

1 **Pathway-Based Subnetworks Enable Cross-Disease Biomarker Discovery**

2 Syed Haider^{1,2,14}, Cindy Q .Yao^{1,3,4}, Vicky S. Sabine³, Michal Grzadkowski¹,
3 Vincent Stimper¹, Maud H.W. Starmans^{1,5}, Jianxin Wang¹, Francis Nguyen^{1,4},
4 Nathalie C. Moon¹, Xihui Lin¹, Camilla Drake³, Cheryl A. Crozier³, Cassandra L.
5 Brookes⁶, Cornelis J.H. van de Velde⁷, Annette Hasenburg⁸, Dirk G. Kieback⁹,
6 Christos J. Markopoulos¹⁰, Luc Y. Dirix¹¹, Caroline Seynaeve¹², Daniel W. Rea⁶,
7 Arek Kasprzyk¹, Philippe Lambin⁵, Pietro Lio², John M.S. Bartlett^{3,14}, Paul C.
8 Boutros^{1,4,13,14}

9 ¹ Informatics and Biocomputing Program, Ontario Institute for Cancer Research,
10 Toronto, M5G 0A3, Canada

11 ² Computer Laboratory, University of Cambridge, Cambridge, CB3 0FD, United
12 Kingdom

13 ³ Diagnostic Development Program, Ontario Institute for Cancer Research,
14 Toronto, M5G 0A3, Canada

15 ⁴ Department of Medical Biophysics, University of Toronto, Toronto, Canada

16 ⁵ Department of Radiation Oncology (Maastr), GROW-School for Oncology and
17 Developmental Biology, Maastricht University Medical Center, Maastricht, The
18 Netherlands

19 ⁶ Cancer Research UK Clinical Trials Unit, University of Birmingham,
20 Birmingham, B15 2TT, United Kingdom

21 ⁷ Leiden University Medical Center, Leiden, The Netherlands

22 ⁸ University Hospital, Freiburg, Germany

23 ⁹ Klinikum Vest Medical Center, Marl, Germany

24 ¹⁰ Athens University Medical School, Athens, Greece

25 ¹¹ St. Augustinus Hospital, Antwerp, Belgium

26 ¹² Erasmus Medical Center-Daniel den Hoed, Rotterdam, The Netherlands

27 ¹³ Department of Pharmacology and Toxicology, University of Toronto, Toronto,
28 M5S 1A8, Canada

29 ¹⁴ Corresponding authors

30

31

32 **Table of Contents**

33 **Supplementary Methods.....4**

34 **1. Univariate analyses reveal outliers and duplicate profiles in breast cancer.....4**

35 **2. SIMMS’ comparison with other machine learning algorithms and genesets/pathway**

36 **scoring methods5**

37 **3. SIMMS’ comparison with breast, colon, NSCLC and ovarian cancer prognostic biomarkers**

38 **.....7**

39 **4. SIMMS-derived PIK3CA signaling residual risk predictor in early breast cancer 8**

40 *4.1 TEAM cohort power calculations..... 8*

41 *4.2 mRNA abundance data processing..... 8*

42 *4.3 Survival modelling 9*

43 *4.4 IHC4 model..... 11*

44 *4.5 Recurrence probabilities..... 11*

45 *4.6 Performance Assessment..... 11*

46 *4.7 Prognostic assessment of SIMMS PI3K modules signature in CT+/- groups 12*

47 **5. Modelling multi-modal datatypes using SIMMS 12**

48 **6. SIMMS R package..... 13**

49 **Supplementary Figure Legends..... 15**

50 **Supplementary Figure 1..... 15**

51 **Supplementary Figure 2..... 16**

52 **Supplementary Figure 3..... 18**

53 **Supplementary Figure 4..... 20**

54 **Supplementary Figure 5..... 21**

55 **Supplementary Figure 6..... 22**

56 **Supplementary Figure 7..... 24**

57 **Supplementary Figure 8..... 24**

58 **Supplementary Figure 9..... 25**

59 **Supplementary Figure 10..... 26**

60 **Supplementary Figure 11..... 27**

61	Supplementary Figure 12.....	28
62	Supplementary Figure 13.....	29
63	Supplementary Figure 14.....	30
64	Supplementary Figure 15.....	31
65	Supplementary Figure 16.....	32
66	Supplementary Figure 17.....	33
67	Supplementary Figure 18.....	34
68	Supplementary Figure 19.....	35
69	Supplementary Figure 20.....	36
70	Supplementary Figure 21.....	37
71	Supplementary Figure 22.....	38
72	Supplementary Figure 23.....	39
73	Supplementary Figure 24.....	40
74	Supplementary Figure 25.....	42
75	Supplementary References	44
76		
77		

78 **Supplementary Methods**

79 **1. Univariate analyses reveal outliers and duplicate profiles in breast cancer**

80 We collated 14 mRNA abundance breast cancer datasets (**Supplementary**
 81 **Table 2**). Since these datasets originate from different studies and array
 82 platforms, comprehensive univariate analyses were performed to identify outlier
 83 datasets and to find patients duplicated across datasets. First, each dataset was
 84 pre-processed independently (**Methods section: mRNA abundance and**
 85 **survival data pre-processing**). Next, genes across all the datasets were
 86 evaluated for their prognostic ability using a univariate Cox proportional hazards
 87 model followed by the Wald-test. All the genes were subsequently ranked by the
 88 Wald-test P value within each study. The top genes across all studies were
 89 compared on multiple criterion as detailed below:

90 1 - Rank Product

91 The Rank Product¹ of each gene was computed as:

$$RP_g = \sum_{i=1}^k \log(r_{gi})^{\frac{1}{k}} \quad (1)$$

92 Here k represents the number of studies which had the mRNA abundance
 93 measure available for gene g . r_i is the rank of gene g in study i . The overall
 94 ranking table was used as a benchmark to identify datasets in which a given
 95 gene was ranked farthest when its rank product was compared to studywise
 96 ranks. The farthest dataset count was computed for the overall top ranked (100,
 97 200, 300, ..., 1000, 2000) genes (**Supplementary Figure 3a-e**).

98 2 - Percentile ranks

99 The P value (Wald-test) based ranking was transformed into percentile ranks
 100 within each study. These ranks were used as a measure of gene's position with
 101 reference to the benchmark rank derived in the step 1 to evaluate deviation of
 102 genes' ranks for each study (**Supplementary Figure 3f-i**).

103 3 - Intra- and inter-study correlation

104 The mRNA abundance profiles of common genes across all studies were
105 extracted and patient wise Spearman rank correlation coefficient was estimated.
106 The correlation coefficient was used to further analyze intra- and inter-study
107 correlation in order to identify any outlier studies (**Supplementary Figure 3j-l**).

108 Using the above three assessment mechanisms datasets Li and Loi were
109 excluded. We also used correlation between individual mRNA abundance
110 profiles in order to identify potentially redundant patients across studies. This
111 caters for patients which might have participated in more than one study or
112 duplicate data used in multiple studies. The survival data of patients with
113 extremely high correlation coefficient (Spearman's $\rho \geq 0.98$) was matched, and
114 we found 22 samples^{2, 3} having identical survival time and status. These
115 patients were removed from further analyses (**Supplementary Figure 3m**).

116 Cohorts of primary colon, lung and ovarian cancer patient mRNA profiles were
117 assembled in similar ways, however, without outlier detection due to relatively
118 small number of publicly available datasets and no (data curation based)
119 evidence of sample sharing between studies (**Supplementary Tables 3-5**).

120 **2. SIMMS' comparison with other machine learning algorithms and**
121 **genesets/pathway scoring methods**

122 In order to benchmark the prognostic ability of subnetworks by SIMMS' model N;
123 using genes in each subnetwork, we fitted a Cox proportional hazard model
124 using forward selection, backward elimination (R package: MASS v7.3-47),
125 LASSO $L1$ regularization, and ridge regularization (R package: glmnet v2.0-10),
126 as well as a random survival forest (R package: randomForestSRC v2.5.0). To
127 tune the hyperparameter of the regularized Cox models and the random survival
128 forest we applied a grid search algorithm based on cross-validation in the training
129 sets. The final models were tested in the validation cohorts and predicted risk
130 scores of the Cox models (and the average cumulative hazards of the random
131 survival forest) were tested for association with patient outcome. The resulting

132 hazard ratios and respective P values (Wald test) are analog to those presented
133 in **Supplementary Tables 6b,7b,8b,9b** obtained by SIMMS.

