imputation, we explored a number of secondary imputation methods, which could be used after mRNA-
374 guided missing value imputation (Figure EV2E). All of these methods shared the common assumption that most
375 missing values in proteomics datasets are “missing not at random” (MNAR), i.e. that the likelihood of
376 “encountering” a missing value increases as the abundance of the respective analyte approaches the detection
377 limit of the instrument (Lazar et al, 2016). We first confirmed that this assumption holds true for both the
378 CRC65 cell line and CPTAC patient datasets by calculating the frequency and fraction of missing values per
379 intensity bin for both datasets after MComBat adjustment but before mRNA-guided missing value imputation,
380 which was the basis for Figure EV2A. The first imputation method we tried out falls into the MinProb category
381 described by Lazar and colleagues and was the method described in Figure EV3 of the recent publication by
382 Tyanova et al (2016), using a down-shift of 1.8 and a width of 0.3. When we applied this method before mRNA-
383 guided missing value imputation, this resulted in strong bimodality of the intensity distribution for both datasets,
384 but even applying it after mRNA-guided missing value imputation did not circumvent this undesirable
385 behaviour completely. When adjusting the parameters of this perseus-type imputation such that the imputed
386 values became more and more (eventually fully) part of the overall distribution, we found an ever increasing and
387 eventually large number of cases in which the imputed values for specific proteins had much higher intensities
388 than the same but experimentally robustly measured protein intensities in a different sample, which is equally
389 undesirable. Since imputing missing values with the sample-wise (MinDet method described by Lazar et al.) or
390 protein-wise minimum introduced too much bias by reducing the variability of the measurements, we decided to
391 implement a new imputation method termed “minimum-guided missing value imputation” on the protein level.
392 This method ensured that no imputed value for a given protein was bigger than a measured value for the same
393 protein by sampling missing values with replacement from the distribution of measured values smaller than the
394 minimum for the respective protein. This was done across the entire dataset, which is why each dataset needed
395 to be normalised and handled separately. For each protein with one or more missing values across the dataset,
396 we first determined the minimum expression value. In a protein-wise fashion, missing values were subsequently
397 replaced with measured values, which were sampled from the truncated distribution of measured values as
398 described above. This resulted in a much more favourable behaviour with reduced bimodality of the intensity

399 distributions of each dataset if applied after mRNA-guided missing value imputation. Imputation on the protein
400 level however should not be favoured over imputation on the peptide level, since aggregating peptide-level
401 information to form protein expression values already represents a form of implicit missing value imputation, as
402 discussed by Lazar et al (2016). We therefore decided to transfer the concept of minimum-guided missing value
403 imputation to the peptide level, but only imputed missing values for peptides if the corresponding protein
404 abundance was missing as well. With a stable protein-group-to-HGNC-symbol mapping at hand (see above), we
405 made use of MaxQuant's relational database output and mapped each peptide and its intensity in the
406 evidence.txt output file to the protein group and by extension also to the gene group to the quantification of
407 which the respective peptide intensity contributed. Since razor peptides per definition map to more than one
408 protein group, we made sure to map their intensity only to those protein and gene groups for the quantification
409 of which MaxQuant actually used them. For each protein with one or more missing values across the dataset, we
410 then found the sample with the lowest expression of this protein and determined which peptides contributed to
411 this value. For each missing value in each of these reference peptides, we then searched in a window of 1.1 min
412 (centred at the recalibrated retention time of the respective reference peptide in the matching fractions) for
413 peptides with a lower intensity than the reference peptide. From the intensities of these low-intensity peptides,
414 we then sampled the missing values for the corresponding reference peptide with replacement. All imputed
415 peptide intensities corresponding to a given protein with missing values were subsequently summed up in a
416 sample-wise fashion to form protein intensities and calculate giBAQ values as described above. By adjusting the
417 retention time window queried for low-intensity peptides, one can control how similar the imputed values will
418 be to one another, thereby reducing the variability across the dataset in the low-intensity range. This is useful
419 especially for clustering algorithms using correlation as a measure of distance between samples, since low-
420 intensity features will not contribute as much to the distance between samples as they would have if they were
421 imputed with a different method, giving more weight to reliably measured features. Here, we decided to use the
422 same window MaxQuant used for transferring identifications by "Match between runs", since the concept is
423 somewhat similar. Missing values in our Kinobeads dataset were imputed using minimum-guided missing value
424 imputation on the gene group level. It is worth noting that the clustering results were the same ($p < 0.0005$, two-
425 sided Fisher's Exact Test) and the fold-changes of significantly differentially expressed proteins were highly
426 correlated (Pearson's $R = 0.997$) when using the perseus-type or minimum-guided missing value imputation.

427 After this two-step imputation, 93.6/74.7% of all values in the CRC65/CPTAC dataset were measured at the
428 protein level, 2.7/10.9% were measured at the mRNA level and 3.7/14.3% were imputed using minimum-guided

429 missing value imputation on the peptide level, respectively. The Kinobeads dataset contained 5.3% missing
430 values before imputation as described above.

431 Combination of the CRC65 and CPTAC datasets

432 After missing value imputation, the CRC65 cell line and CPTAC patient datasets had to be combined into one
433 expression matrix in order to be able to identify integrated proteomic subtypes of CRC consisting of cell lines
434 and patients. However, the imputation we performed might influence downstream analyses, since some of the
435 proteomics samples had more missing values than others and because we did not have transcriptomics data for
436 all samples. Missing values in samples for which we did not have data on the mRNA level were imputed using
437 minimum-guided missing value imputation only, possibly generating outlier samples, which will distort any
438 unsupervised clustering. We therefore first removed outlier samples from both datasets separately using
439 arrayQualityMetrics, only keeping samples not marked as an outlier by any of the outlier detection methods
440 performed (“Distance between arrays”, “Boxplots” and “MA plots”). Afterwards, we merged the two expression
441 matrices into one joined matrix. Since we identified substantially more peptides per sample in the CRC65
442 dataset than in the CPTAC dataset (Fig. 2 and EV2C) and because MaxQuant calculates protein intensities as
443 the sum of peptide intensities, we expected strong systematic differences between the two datasets, which would
444 have an unfavourable impact on any clustering algorithm applied to the joined expression matrix. Therefore, we
445 adjusted systematic differences between the two datasets using ComBat (Johnson et al, 2007) with MSI status as
446 a covariate as described above.

447 Identification & characterisation of subtypes

448 Identification of subtypes

449 After combining the full proteome datasets into one expression matrix and following the removal of systematic
450 differences between them, we used consensus clustering (Monti et al, 2003) as implemented in the
451 ConsensusClusterPlus R package v1.24.0 to determine Full Proteome Subtypes (FPSs) and Kinobeads Subtypes
452 (KSs) of CRC. The method we used was in many aspects similar to the method used by Zhang et al (2014), but
453 made use of a different agglomeration method, while being more stringent during the identification of core
454 samples. The main difference however was that we used the entire expression matrix for the discovery of FPSs
455 and only restricted our Kinobeads data to kinases during the identification of KSs. We used hierarchical
456 clustering to assign samples to one of $k=4$ clusters, with “1-Pearson correlation” as the distance metric,

457 “ward.D2” as the agglomeration method and 1,000 resampling iterations using 80% of all gene groups and
458 samples, while varying k from two to eight clusters during exploratory data analysis. We decided to use four
459 clusters based on the four known Consensus Molecular Subtypes published by Guinney et al (2015). Clusters
460 with less than six members were removed, followed by the calculation of item-consensus values and silhouette
461 widths (Rousseeuw, 1987) for each sample using the calcICL function from the ConsensusClusterPlus R package
462 v1.24.0 and the silhouette function of the R package cluster v2.0.3, respectively. Being more stringent than
463 Zhang et al., we defined core samples as samples with a positive silhouette width, provided that they also show
464 the highest item-consensus with the cluster they were finally assigned to. Clusters only consisting of core
465 samples for which we did not have mRNA data were excluded in order to avoid artefacts due to missing value
466 imputation. The columns in the final heat maps were restricted to core samples, which were ordered according
467 to the final consensus tree. This method was also used to cluster the reduced mRNA expression matrix after the
468 exclusion of cell lines not in the CRC65 dataset and transcripts not in the CMS classifier by Guinney et al.
469 (2015), in order to produce Figure 1B.

470 Differential expression analysis

471 We used significance analysis of microarrays (SAM; Tusher et al, 2001) as implemented in the samr R package
472 v2.0 in order to discover gene groups differentially expressed between the different subtypes we identified,
473 restricting the analysis to core samples while using 100 permutations, the Wilcoxon test statistic and a target
474 FDR of 0.01 (1%). The rows in the final heat maps were ordered according to a similar metric, which was
475 intended to visualise which gene groups were higher or lower expressed in one subtype compared to the
476 remaining subtypes. First, we standardised the expression matrix in a row-wise fashion so that each gene group
477 had a mean of zero and a standard deviation of one. For each gene group, we then calculated all pairwise
478 combinations of two-sided Wilcoxon test statistics between each subtype and the remaining subtypes. Each gene
479 group was then assigned to the subtype with the lowest p-value. For each subtype, gene groups associated with
480 it were then ordered in increasing order of their summed expression across all samples in this subtype, starting
481 with the gene groups with higher expression in the respective subtype compared to the others, followed by gene
482 groups with lower expression in the respective subtype compared to the others. These blocks of gene groups
483 were then sorted according to the order of the subtypes they were associated with in the final consensus tree,
484 creating the checkerboard-like pattern seen in Figures 3A and Figure EV5G.

485 Functional annotation enrichment & clustering

486 We used MetaCore v6.26.68498 from Thomson Reuters to perform enrichment analyses of functional
487 annotations for the gene groups differentially expressed between the different FPSs. In order to be more
488 stringent, we set the background MetaCore uses to calculate its p-values for the enrichment of functional
489 annotations using the hypergeometric test to the list of network objects defined by the HGNC symbols detected
490 in at least two samples in both the CRC65 cell line and CPTAC patient datasets. This ensures that the
491 enrichment analysis does not suffer from acquisition bias. The standardized mean difference between subtype-
492 specific gene group expression and overall mean expression as computed using SAM served as the raw input for
493 the enrichment analysis, which was further filtered to remove values between the 2.5% and 97.5% quantiles of
494 these values across all permutations, in order to restrict the analysis to gene groups with large class contrasts.
495 This filtered matrix was uploaded to <https://portal.genego.com/>, followed by the calculation of enrichment
496 analyses with respect to “Pathway Maps”, “Process Networks”, “GO Processes”, “GO Molecular Functions”
497 and “GO Localizations” for up- and down-regulated gene groups, respectively. We exported the entire table of
498 annotations for each of these functional annotation enrichments and further processed them in R. For “Pathway
499 Maps” and “Process Networks”, we only imported annotations with a minimum $FDR \leq 0.05$ (5%) as provided by
500 MetaCore, merged the tables of annotations enriched in up- and down-regulated gene groups from any subtype
501 into a single table, excluded rodent-specific annotations and subsequently assembled bar charts of $-\log_{10}(FDR)$
502 values for up- and down-regulated gene groups into a single figure. For each of the different Gene Ontology
503 (GO) categories, we also only imported annotations with a minimum $FDR \leq 0.05$ (5%) as provided by MetaCore.
504 Subsequently, tables for the different GO categories were combined while keeping GO-terms enriched in up-
505 and down-regulated gene groups specific to a certain subtype separate from one another. This resulted in six
506 different tables; two for each subtype containing GO-terms associated with up- and down-regulated gene
507 groups, respectively. Afterwards, GO-term names were re-mapped to GO Accessions using a mapping
508 downloaded from AmiGO using GOOSE (<http://amigo.geneontology.org/goose>, 19 May 2016) and GO
509 Accessions with an FDR of more than 0.05 (5%) were removed from the corresponding table. Finally, these
510 tables were sequentially uploaded to REVIGO (<http://revigo.irb.hr/>) for summarisation and visualisation of
511 significant GO-terms using treemaps, with an allowed similarity of 0.7, interpreting associated numbers as p-
512 values, restricting the database to Homo sapiens and employing the SimRel semantic similarity measure (01
513 July 2016; Schlicker et al, 2006).