134 To compare SIMMS prognostic performance against other genesets/pathway
135 summary scoring methods, we chose four methods representing three different
136 classes of scoring (CORGs: *t-test* based feature selection yielding pathway
137 activation scores⁴, Guo: geneset's summary scores⁵ and PCA score: geneset's
138 principal component scores⁶). Briefly, CORGs' pathway activation scores with
139 embedded t-statistic based feature selection were estimated using training
140 datasets for each cancer type. Using selected features from training set, pathway
141 activation scores were estimated for validation datasets. For CORGs, good and
142 poor outcome samples were determined using survival time cut-off of: breast = 5
143 years, colon = 5 years, NSCLC = 3 years and ovarian = 3 years removing any
144 samples censored prior to the cut-off time; consistent with other analyses in this
145 manuscript. Guo et al. scores were estimated using mean and median
146 expression levels of genes in a given subnetwork, yielding effectively two
147 different scores which were treated as two independent methods. PCA scores
148 were estimated by using the first principal component as representative summary
149 measure of genes in a subnetwork, which is analogous to estimates used in Bild
150 et al.⁶ and gsdScore⁷.

151 For each subnetwork scoring method, Cox proportional hazards model was fitted
152 on training datasets and applied to predict risk scores using validation datasets.
153 These predicted risk scores were dichotomised on training set median risk score
154 and resulting groups were tested for association with patient outcome using Cox
155 model. The results of Cox model were compared across various methods and
156 SIMMS Model N.

157 For sensitivity, 'positive' subnetworks were defined as those having at least three
158 genes significantly associated with patient outcome in the training datasets
159 (Wald-test $P < 0.05$). Here, mRNA abundance of each gene was dichotomised
160 into low- and high-risk groups and tested for survival association using a
161 univariate Cox proportional hazards model. Using validation datasets, the

162 proportion of correctly recovered subnetworks (Wald-test $P < 0.05$) by each
163 method were regarded as true positive rate.

164 **3. SIMMS' comparison with breast, colon, NSCLC and ovarian cancer prognostic** 165 **biomarkers**

166 In order to compare the performance of SIMMS' with existing gene expression-
167 based breast⁸, colon^{9, 10}, NSCLC¹¹⁻¹⁵ and ovarian¹⁶⁻¹⁹ cancer prognostic
168 biomarkers, we limited our search to the studies which shared the validation
169 datasets with those included in our analysis as validation datasets. This selection
170 criterion enabled unbiased comparison of hazard ratios and P values between
171 published markers and those identified by SIMMS for the same cohorts unless
172 specified otherwise. To maintain parity, strictly gene expression-based predictors
173 estimating hazard ratios were included for comparison with SIMMS. These
174 results are presented in **Supplementary Table 14**. For breast cancer biomarker,
175 previously published⁸ assessment of 9 breast cancer risk predictors were
176 compared against the same set of ER+ breast cancer patients in Metabric
177 Training cohort (n=801). For consistent comparison, SIMMS classifier was
178 trained on Metabric Validation cohort, and validated on Metabric Training cohort
179 predicting exactly the same (5-year) overall survival end-point as used by the
180 Zhao *et al.*⁸. To test the colon cancer 34-gene signature¹⁰ on TCGA cohort, this
181 signature was re-implemented following the original protocol. Briefly, VMC and
182 Moffitt sub-cohorts were treated as training and validation sets respectively. The
183 validation results on the Moffitt cohort (Smith) and TCGA cohort are recorded in
184 **Supplementary Table 14**. NSCLC validation was limited to lung
185 adenocarcinomas only. Both NSCLC and ovarian cancer comparisons were
186 performed in the similar way maintaining the validation cohorts for coherent
187 comparison. TCGA RNA-Seq data was used as colon and ovarian cancers
188 validation cohorts, however, panel of other markers used microarray-based
189 profiles for these two cohorts.
190 SIMMS identified markers of ER+ breast cancer compared favourably to nine
191 other breast cancer markers of clinical outcome (**Supplementary Table 14**).

192 SIMMS produced the best prognostic marker for colon cancer by a wide margin
193 compared to two other markers of patient outcome (**Supplementary Table 14**).
194 Similar trend of enhanced performance was observed for NSCLC
195 (adenocarcinomas) markers where SIMMS outperformed seven other markers in
196 3/4 independent validation studies (**Supplementary Table 14**).

197 **4. SIMMS-derived PIK3CA signaling residual risk predictor in early breast cancer**

198 *4.1 TEAM cohort power calculations*

199 Power calculations were performed on complete TEAM cohort (n = 3,476; events
200 = 507) and for each of the training (n = 1,734; events = 250) and validation (n =
201 1,742; events = 257) subsets separately. Power estimates representing the
202 likelihood of observing a specific HR against the above-mentioned events,
203 (assuming equal-sized patient groups) were derived using the following formula
204 (2):

$$z_{power} = \frac{\sqrt{E} \times \ln(HR)}{2} - z(1 - \frac{\alpha}{2}) \quad (2)$$

205 where E represents the total number of events (DRFS) and α represents the
206 significance level which was set to 10^{-3} . z_{power} was calculated for HR ranging from
207 1 to 3 with steps of 0.01.

208 *4.2 mRNA abundance data processing*

209 Raw mRNA abundance counts data were pre-processed using R package
210 NanoStringNorm²⁰ (v1.1.16). In total, 252 pre-processing schemes were
211 evaluated; parameterising normalization with respect to six positive controls,
212 eight negative controls and six housekeeping genes (GUSB, PUM1, SF3A1,
213 TBP, TFRC and TMED10) followed by global normalization. To identify the

214 optimal pre-processing parameters, two criteria were defined. First, each of the
215 252 pre-processing schemes was ranked based on their ability to maximize
216 Euclidean distance of ERBB2 mRNA abundance between HER2-positive and
217 HER2-negative samples. The process was repeated for 1000 random subsets of
218 HER2-positive and HER2-negative samples for each of the pre-processing
219 schemes. Second, using 37 replicates of an RNA pool extracted from 5 randomly
220 selected anonymized FFPE breast tumour samples, pre-processing schemes
221 were ranked based on inter-batch variation. To this end, mixed effects linear
222 models were used and residual estimates were used as a measure of inter-batch
223 variation (R package: nlme v3.1-113). Cumulative ranks based on these two
224 criteria were estimated using RankProduct¹ resulting in selection of an optimal
225 pre-processing scheme of normalisation to the *geometric mean* derived from all
226 genes followed by *rank normalisation*. Samples with RNA content |z-score| > 6
227 were discarded as being potential outliers. Only one sample was removed from
228 the top pre-processing scheme. Six samples were run in duplicates, and their
229 raw counts were averaged and subsequently treated as a single sample. Training
230 and validation cohorts were created by randomly splitting 297 NanoString
231 nCounter cartridges into two groups (**Supplementary Table 20**), which ensures
232 that there are no batch-effects shared between the two cohorts.

233 *4.3 Survival modelling*

234 Univariate survival analysis of mRNA abundance profiles was performed by
235 median-dichotomizing every gene's mRNA abundance into high- and low-
236 abundance groups (**Supplementary Table 16**), except for *ERBB2* where risk

237 groups were determined via expectation-maximization clustering ($k=2$) because
238 of the presence of a well-established sub-population of *ERBB2* expressing
239 cancers (<15%) which are regarded as HER2/*ERBB2* positive tumours. Survival
240 analysis of clinical variables modelled age as a binary variable (dichotomized at
241 age ≥ 55 as a surrogate for menopausal status), while grade, nodal status and
242 tumour size were modelled as ordinal variables (**Supplementary Table 17**). For
243 mRNA and IHC4 models, tumour size was treated as a continuous variable.
244 Univariate survival analysis of mutational profiles (*AKT1*, *PIK3CA* and *RAS*;²¹)
245 was performed by dichotomizing patients into mutant and wild-type groups.
246 Risk score profiles (**Methods**) of patients in the Training cohort were used to fit a
247 multivariate Cox proportional hazards model alongside clinical variables. Given
248 the small number of variables to select from (continuous = 9, factors = 3) and a
249 mix of continuous and ordinal variables, we chose backwards step-wise
250 refinement algorithm (AIC penalty term: $k = 1$ degrees of freedom) and created a
251 module-based risk model (**Supplementary Table 19**). The parameters estimated
252 by the multivariate model (Training cohort) were applied to the patients in the
253 Validation cohort generating per-patient risk score. These risk scores
254 (continuous) were grouped into quartiles using the thresholds derived from the
255 Training cohort, and resulting groups were subsequently evaluated through
256 Kaplan-Meier analysis. All models were trained and validated using DRFS
257 truncated to 10 years as an end-point. All survival modelling was performed in
258 the R statistical environment (R package: survival v2.37-4).