514 Reanalysis of public dose-response datasets

515 Pre-processing

516 The same cell lines and drugs have inconsistent names in the different drug sensitivity datasets, which is why
517 the first pre-processing step always involved the mapping of all cell line and drug names to a set of consistent
518 sample and drug names, respectively (see also Tables EV6A and 6B). For reasons of transparency, we also
519 annotated each drug with the corresponding dataset it was derived from. Concordant cell lines were treated as
520 replicates during model fitting, which was kept as consistent as possible between the different datasets. In
521 essence, we always modelled a relative response measure R as a function of the final drug concentration on the
522 assay plate as described below.

523 GDSC

524 We obtained raw dose-response data (`gdsc_drug_sensitivity_raw_data_w5.csv`), annotation information
525 (`gdsc_tissue_output_w5.csv`) and screening concentrations (`gdsc_compounds_conc_w5.csv`) from
526 <ftp://ftp.sanger.ac.uk/pub4/cancerrxgene/releases/release-5.0/> (04 September 2015) prior to the publication by
527 Iorio et al (2016), which is why all our analyses were based on release v5.0. First, we annotated the raw data
528 with the corresponding drug names and “fold-dilution” information, as well as data on the maximum screening
529 concentration used for each drug. This allowed us to calculate the actual concentrations at which each
530 compound was used in the drug screen. We then calculated the mean background signal for each compound and
531 cell line (blank wells not containing cells) and subtracted it from each corresponding raw signal, setting negative
532 values to zero afterwards. Next, we determined the mean control signal for each compound and cell line and
533 afterwards divided each corresponding background-corrected raw signal by its respective control signal,
534 generating normalised and background-corrected viability data. Finally, for each compound, compound
535 concentration and cell line, we then calculated relative response measures R as the mean of these values and
536 restricted the dataset to cell lines from the large intestine, as well as C32 and SW626.

537 CCLE

538 We downloaded Supplementary Table 11 from the original publication by Barretina et al (2012) and restricted
539 the dataset first to cell lines from the large intestine, as well as C32 and Colo 741. We then restricted the dataset
540 further to only include drugs for which the authors were able to fit at least two sigmoid dose-response models
541 among our selection of cell lines as indicated in the `FitType` column, removing dose-response data for both
542 “Nutlin-3” and “PHA-665752”. Subsequently, we converted raw activity values A (“Activity Data (median)”) to

543 back to relative response measures $R=T/U=1+A/100$ (T represents the response for the compound-treated well,
544 U represents the median response of the untreated wells across the plate) by reverting the corresponding
545 equation in the addendum to the original publication, in order to make the CCLE dataset more comparable to the
546 other drug sensitivity datasets.

547 CTRP

548 We obtained the expanded CTRP v2.0 dataset from the FTP server of the CTD² data portal at
549 ftp://caftpd.nci.nih.gov/pub/dcc_ctd2/Broad/CTRPv2.0_2015_ctd2_ExpandedDataset/ (23 December 2015).
550 The results published here are therefore partially based upon data generated by the Cancer Target Discovery and
551 Development (CTD2) Network (<https://ctd2.nci.nih.gov/dataPortal/>) established by the National Cancer
552 Institute's Office of Cancer Genomics. First, we annotated the raw data with the corresponding compound and
553 cell line names, as well as information on the tissue of origin for the cell lines used in the screen. Based on these
554 annotations, the dataset was restricted to cell lines from the large intestine, as well as C32 and Colo 741.
555 Viability data annotated with the same cell line name, compound concentration and compound name were
556 treated as replicates during model fitting. The weighted percent-viability with error propagation ("cpd_avg_pv")
557 was used as a relative response measure R.

558 Cetuximab

559 We downloaded (03 January 2016) Supplementary Data 1 from the original publication by Medico et al (2015)
560 and converted the percentage of growth inhibition P to a relative response measure $R=1-P/100$ in order to make
561 the data comparable to the other drug sensitivity datasets.