259 *4.4 IHC4 model*

260 IHC4-protein risk scores were calculated as described by Cuzick *et al.*²², and
261 then adjusted for clinical covariates (age, nodal status, grade and tumour size).
262 Model predictions (continuous risk scores) were separated into quartiles (**Figure**
263 **5c**) and analysed using Kaplan-Meier analysis and multivariate Cox proportional
264 hazards model adjusted for clinical variables.

265 *4.5 Recurrence probabilities*

266 Recurrence probabilities at 5 years were estimated by binning the predicted risk
267 scores in 25 equal groups. For each group, recurrence probability $R_{(t)}$ was
268 estimated as $1-S_{(t)}$, where $S_{(t)}$ is the Kaplan-Meier survival estimate at year 5.
269 The $R_{(t)}$ estimates of 25 groups were smoothed using local polynomial regression
270 fit. The predicted estimates were plotted against the median risk score of each
271 group except the first and last group, where the lowest risk score and 99th
272 percentile were used, respectively.

273 *4.6 Performance Assessment*

274 Performance of survival models was compared through area under the *receiver*
275 *operating characteristic* (ROC) curve. Significance of difference between the
276 ROC curves was assessed through permutation analysis (10,000 permutations
277 by shuffling the risk scores while maintaining the order of survival objects).
278 Patients censored before 5 years (Training cohort: n = 192, Validation cohort: n =
279 181) were eliminated from sampling. For percentage concordance analysis,
280 patients with a relapse (after removing the afore-mentioned patients) were
281 considered as high risk and the rest of the patients were classed as low risk

282 patients. Median dichotomised risk groups determined by PIK3CA predictor and
283 IHC4 predictor were compared against the high and low risk patients. The
284 percentage of predictions matching the high and low risk groups were regarded
285 as concordant. ROC analysis was implemented using R packages pROC
286 (v1.6.0.1) and survivalROC (v1.0.3). Using the same median dichotomised risk
287 groups and actual high and low risk groups, Net reclassification improvement for
288 PIK3CA predictor over IHC4 predictor was estimated using the R package
289 PredictABEL (v1.2-1).

290 *4.7 Prognostic assessment of SIMMS PI3K modules signature in CT+/- groups*

291 SIMMS-derived PI3K modules signature was evaluated in chemotherapy-
292 stratified groups without the prior knowledge of nodal status. Patients in the
293 highest risk quartile (Q4) showed significantly decreased survival rate compared
294 to low risk patients, independent of whether they received chemotherapy (Q4 vs.
295 Q1 HR=11.07, 95%CI: 3.47-35.26; P=5.29x10⁻¹¹) (**Supplementary Figure 24e**)
296 or did not (Q4 vs. Q1 HR=9.74 95%CI: 5.58-17.02; P=1.66x10⁻²⁹)
297 (**Supplementary Figure 24f**).

298 **5. Modelling multi-modal datatypes using SIMMS**

299 Recent studies conducted by TCGA have generated datasets on matched
300 genomic and transcriptomic profiles including mutations, copy-number aberration
301 (CNA), DNA methylation and mRNA abundance^{17, 23}. These datasets can
302 potentially lead to the discovery of new biomarkers bridging the gap between
303 multi-modal molecular features and clinical covariates. To test this, we curated
304 previously published pathway modules (MEMo²⁴) from TCGA studies harbouring
305 multiple aberrations (e.g. somatic mutations, somatic copy-number aberrations,
306 dysregulated mRNA abundance levels, and DNA methylation levels)^{17, 25-27}. The

307 combined database was composed of 23 breast, 1 colorectal, 8 kidney renal
308 clear cell and 3 ovarian cancer modules (**Supplementary Table 21**). Using these
309 modules, SIMMS' (Model N) breast cancer risk predictors were created using
310 1000 randomly generated subsets (50% samples) of Metabric cohort and
311 validated on the held-out Metabric subsets as well as TCGA breast cancer
312 cohort. Similarly, 1000 randomly generated subsets (50% samples) of TCGA
313 colorectal, kidney renal clear cell carcinoma and ovarian cancers were used to
314 train and validate (50% held-out samples) the prognostic ability of each of the
315 subnetwork modules. The results of 1000 models per cancer type were
316 summarised using Fisher's method (Fisher's combined probability test) resulting
317 in a chi-square estimate and a P value. We used molecular features based on
318 mRNA and CNA as gene-level properties. Copy number levels -2 (homozygous
319 deletion) and -1 (heterozygous deletion) were collapsed into one group, whilst
320 gene copy-number levels 1 (gain) and 2 (amplification) were collapsed into a
321 single group. Copy-number levels were modelled using Cox proportional hazards
322 model and compared against the baseline copy number of 0 (diploid). Additional
323 filter of minimum 3% copy-number aberration frequency in the cohort in at least
324 one group (gain/amplification and deletion) was applied prior to estimating
325 parameters for each gene, failing which would mean gene's copy-number
326 changes would not contribute to SIMMS' risk scores. Overall survival was used
327 as survival end-point for all cancer types analysed in the multi-modal modelling.

328 **6. SIMMS R package**

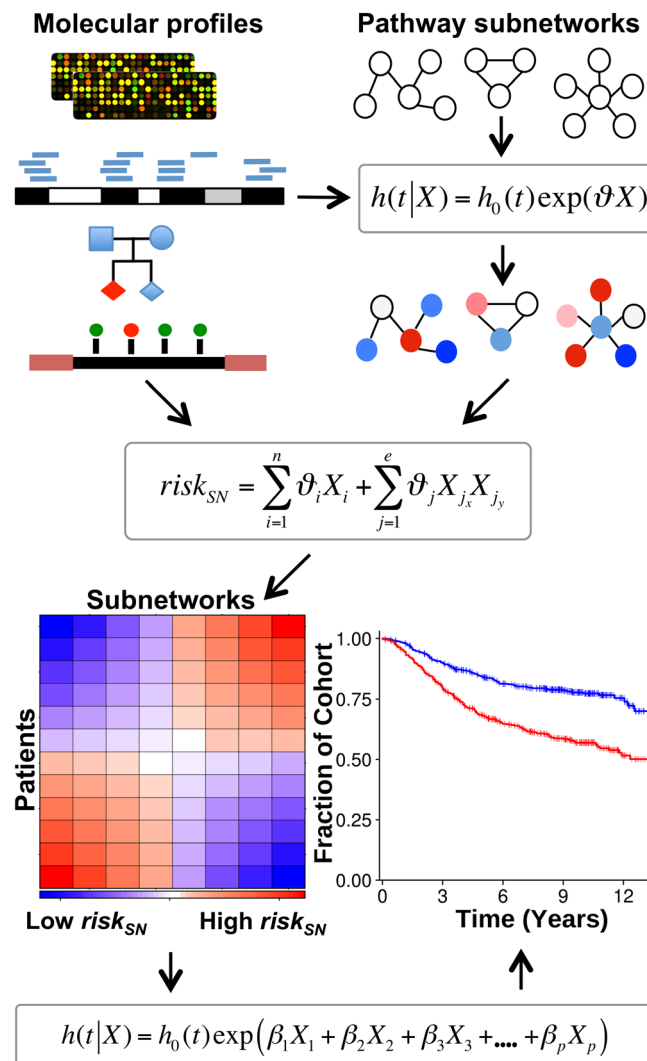
329 SIMMS is implemented in R and is available under the GNU General Public
330 License (GPL) version 2 through CRAN: [https://cran.r-](https://cran.r-project.org/web/packages/SIMMS)
331 [project.org/web/packages/SIMMS](https://cran.r-project.org/web/packages/SIMMS). SIMMS is generic and can work with any
332 combination of molecular features and interaction networks. It provides an
333 extendible framework to support user-defined parameter estimation and
334 classification algorithms. The R package of SIMMS offers three key features: (i)
335 support for multiple datatypes (mRNA, methylation, CNA etc), (ii) support for
336 user-defined networks, and (iii) support for user-defined methods for quantifying
337 dysregulation of a subnetwork. For (i), users can supply the location and names

338 of the files they would like to analyze with SIMMS. For (ii), a text file describing
339 networks in a tab-delimited format can be supplied as an input to SIMMS, see
340 `pathway_based_networks*.txt` files that comes as a part of R package. For (iii),
341 the package offers an interface function `'derive.network.features'` that accepts a
342 parameter `'feature.selection.fun'` for user-defined function name (see code
343 snippet below). By default, the function `'calculate.network.coefficients'` is called to
344 estimate MDS and risk scores for Mode N, Model E and Mode N+E as described
345 in this paper. However, users can easily write their own algorithms and simply
346 use them with SIMMS as a plug and play component. For details, see package
347 manual and vignettes.