562 Dose-response models & parameter extraction

563 Using the drc package v2.5-12 in R, we fitted the classical symmetric four-parameter log-logistic model to each
564 drug in each dataset:

$$f(x, (b, c, d, e)) = c + \frac{d - c}{1 + \exp(b(\log(x) - \log(e)))}$$

565 The four parameters c="Lower Limit" (i.e. lower limit of the relative response as the drug concentration
566 approaches infinity), d="Upper Limit" (i.e. upper limit of the relative response as the drug concentration
567 approaches zero), b="Slope" (Hill slope, i.e. slope at the inflection point of the dose-response curve), and
568 e="ED50" (i.e. the dose required to reduce the relative response to $c+(d-c)/2$) were estimated from the data and
569 subsequently extracted from the fitted model. We also modified the computeAUC function from the drexplorer

570 package v1.1.2 to accept the output of the `drm` function from the `drc` R package as an input and used it to extract
571 the standardised area under the dose-response curve (AUC) for each of the fitted models across the tested dose-
572 range. Here, the AUC across the tested dose-range was defined as the area under the dose-response curve
573 between zero and one, divided by the area under $y=1$ from the lowest to the highest concentration tested. It is
574 worth noting that the relative response measures described above varied between approximately zero and one,
575 enabling unified model fitting and parameter extraction.

576 Modelling of Drug Sensitivity

577 In order to find proteins and kinases associated with drug sensitivity or resistance, we used elastic net regression
578 (Zou & Hastie, 2005) to model drug sensitivity as a function of protein and kinase profile, respectively.
579 Together with a bootstrapping approach, this regularised multivariate linear regression method was previously
580 successfully applied to identify genomic, transcriptomic and proteomic markers of drug sensitivity and
581 resistance (Barretina et al, 2012; Garnett et al, 2012; Gholami et al, 2013; Iorio et al, 2016). The advantages of
582 the elastic net are discussed in detail in the Supplementary Methods associated with the publication by Barretina
583 et al (2012), as well as in the original publication by Zou and Hastie (2005). Briefly, the algorithm is especially
584 useful for “large p , small N ” ($p \gg N$) regression problems involving for example data generated using current
585 omics technologies of various kind, which measure many, possibly highly correlated analytes (p) across a small
586 number of samples (N). In settings like these with varying degrees of multicollinearity, the ordinary least
587 squares (OLS) estimates of the regression coefficients and the resulting models frequently perform poorly both
588 in terms of prediction accuracy on new data, as well as regarding the interpretability of the model itself (Zou &
589 Hastie, 2005). The elastic net aims at improving both prediction performance and model interpretability by
590 consolidating L2-regularised Ridge regression (Hoerl & Kennard, 1970) and L1-regularised Lasso regression
591 (Tibshirani, 1996), in order to combine the advantages of the different penalty terms they use to constrain the
592 coefficient vector. On its own, Ridge regression often has a higher predictive accuracy than OLS in the above-
593 mentioned situations by striking a balance between bias and variance of predicted values via its bound on the
594 L2-norm of the coefficient vector, thereby aiming at improving the overall prediction accuracy (Tibshirani,
595 1996). However, it always keeps all predictors in the model, which is not useful if the goal is to increase its
596 interpretability compared to OLS. Lasso regression on the other hand does automatic variable selection via its
597 bound on the L1-norm of the coefficient vector, thereby shrinking some coefficients while forcing others to be
598 exactly zero, which results in a more parsimonious model consisting of a reduced subset of good predictors,

599 which is easier to interpret (Tibshirani, 1996). However, this comes at the cost of only selecting at most N
600 variables in settings where $p \gg N$, as well as only selecting one variable from a group of correlated variables to
601 be included in the model while excluding the others in settings with multicollinearity (Zou & Hastie, 2005). The
602 authors also observed that Ridge regression performs far better than Lasso regression with respect to predictions
603 on new data in settings with multicollinearity where $N > p$. By combining these two penalty terms, the elastic net
604 now promotes parsimony of models through the L1-penalty, while at the same time encouraging a grouping
605 effect through the L2-penalty (Barretina et al, 2012), meaning that highly collinear predictors are usually either
606 in the model or dropped from it together. Double shrinkage of coefficients is prevented by the introduction of a
607 scaling factor (Zou & Hastie, 2005) and the higher level parameter α with $0 \leq \alpha \leq 1$ is used to control the
608 balance between the L2-penalty ($\alpha=0$) and the L1-penalty ($\alpha=1$), while λ controls the degree of regularisation
609 (Friedman et al, 2010).

610 If $y_i \in \mathbb{R}^N$ are drug responses expressed as AUC for $i=1, \dots, N$ cell lines and $x_i \in \mathbb{R}^N$ are standardised
611 measurements for $i=1, \dots, N$ cell lines across $j=1, \dots, p$ proteins or kinases with zero mean and unit variance
612 forming N observation pairs (x_i, y_i) , then the cost function we minimised using glmnet v2.0-5 (Friedman et al,
613 2010) is given by

$$\min_{(\beta_0, \beta) \in \mathbb{R}^{p+1}} \frac{1}{2N} \sum_{i=1}^N (y_i - \beta_0 - x_i^T \beta)^2 + \lambda \left[(1 - \alpha) \frac{1}{2} \|\beta\|_{l_2}^2 + \alpha \|\beta\|_{l_1} \right]$$