348

```
349 derive.network.features <- function(  
350     data.directory = ".",  
351     output.directory = ".",  
352     data.types = c("mRNA"),  
353     feature.selection.fun = "calculate.network.coefficients",  
354     feature.selection.datasets = NULL,  
355     feature.selection.p.thresholds = c(0.05),  
356     subset = NULL, ...  
357 );
```

358

359 **Supplementary Figure Legends**

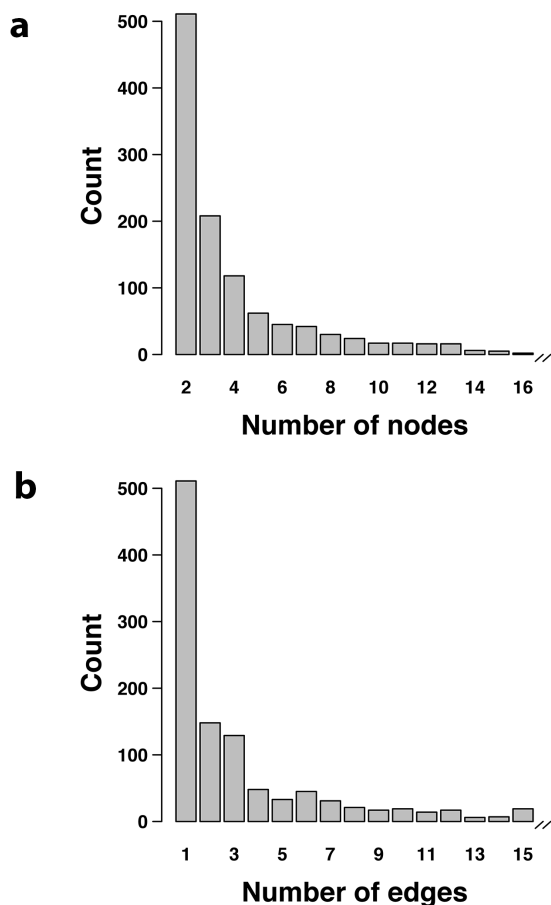
360

361 **Supplementary Figure 1**

362 **Schematic overview of SIMMS.** Subnetwork modules were extracted from NCI-
 363 Nature/Biocarta/Reactome curated pathways by isolating protein-protein
 364 interaction networks within a pathway. Molecular profiles were systemised and
 365 split into independent training and validation sets. Each extracted subnetwork
 366 was scored (module-dysregulation score) using 3 different models and ranked.
 367 High-ranking subnetworks were used to compute a patient-wise risk score. Most
 368 optimal combination of predictive subnetworks was selected using a machine

369 learning algorithm with built-in options of generalized linear models with elastic-
370 nets parameter alpha (α) supporting ridge to LASSO $L1$ -regularization ($\alpha \in [0,1]$),
371 Backward elimination and Forward selection algorithms, resulting in a
372 multivariate subnetwork-based classifier. The classifier is then tested on
373 independent validation sets.

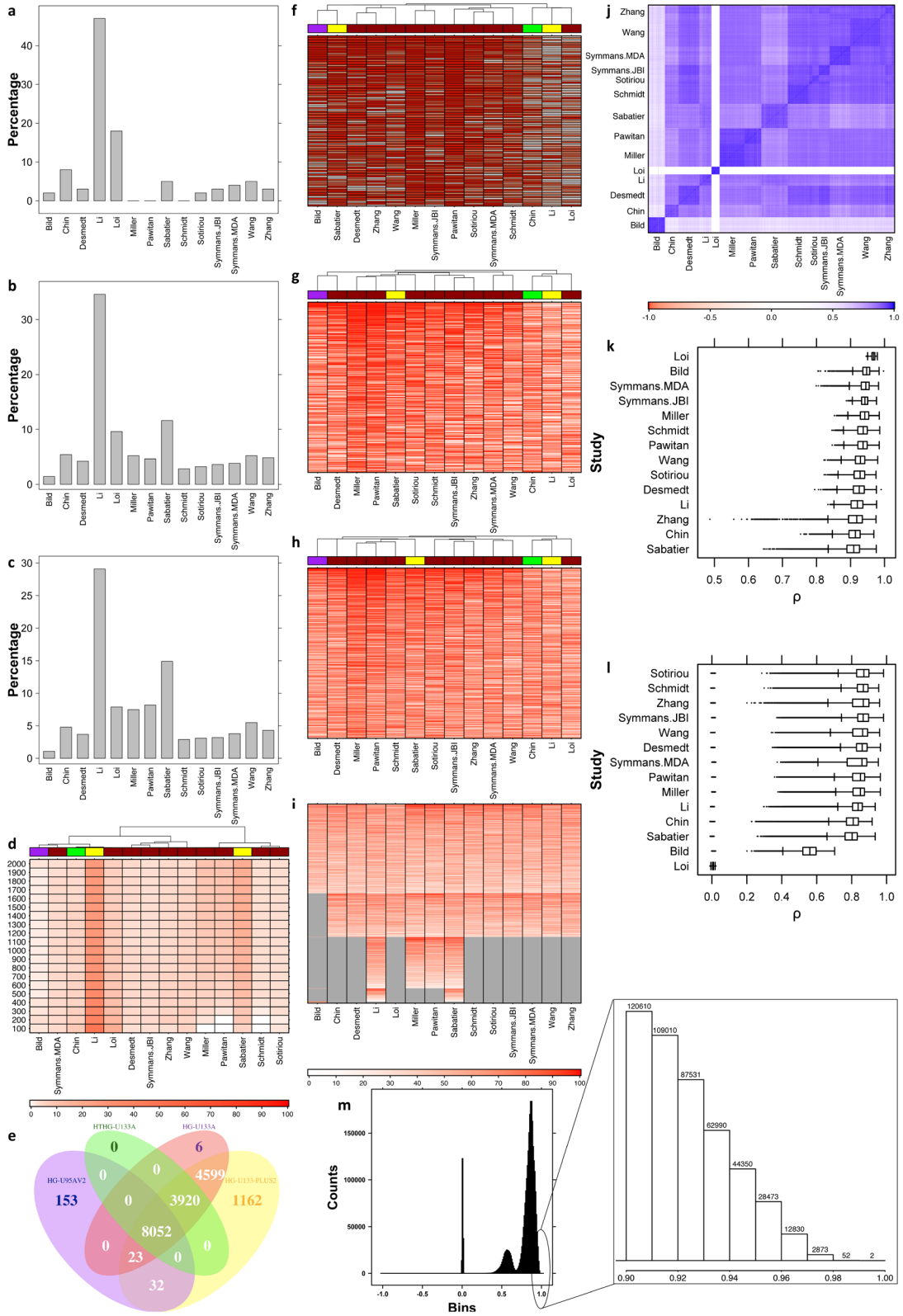
374



375

376 [Supplementary Figure 2](#)

377 **Summary of pathways database.** Distribution of nodes (a) and edges (b)
378 across all subnetwork modules extracted from NCI-Nature curated pathways
379 (Reactome and Biocarta inclusive).



380

381 **Supplementary Figure 3**

382 **Quality assessment and identification of repeated patient profiles. (a,b,c)** A
383 univariate Cox model was fit to each gene in each study in the breast cancer
384 cohort. Genes were ranked according to their P value (Wald-test), and a
385 cumulative rank for all the genes was estimated using the *rank product* for each
386 gene. The top ranked 100 (a), 500 (b) and 1,000 (c) genes were used to identify
387 the study in which each gene was farthest away from the cumulative rank. The
388 frequency of a study being farthest was recorded for each of the top ranked 100,
389 500 and 1,000 genes. Li and Loi datasets seem to be notable outliers. As the
390 threshold is relaxed, Sabatier dataset also begins to show deviation compared to
391 other datasets.

392

393 **(d)** Heatmap showing a summary of barplots (a-c) of the top ranked (rank
394 product) 100 to 2000 genes with the percentage measure as the frequency of
395 each dataset being the farthest from the rank product of top *n* genes. The
396 covariates represent different microarray platforms: HG-U95AV2=purple, HTHG-
397 U133A=green, HG-U133A=red, HG-U133-PLUS2=yellow.

398

399 **(e)** 4-way Venn diagram representing overlap of genes across the four Affymetrix
400 array platforms used in the 14 breast cancer datasets included in this study. Note
401 that the Bild dataset (array platform: HG-U95AV2) has the least number of genes
402 (8,260) with 8,052 genes that exist across all array platforms. The analysis in a-d
403 was done on this common gene set only.

404

405 **(f,g,h)** Gene ranks transformed into percentile ranks within all studies. The rank
406 product based top 100 (f), 500 (g), and 1,000 (h) genes shown in terms of their
407 percentile rank within each study. Li, Loi and Chin datasets clustered together
408 and had lower percentile ranks compared to other datasets. However, Sabatier's
409 percentile ranks were similar to other datasets thereby deemed suitable for
410 inclusion in this study.

411

412 **(i)** Summary heatmap of percentile ranks across all studies, ordered by groups of
413 genes common across studies, thereby maintaining coherent comparison of
414 ranks.

415

416 **(j)** Heatmap of Spearman correlation between patients' mRNA abundance
417 profiles. Loi dataset quite clearly shows weak correlation with the other datasets,
418 again reflecting unusual expression patterns compared to other datasets.

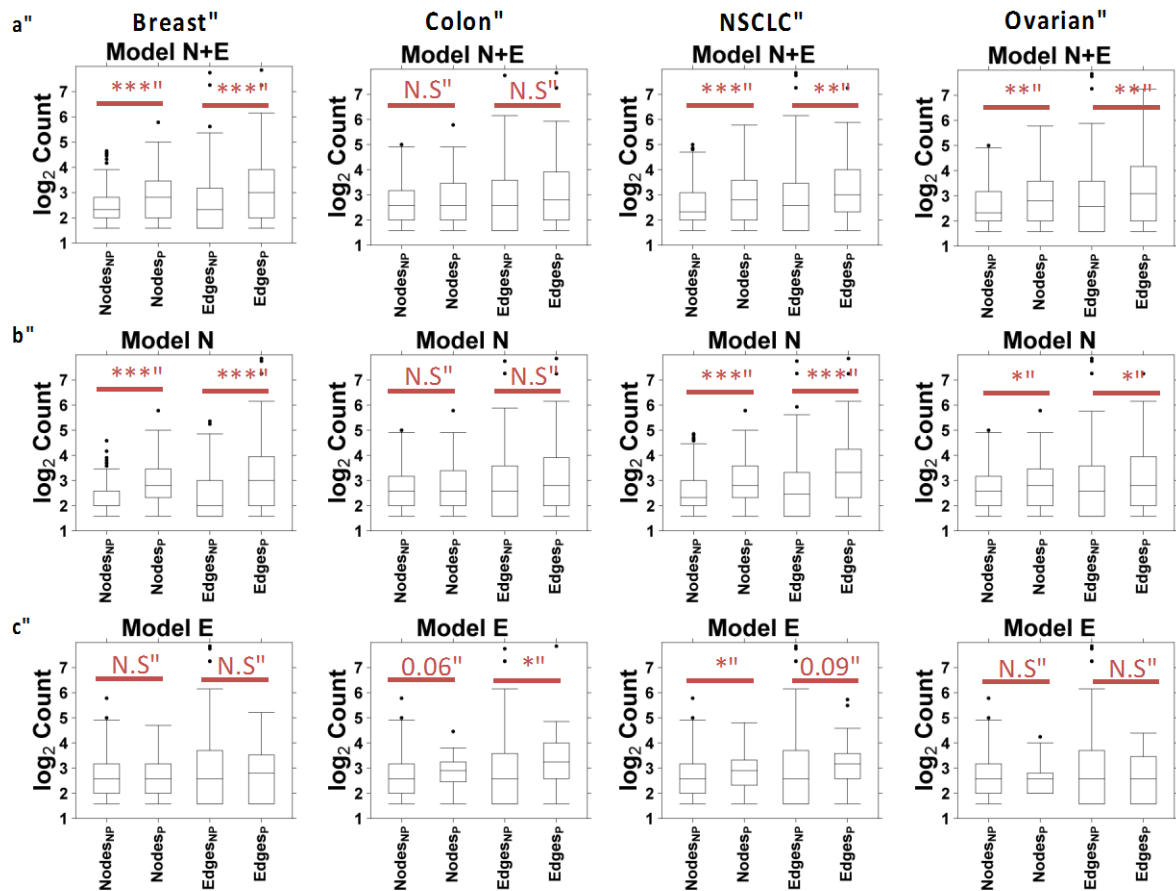
419

420 **(k,l)** Box-whisker plots of intra- (k) and inter-study (l) correlation between
421 patients' mRNA abundance profiles. The results show distinctively strong
422 correlation within Loi dataset (k) and weak correlation between Loi and other
423 datasets (l). Boxplot lines show lower quartile, median and upper quartile.
424 Whiskers extend to the point closest to the upper/lower quartile $\pm (1.5 \times \text{IQR})$.

425

426 **(m)** Histogram of Spearman correlation of patients' mRNA abundance profiles.
427 From left to right, the first peak represents correlation between Loi and other
428 datasets. The second peak represents correlation between Bild and other
429 datasets, while the third peak constitutes the correlation between the remaining
430 datasets. The survival data of highly correlated profiles (zoomed in panel, $0.98 \leq$
431 $\rho \leq 1.00$) was further inspected, resulting in 22 patients that were found in both
432 Sotiriou and Symmans (JBI) datasets having identical survival data. These were
433 removed from Symmans (JBI) dataset for further analysis.

434

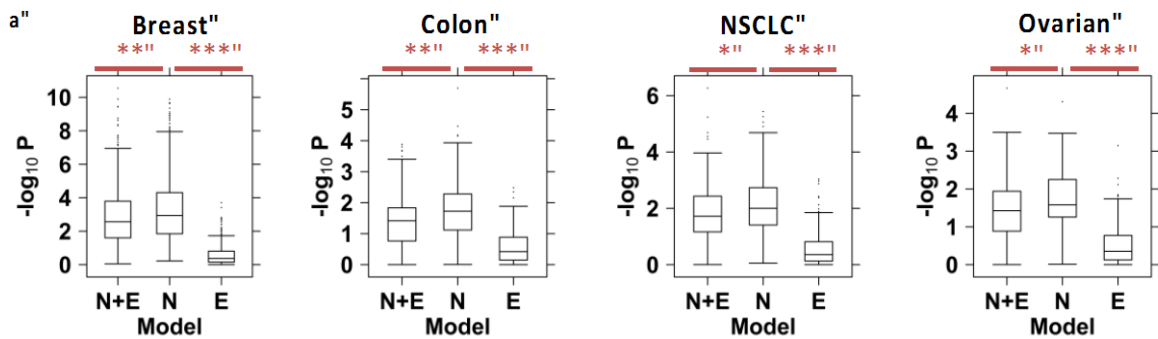


435

436 **Supplementary Figure 4**

437 **Distribution of prognostic ability versus the size of subnetworks. (a-c)** For
 438 each of the three scoring schemes *i.e.* Model N+E, Model N and Model E (see
 439 Methods), distribution of subnetwork size for prognostic (**P**) (Wald test $P < 0.05$;
 440 validation cohorts) and not prognostic (**NP**) subnetwork modules. Size of a
 441 subnetwork was defined in terms of number of nodes and number of edges.
 442 Pairwise comparisons were performed using Wilcoxon rank sum test (* $P < 0.05$, **
 443 $P < 0.01$, *** $P < 0.001$, N.S. $P > 0.1$). Boxplot lines show lower quartile, median and
 444 upper quartile. Whiskers extend to the point closest to the upper/lower quartile \pm
 445 $(1.5 \times \text{IQR})$.