614 We used leave-one-out cross-validation over a grid of 100 equally-spaced values of α with $0 \leq \alpha \leq 1$ by 100
615 values of λ automatically selected by the algorithm for each α to find the combination of α and λ resulting in the
616 lowest mean-squared error between the fitted mean and the response. We then turned to a bootstrapping
617 approach similar to the ones used before in order to be able to sort the different proteins or kinases according to
618 their importance in predicting the response to a given drug (Barretina et al, 2012; Garnett et al, 2012; Gholami et
619 al, 2013; Iorio et al, 2016). Briefly, we generated 1000 bootstrap datasets by sampling the above-mentioned
620 observation pairs with replacement and solved the elastic net optimisation problem for each of these datasets
621 using the optimal values of α and λ . This resulted in a $1000 \times (p + 1)$ matrix of regression coefficients
622 capturing the solutions for the bootstrap datasets as rows and the regression coefficients as columns. The
623 intercept column $\beta_{int}^{BS} \in \mathbb{R}^{1000,1}$ was stored separately from the remainder of the matrix, which left us with a
624 $1000 \times p$ matrix of coefficients $\beta_{coef}^{BS} \in \mathbb{R}^{1000,p}$ for $j=1, \dots, p$ proteins or kinases across 1000 bootstrap datasets.
625 Each column of β_{coef}^{BS} was subsequently multiplied by the standard deviation of the corresponding protein or
626 kinase across all x_i , as well as by its corresponding weight ω described in Equation 11 in the publication by

627 Moghaddas Gholami et al (2013) if $\omega \in \{-1,1\}$ or by zero if $\omega \notin \{-1,1\}$. This resulted in a $1000 \times p$ matrix of
628 effect sizes β_{eff}^{BS} analogous to Equation 12 in the publication by Moghaddas Gholami et al. It is worth stressing
629 that columns of effect sizes corresponding to proteins or kinases whose coefficients did not have a consistent
630 sign across all 1000 bootstrap solutions were set to zero through the multiplication with the weight ω . For each
631 column in β_{int}^{BS} , β_{coef}^{BS} and β_{eff}^{BS} , we then calculated the mean across all bootstrap solutions to form β_{int}^{mean} ,
632 β_{coef}^{mean} and β_{eff}^{mean} , respectively, and also counted the number of times the absolute effect size of a given protein
633 or kinase was bigger than zero to form $F_{eff} \in \mathbb{N}^p$. For each protein or kinase $j=1, \dots, p$, we then determined
634 whether (a) $F_{eff,j}$ was bigger than or equal to the 9th decile of F_{eff} and also checked whether $\beta_{eff,j}^{mean}$ was either
635 (b₁) bigger than or equal to the 19th ventile of β_{eff}^{mean} or (b₂) smaller than or equal to the 1st ventile of β_{eff}^{mean} . We
636 only considered proteins or kinases for which both condition (a) and one of condition (b₁) or (b₂) were true to be
637 robust predictive markers of drug response, while all proteins or kinases $j=1, \dots, p$ with $|\beta_{eff,j}^{mean}| > 0$ were
638 designated as being associated with drug response. In effect-size heat maps, we only displayed the top 5 and
639 bottom 5 robust predictive markers of drug response with respect to β_{eff}^{mean} , while restricting the columns to
640 only the top 10 and bottom 10 cell lines with respect to their drug response as measured by AUC. For each drug,
641 we always used all N available observation pairs (x_i, y_i) to train the model with three different input matrices
642 each:

- 643 1) The combined full proteome expression matrix after missing value imputation
- 644 2) All proteins in the Kinobeads expression matrix after missing value imputation
- 645 3) All kinases in the Kinobeads expression matrix after missing value imputation

646 For each drug, we then used the corresponding β_{int}^{mean} and β_{coef}^{mean} from scenario 1) to predict the drug response
647 \hat{y}_i of each cell line and patient in the combined dataset based on the corresponding full proteome expression
648 profile x_i using the following formula:

$$\hat{y}_i = \beta_{int}^{mean} + \sum_{j=1}^p \beta_{coef,j}^{mean} \times x_{i,j}$$

649 These predictions were subsequently standardised to have zero mean and unit variance and can be found in
650 Table EV3A. We also predicted \hat{y}_i using β_{int}^{mean} and β_{coef}^{mean} from scenario 3) and subsequently standardised
651 these predictions to have zero mean and unit variance in order to generate kinase-centric drug sensitivity
652 hypotheses. Cell lines with a standardised AUC bigger than zero were designated as “resistant”, while cell lines
653 with a standardised AUC smaller than zero were designated as “sensitive”. The type of figure used to visualise

654 the result of an elastic net regression is termed “effect-size heat map”. In effect-size heat maps, up to five of the
655 most robust predictive markers of drug resistance and sensitivity with respect to β_{eff}^{mean} were displayed,
656 respectively, while restricting the columns to at most the top 10 and bottom 10 cell lines with respect to their
657 drug response as measured by AUC in order to increase visual clarity. In such a figure, these top and bottom ten
658 cell lines are shown as columns ordered from left to right in increasing order of drug resistance, while the
659 predictors (proteins or kinases) with respect to absolute effect-size are shown as rows. In the bar plot to the left,
660 predictors associated with drug resistance are visualised by yellow bars, while predictors associated with drug
661 sensitivity are shown as dark blue bars in increasing or decreasing order of absolute effect-size from top to
662 bottom, respectively. The heat map itself shows standardized expression values (z-scores) of the predictors with
663 blue=“low expression” and red=“high expression”. Below the heat map, a heat strip visualises the AUC of the
664 drug in question from dark-blue=“low AUC” to yellow=“high AUC” (see Figure 1A for an example of an
665 effect-size heat map).

666 Cox proportional hazards regression

667 In order to assess whether or not the cytoplasmic/membranous expression level of MERTK is a prognostic
668 biomarker in CRC, we modelled the different outcome variables of the QUASAR 2 trial as a function of
669 MERTK expression using univariate and multivariate Cox proportional hazards regression (COXPH), restricted
670 to the first 5 years of follow-up. First, univariate COXPH was carried out using the `coxph` function of the R
671 package `survival` v2.38-3 with default parameters and “Gender”, “Age”, “BMI”, “Treatment” (capecitabine ±
672 bevacizumab), “Location” (right colon or left colon), “T” (from TNM staging system; T2/T3 or T4), “N” (from
673 TNM staging system; N0, N1 or N2), “MSI” (microsatellite instability), “CIN” (chromosomal instability) and
674 “MERTK” (cytoplasmic/membranous MERTK expression; ≤5% or >5%) as covariates. The threshold of 5%
675 was selected in order to balance the groups while avoiding false positives due to weak staining. Significant
676 ($p < 0.05$) covariates from univariate COXPH were then used as explanatory variables in multivariate COXPH.
677 We made use of the `stepAIC` function from the R package `MASS` v7.3-45 to perform stepwise backward model
678 selection starting with the full model and removing explanatory variables in order to minimise Akaike’s
679 Information Criterion (AIC). For each outcome variable, final models were then built using significant ($p < 0.05$)
680 covariates from the stepwise selected models. These final models all fulfilled the proportional hazards
681 assumption. Results from the COXPH analyses can be found in Table EV5.

682 Annotation information

683 We compiled annotation information for the CRC65 and CPTAC datasets, as well as target annotations for
684 drugs in the four drug sensitivity datasets (GDSC, CCLE, CTRP and Cetuximab) from several sources
685 (Supplementary Table 11 & GSE36133, Barretina et al., 2012; Supplementary Table 1, Bracht et al., 2010;
686 Supplementary Table 3, De Sousa et al., 2013; COSMIC, Forbes et al., 2015; GDSC release 5.0 & E-MTAB-
687 783, Garnett et al., 2012; cms_labels_public_all.txt, Guinney et al., 2015; GDSC release 5.0, Iorio et al., 2016;
688 Figure 1 & E-MTAB-2706, Klijn et al., 2015; Table 1, Ku et al., 1999; Table S1, Liu and Zhang, 2016;
689 GSE28567, Loboda et al., 2011; Supplementary Data 1 & 2 & GSE59857, Medico et al., 2015; Supplementary
690 Table 1 & 6, Mouradov et al., 2014; CTRP v2.0 expanded dataset, Rees et al., 2016; Supplementary Table 2,
691 Sadanandam et al., 2013; Supplementary Table 1 & platform codes "IlluminaGA_RNASeq" &
692 "AgilentG4502A_07_3", TCGA Network, 2012; GSE8332, Wagner et al., 2007; Table 1, Wheeler et al., 1999;
693 GSE24795, Wilding et al., 2010; Supplementary Table 1, Zhang et al., 2014) and summarised them in Table
694 EV6. Entries for cell lines in the columns "Sadanandam subtype", "Marisa subtype", "De Sousa E Melo
695 subtype" and "Roepman subtype" were based on the publication by Medico and colleagues and were not
696 restricted by an FDR filter. The "Goblet-like", "Goblet" and "Enterocyte" subtype of the classifier by
697 Sadanandam and colleagues were merged into the "Enterocyte/Goblet-like" subtype, while "MSI-L" was
698 recoded to "MSI-" and "MSI-H" was recoded to "MSI+" according to the publication by Liu and Zhang. Some
699 of the other annotations were also recoded to fit the unified nomenclature. For drug targets, we applied the
700 following heuristic: After mapping compounds to a consensus set of names, the annotation from the CTRP
701 dataset was used as the first source of target information. Drugs for which this dataset lacked annotation were
702 assigned targets based on the annotation from the GDSC dataset. Drugs the target information of which was not
703 included in the GDSC dataset were then assigned targets based on the annotation from the CCLE dataset. Drugs
704 still left without target information were finally assigned labels based on the information contained in the
705 "target_or_activity_of_compound" column of the "v20.meta.per_compound.txt" file from the expanded CTRP
706 v2.0 dataset. Amendments to these mappings in order to reduce redundancy are documented in Table EV6B.

707

708 Appendix Supplementary References

709

710 Barretina J, Caponigro G, Stransky N, Venkatesan K, Margolin AA, Kim S, Wilson CJ, Lehar J, Kryukov GV,
711 Sonkin D, Reddy A, Liu M, Murray L, Berger MF, Monahan JE, Morais P, Meltzer J, Korejwa A, Jane-
712 Valbuena J, Mapa FA et al (2012) The Cancer Cell Line Encyclopedia enables predictive modelling of
713 anticancer drug sensitivity. *Nature* **483**: 603-607

714

715 Bracht K, Nicholls AM, Liu Y, Bodmer WF (2010) 5-Fluorouracil response in a large panel of colorectal cancer
716 cell lines is associated with mismatch repair deficiency. *Br J Cancer* **103**: 340-346

717

718 De Sousa EMF, Wang X, Jansen M, Fessler E, Trinh A, de Rooij LP, de Jong JH, de Boer OJ, van Leersum R,
719 Bijlsma MF, Rodermond H, van der Heijden M, van Noesel CJ, Tuynman JB, Dekker E, Markowitz F, Medema
720 JP, Vermeulen L (2013) Poor-prognosis colon cancer is defined by a molecularly distinct subtype and develops
721 from serrated precursor lesions. *Nat Med* **19**: 614-618

722

723 Edgar R, Domrachev M, Lash AE (2002) Gene Expression Omnibus: NCBI gene expression and hybridization
724 array data repository. *Nucleic Acids Res* **30**: 207-210