446

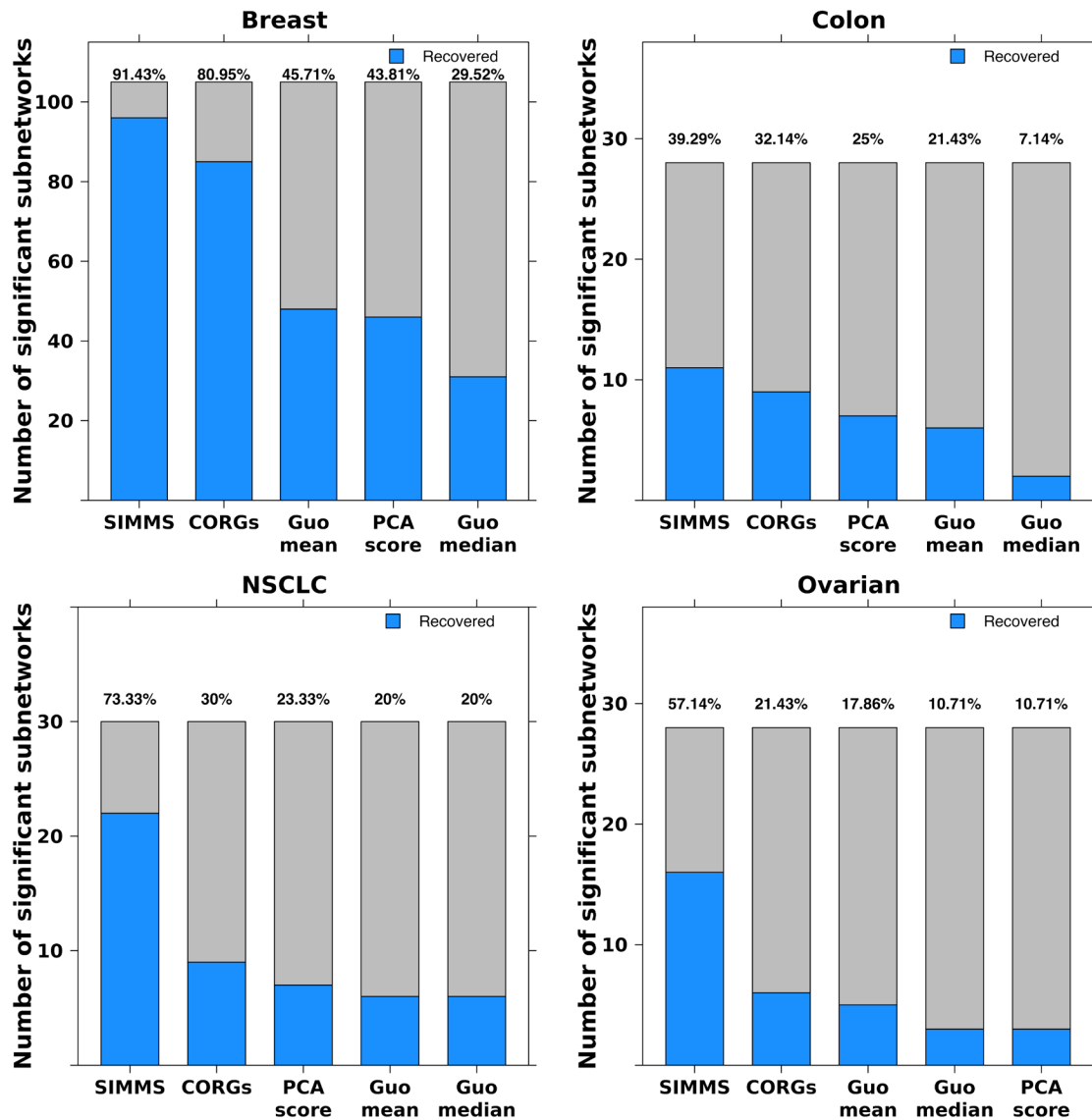


447

448 **Supplementary Figure 5**

449 **Prognostic ability of SIMMS' models (a)** Distribution of prognostic ability ($-\log_{10} P$) of subnetwork modules which were significant (Wald test $P < 0.05$) in at
 450 least one scoring scheme (Model N+E, Model N and Model E), in respective
 451 cancer type. $-\log_{10} P$ values were compared using one-way ANOVA ($P < 0.05$)
 452 followed by Tukey HSD test. Tukey HSD test's adjusted P values for only Model
 453 N vs Model N+E and Model E are displayed (* $P < 0.05$, ** $P < 0.01$, *** $P < 0.001$).
 454 Boxplot lines show lower quartile, median and upper quartile. Whiskers extend to
 455 the point closest to the upper/lower quartile $\pm (1.5 \times \text{IQR})$.
 456