725

726 Edwards D (2003) Non-linear normalization and background correction in one-channel cDNA microarray
727 studies. *Bioinformatics* **19**: 825-833

728

729 Forbes SA, Beare D, Gunasekaran P, Leung K, Bindal N, Boutselakis H, Ding M, Bamford S, Cole C, Ward S,
730 Kok CY, Jia M, De T, Teague JW, Stratton MR, McDermott U, Campbell PJ (2015) COSMIC: exploring the
731 world's knowledge of somatic mutations in human cancer. *Nucleic Acids Res* **43**: D805-811

732

733 Friedman J, Hastie T, Tibshirani R (2010) Regularization Paths for Generalized Linear Models via Coordinate
734 Descent. *J Stat Softw* **33**: 1-22

735

736 Garnett MJ, Edelman EJ, Heidorn SJ, Greenman CD, Dastur A, Lau KW, Greninger P, Thompson IR, Luo X,
737 Soares J, Liu Q, Iorio F, Surdez D, Chen L, Milano RJ, Bignell GR, Tam AT, Davies H, Stevenson JA,
738 Barthorpe S et al (2012) Systematic identification of genomic markers of drug sensitivity in cancer cells. *Nature*
739 **483**: 570-575

740

741 Gholami AM, Hahne H, Wu Z, Auer FJ, Meng C, Wilhelm M, Kuster B (2013) Global proteome analysis of the
742 NCI-60 cell line panel. *Cell Rep* **4**: 609-620

743

744 Guinney J, Dienstmann R, Wang X, de Reynies A, Schlicker A, Soneson C, Marisa L, Roepman P,
745 Nyamundanda G, Angelino P, Bot BM, Morris JS, Simon IM, Gerster S, Fessler E, De Sousa EMF, Missiaglia
746 E, Ramay H, Barras D, Homicsko K et al (2015) The consensus molecular subtypes of colorectal cancer. *Nat*
747 *Med* **21**: 1350-1356

748

749 Hoerl AE, Kennard RW (1970) Ridge Regression - Biased Estimation for Nonorthogonal Problems.
750 *Technometrics* **12**: 55-&

751

752 Iorio F, Knijnenburg TA, Vis DJ, Bignell GR, Menden MP, Schubert M, Aben N, Goncalves E, Barthorpe S,
753 Lightfoot H, Cokelaer T, Greninger P, van Dyk E, Chang H, de Silva H, Heyn H, Deng X, Egan RK, Liu Q,
754 Mironenko T et al (2016) A Landscape of Pharmacogenomic Interactions in Cancer. *Cell* **166**: 740-754

755

756 Johnson WE, Li C, Rabinovic A (2007) Adjusting batch effects in microarray expression data using empirical
757 Bayes methods. *Biostatistics* **8**: 118-127

758

759 Klijn C, Durinck S, Stawiski EW, Haverty PM, Jiang Z, Liu H, Degenhardt J, Mayba O, Gnad F, Liu J, Pau G,
760 Reeder J, Cao Y, Mukhyala K, Selvaraj SK, Yu M, Zynda GJ, Brauer MJ, Wu TD, Gentleman RC et al (2015)
761 A comprehensive transcriptional portrait of human cancer cell lines. *Nat Biotechnol* **33**: 306-312

762

763 Ku JL, Yoon KA, Kim DY, Park JG (1999) Mutations in hMSH6 alone are not sufficient to cause the
764 microsatellite instability in colorectal cancer cell lines. *European Journal of Cancer* **35**: 1724-1729

765

766 Lazar C, Gatto L, Ferro M, Bruley C, Burger T (2016) Accounting for the Multiple Natures of Missing Values
767 in Label-Free Quantitative Proteomics Data Sets to Compare Imputation Strategies. *J Proteome Res* **15**: 1116-
768 1125

769

770 Liu Q, Zhang B (2016) Integrative Omics Analysis Reveals Post-Transcriptionally Enhanced Protective Host
771 Response in Colorectal Cancers with Microsatellite Instability. *J Proteome Res* **15**: 766-776

772

773 Loboda A, Nebozhyn MV, Watters JW, Buser CA, Shaw PM, Huang PS, Van't Veer L, Tollenaar RA, Jackson
774 DB, Agrawal D, Dai H, Yeatman TJ (2011) EMT is the dominant program in human colon cancer. *BMC Med
775 Genomics* **4**: 9

776

777 Medico E, Russo M, Picco G, Cancelliere C, Valtorta E, Corti G, Buscarino M, Isella C, Lamba S, Martinoglio
778 B, Veronese S, Siena S, Sartore-Bianchi A, Beccuti M, Mottolese M, Linnebacher M, Cordero F, Di
779 Nicolantonio F, Bardelli A (2015) The molecular landscape of colorectal cancer cell lines unveils clinically
780 actionable kinase targets. *Nat Commun* **6**: 7002

781

782 Monti S, Tamayo P, Mesirov J, Golub T (2003) Consensus clustering: A resampling-based method for class
783 discovery and visualization of gene expression microarray data. *Mach Learn* **52**: 91-118

784

785 Mouradov D, Sloggett C, Jorissen RN, Love CG, Li S, Burgess AW, Arango D, Strausberg RL, Buchanan D,
786 Wormald S, O'Connor L, Wilding JL, Bicknell D, Tomlinson IP, Bodmer WF, Mariadason JM, Sieber OM
787 (2014) Colorectal cancer cell lines are representative models of the main molecular subtypes of primary cancer.
788 *Cancer Res* **74**: 3238-3247

789

790 Rees MG, Seashore-Ludlow B, Cheah JH, Adams DJ, Price EV, Gill S, Javaid S, Coletti ME, Jones VL,
791 Bodycombe NE, Soule CK, Alexander B, Li A, Montgomery P, Kotz JD, Hon CS, Munoz B, Liefeld T, Dancik
792 V, Haber DA et al (2016) Correlating chemical sensitivity and basal gene expression reveals mechanism of
793 action. *Nat Chem Biol* **12**: 109-116

794

795 Ritchie ME, Silver J, Oshlack A, Holmes M, Diyagama D, Holloway A, Smyth GK (2007) A comparison of
796 background correction methods for two-colour microarrays. *Bioinformatics* **23**: 2700-2707

797

798 Rousseeuw PJ (1987) Silhouettes - a Graphical Aid to the Interpretation and Validation of Cluster-Analysis.
799 *Journal of Computational and Applied Mathematics* **20**: 53-65

800

801 Sadanandam A, Lyssiotis CA, Homicsko K, Collisson EA, Gibb WJ, Wullschleger S, Ostos LC, Lannon WA,
802 Grotzinger C, Del Rio M, Lhermitte B, Olshen AB, Wiedenmann B, Cantley LC, Gray JW, Hanahan D (2013)
803 A colorectal cancer classification system that associates cellular phenotype and responses to therapy. *Nat Med*
804 **19**: 619-625

805

806 Schlicker A, Domingues FS, Rahnenfuhrer J, Lengauer T (2006) A new measure for functional similarity of
807 gene products based on Gene Ontology. *BMC Bioinformatics* **7**: 302

808

809 Schwanhaussner B, Busse D, Li N, Dittmar G, Schuchhardt J, Wolf J, Chen W, Selbach M (2011) Global
810 quantification of mammalian gene expression control. *Nature* **473**: 337-342

811

812 Shi W, Oshlack A, Smyth GK (2010) Optimizing the noise versus bias trade-off for Illumina whole genome
813 expression BeadChips. *Nucleic Acids Res* **38**: e204

814

815 Silver JD, Ritchie ME, Smyth GK (2009) Microarray background correction: maximum likelihood estimation
816 for the normal-exponential convolution. *Biostatistics* **10**: 352-363

817

818 Stein CK, Qu P, Epstein J, Buros A, Rosenthal A, Crowley J, Morgan G, Barlogie B (2015) Removing batch
819 effects from purified plasma cell gene expression microarrays with modified ComBat. *BMC Bioinformatics* **16**:
820 63

821

822 TCGA Network (2011) Integrated genomic analyses of ovarian carcinoma. *Nature* **474**: 609-615

823

824 TCGA Network (2012) Comprehensive molecular characterization of human colon and rectal cancer. *Nature*
825 **487**: 330-337

826

827 Tibshirani R (1996) Regression shrinkage and selection via the Lasso. *J Roy Stat Soc B Met* **58**: 267-288

828

829 Tusher VG, Tibshirani R, Chu G (2001) Significance analysis of microarrays applied to the ionizing radiation
830 response. *Proc Natl Acad Sci U S A* **98**: 5116-5121

831

832 Tyanova S, Temu T, Sinitcyn P, Carlson A, Hein MY, Geiger T, Mann M, Cox J (2016) The Perseus
833 computational platform for comprehensive analysis of (prote)omics data. *Nat Methods* **13**: 731-740

834

835 Wagner KW, Punnoose EA, Januario T, Lawrence DA, Pitti RM, Lancaster K, Lee D, von Goetz M, Yee SF,
836 Totpal K, Huw L, Katta V, Cavet G, Hymowitz SG, Amler L, Ashkenazi A (2007) Death-receptor O-
837 glycosylation controls tumor-cell sensitivity to the proapoptotic ligand Apo2L/TRAIL. *Nat Med* **13**: 1070-1077

838

839 Wheeler JM, Beck NE, Kim HC, Tomlinson IP, Mortensen NJ, Bodmer WF (1999) Mechanisms of inactivation
840 of mismatch repair genes in human colorectal cancer cell lines: the predominant role of hMLH1. *PNAS* **96**:
841 10296-10301

842

843 Wilding JL, McGowan S, Liu Y, Bodmer WF (2010) Replication error deficient and proficient colorectal cancer
844 gene expression differences caused by 3'UTR polyT sequence deletions. *Proc Natl Acad Sci U S A* **107**: 21058-
845 21063

846

847 Wilhelm M, Schlegl J, Hahne H, Gholami AM, Lieberenz M, Savitski MM, Ziegler E, Butzmann L, Gessulat S,
848 Marx H, Mathieson T, Lemeer S, Schnatbaum K, Reimer U, Wenschuh H, Mollenhauer M, Slotta-Huspenina J,
849 Boese JH, Bantscheff M, Gerstmair A et al (2014) Mass-spectrometry-based draft of the human proteome.
850 *Nature* **509**: 582-587

851

852 Zhang B, Wang J, Wang X, Zhu J, Liu Q, Shi Z, Chambers MC, Zimmerman LJ, Shaddox KF, Kim S, Davies
853 SR, Wang S, Wang P, Kinsinger CR, Rivers RC, Rodriguez H, Townsend RR, Ellis MJ, Carr SA, Tabb DL et al
854 (2014) Proteogenomic characterization of human colon and rectal cancer. *Nature* **513**: 382-387

855

856 Zou H, Hastie T (2005) Regularization and variable selection via the elastic net. *Journal of the Royal Statistical*
857 *Society: Series B (Statistical Methodology)* **67**: 301-320

858

859