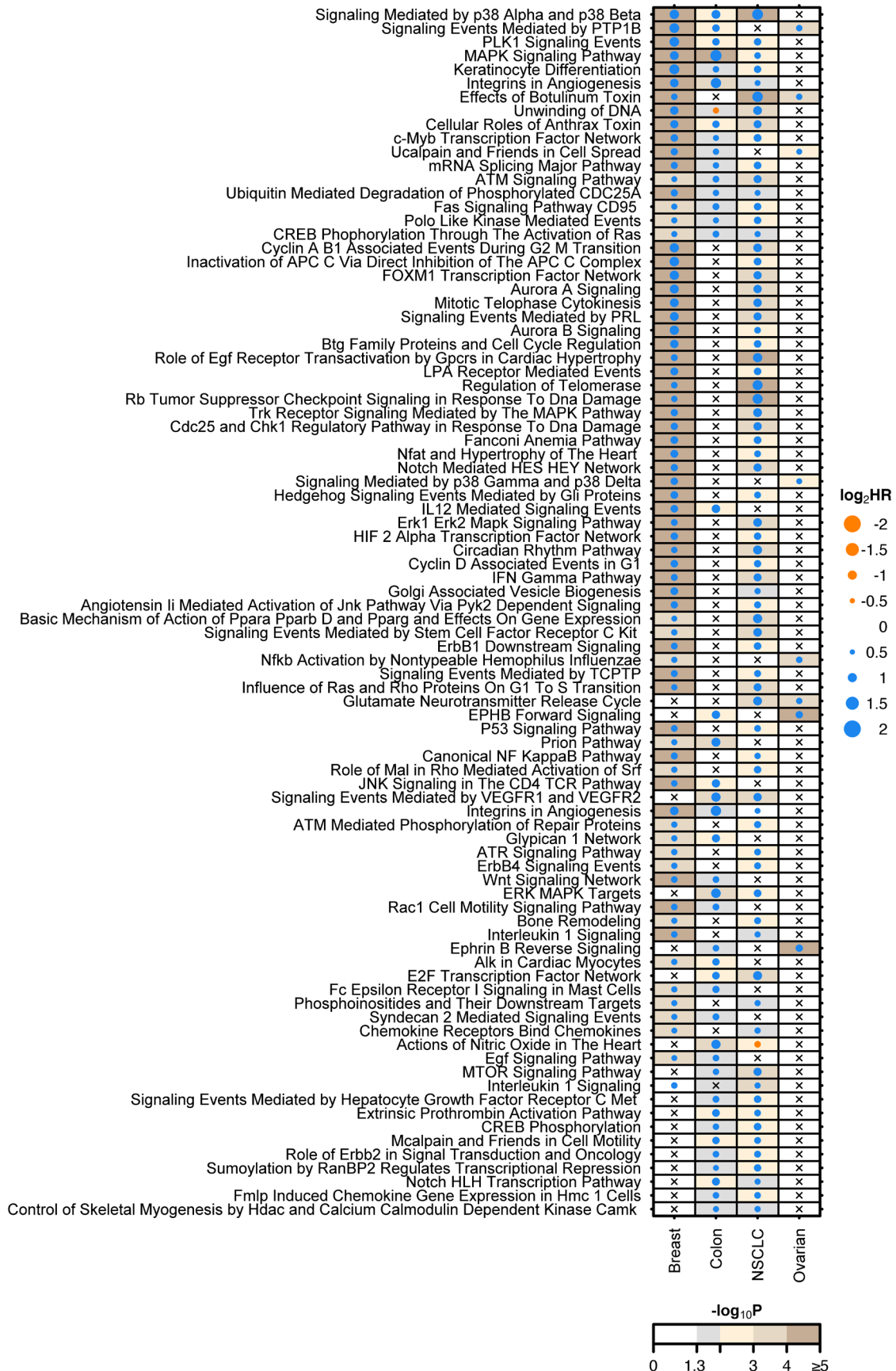
457



458

459 **Supplementary Figure 6**

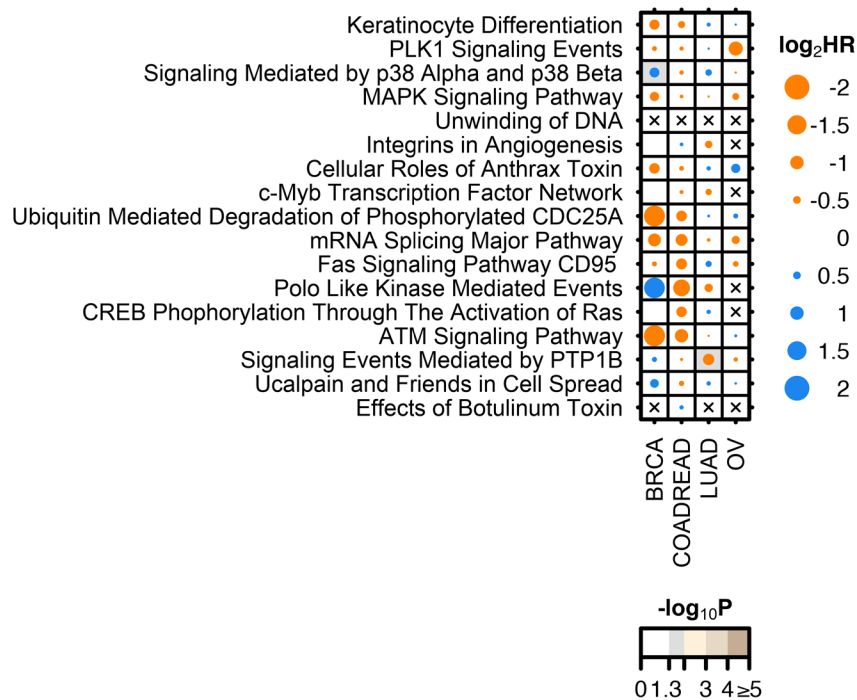
460 **Comparison of subnetwork scoring methods.** Sensitivity assessment of
 461 correctly recovered 'positive' subnetwork modules (those likely to be associated
 462 with patient outcome) by various subnetwork/pathway scoring methods. Height of
 463 each bar represents total number of 'positive' subnetworks, while the blue colour
 464 shows proportion of correctly recovered 'positive' subnetworks. Numbers above
 465 the bars represent % true positive rate.



467 **Supplementary Figure 7**

468 **Prognostic assessment of SIMMS' predicted risk scores.** Dot plot of hazard
 469 ratios and P values of subnetwork modules significant in at least 2/4 cancer
 470 types. A Cox proportional hazards model was fitted to dichotomous risk scores
 471 (threshold derived from the training cohort) across the entire validation cohort.
 472 Crosses represent absence of subnetwork module from a particular cancer type.

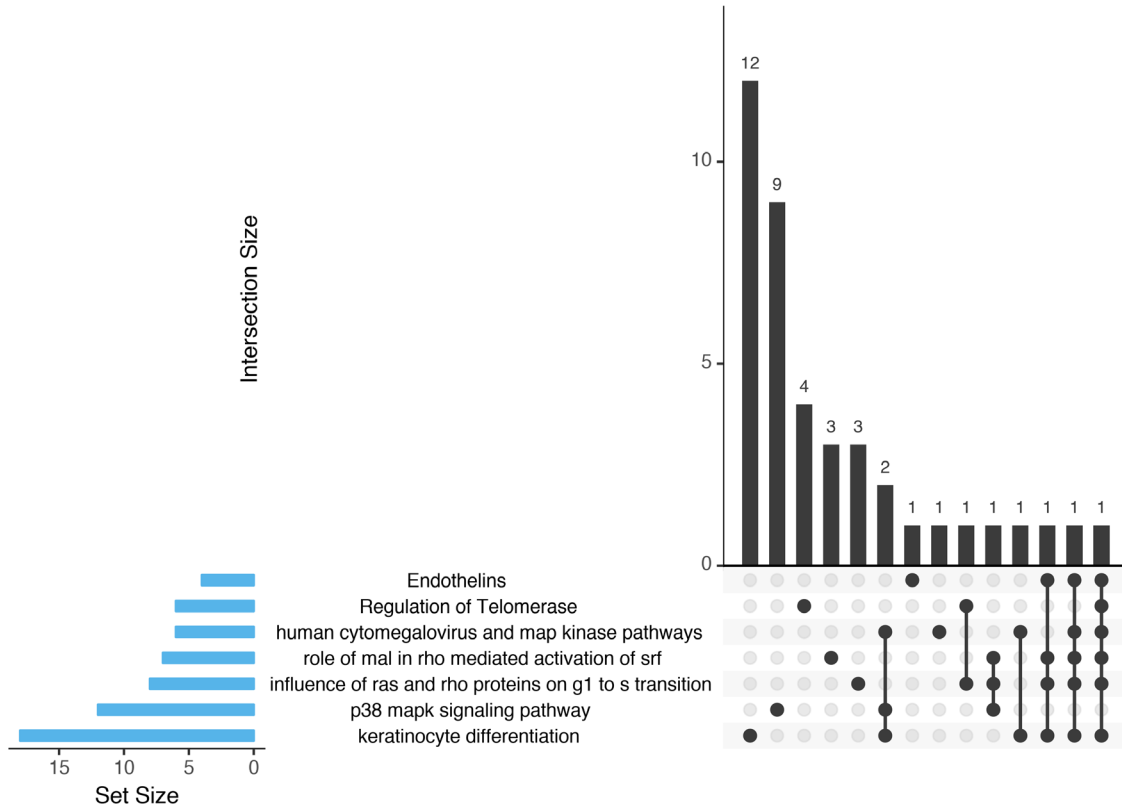
473



474

475 **Supplementary Figure 8**

476 **Prognostic assessment of mutation burden.** Dot plot of hazard ratios and P
 477 values of subnetwork modules in **Figure 1i**. Using TCGA datasets for breast,
 478 colorectal, lung adenocarcinoma and ovarian cancers; for each of these
 479 subnetwork modules (using mutations in genes involved), patients were assigned
 480 to mutant group if any gene in the subnetwork was mutated, otherwise to non-
 481 mutant group. A Cox proportional hazards model was fitted to test association of
 482 these groups with patient outcome.

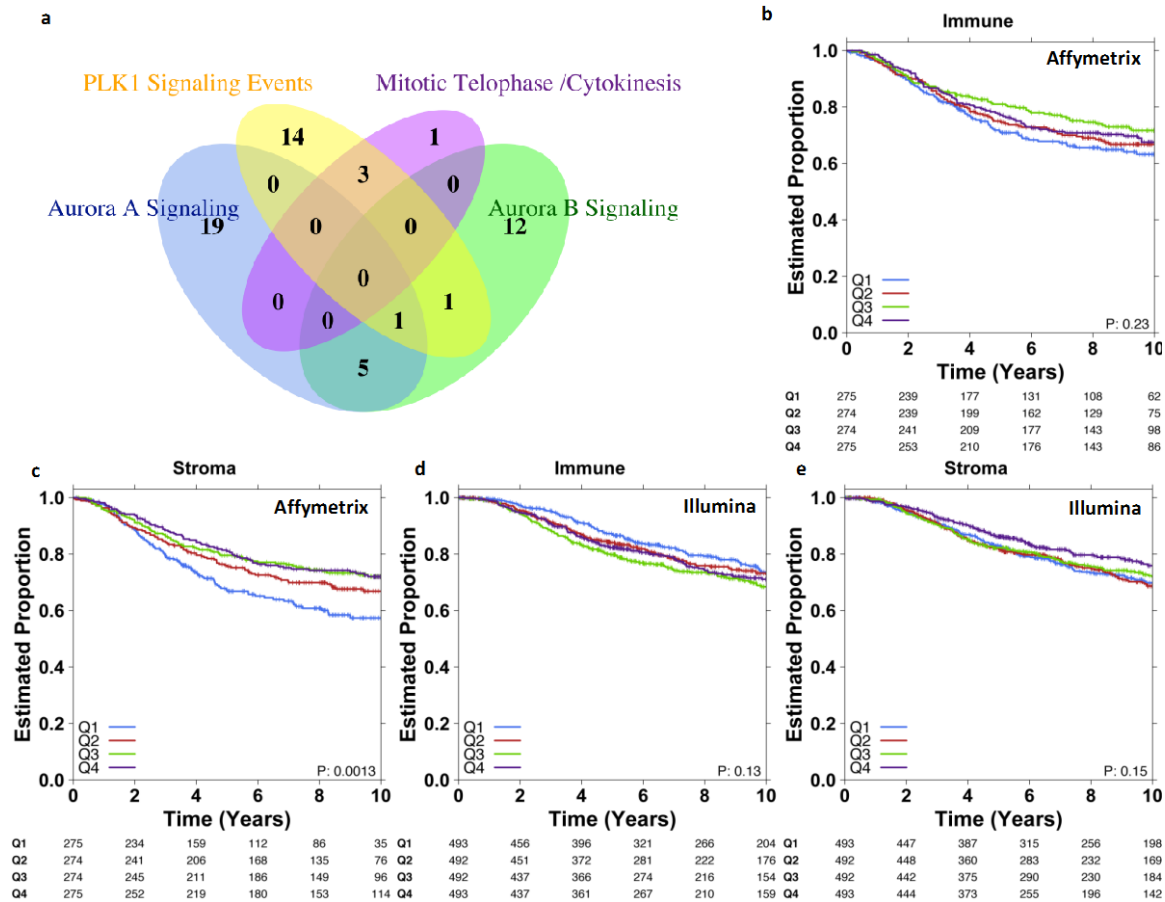


483

484 **Supplementary Figure 9**

485 **Overlap of genes in subnetworks with both prognostic and predictive**
 486 **ability.** Upset plot showing overlap of genes between subnetworks which
 487 showed significant prognostic as well as predictive (platinum response)
 488 association in TCGA ovarian cancer cohort.

489



490

491 **Supplementary Figure 10**

492 **Overlap of genes in cell cycle subnetwork modules, and prognostic**

493 **assessment of immune and stromal scores. (a)** Venn diagram showing

494 overlapping genes between proliferation subnetwork modules derived from the

495 pathways of Aurora A signaling (module 1), Aurora B signaling (module 1), *PLK1*

496 signaling events (module 1) and Mitotic Telophase/Cytokinesis (module 1). The

497 maximal overlap was of a single gene (*AURKA*) common across three modules

498 (Aurora A, Aurora B and *PLK1* modules). Module number in parenthesis refers to

499 unique module number within a pathway in SIMMS' network database (SIMMS R

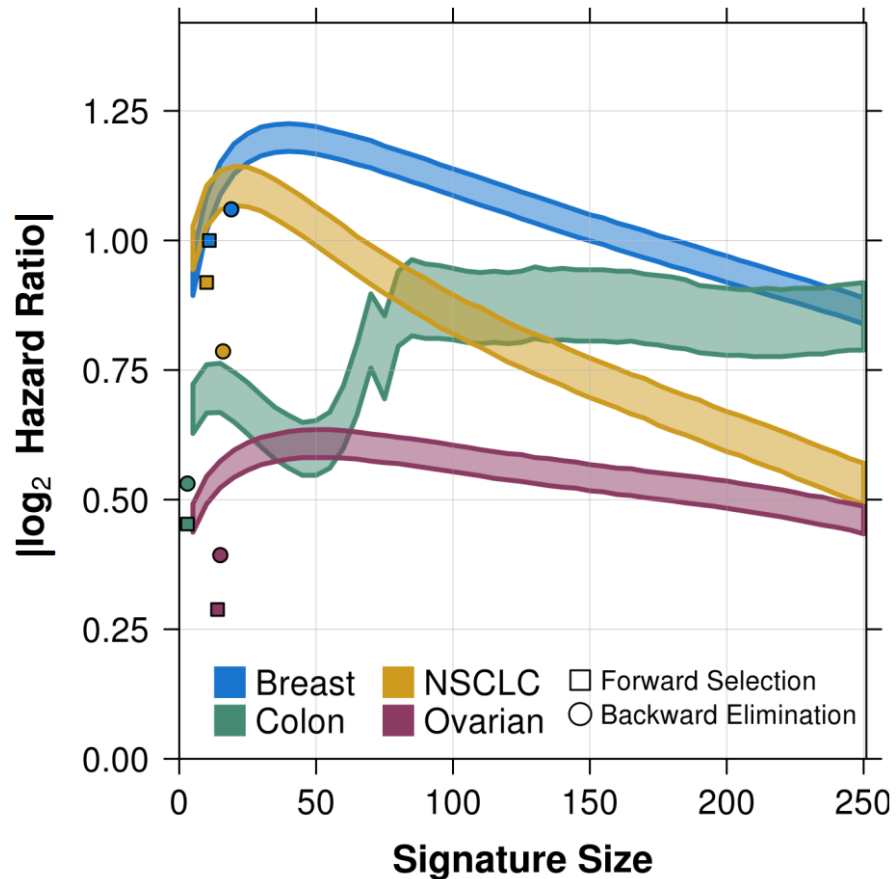
500 package). **(b, c)** Prognostic assessment of Immuno and Stromal scores

501 estimated using ESTIMATE in Affymetrix based breast cancer validation cohorts

502 **(Supplementary Table 2).** **(d, e)** Prognostic assessment of Immuno and Stromal

503 scores estimated using ESTIMATE in Illumina based Metabric breast cancer

504 cohort. For b-e, patient groups (Q1-Q4) were created using quantiles of
 505 Immuno/Stromal scores.



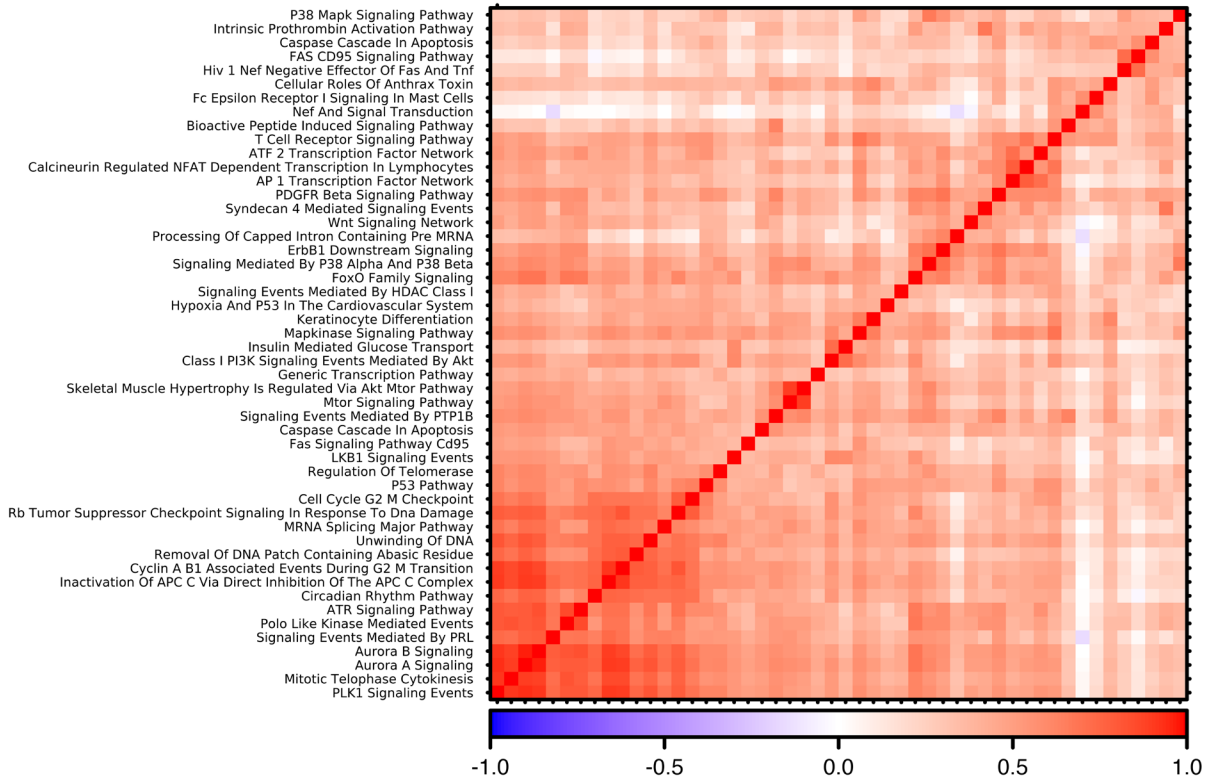
506
 507

508 **Supplementary Figure 11**

509 **Resampling of subnetworks database assessing sensitivity to initialisation**
 510 **size of SIMMS' multivariate models.** Performance (SIMMS Model N) of breast,
 511 colon, NSCLC and ovarian cancer candidate biomarkers represented as a
 512 function of marker size. Jackknifing was performed over the subnetwork marker
 513 space for various tumour types. Ten million unique markers (200,000 for each
 514 marker size $n=5,10,15,\dots,250$) were randomly sampled using all 500
 515 subnetworks regardless of their size. All biomarkers were generated using two
 516 independent machine learning paradigms; backward elimination and forward
 517 selection. The prognostic performance of each candidate biomarker was
 518 measured by taking the absolute value of the \log_2 -transformed hazard ratio

519 estimated with a multivariate Cox proportional hazards model based on SIMMS
 520 Model N scores. These randomization results depict a range of prognostic
 521 performance between 75th and 95th percentiles at each marker size and were
 522 used as a guide to estimate the optimal top *n* number of subnetwork modules
 523 required to establish a multivariate classifier for a given tumour type.

524

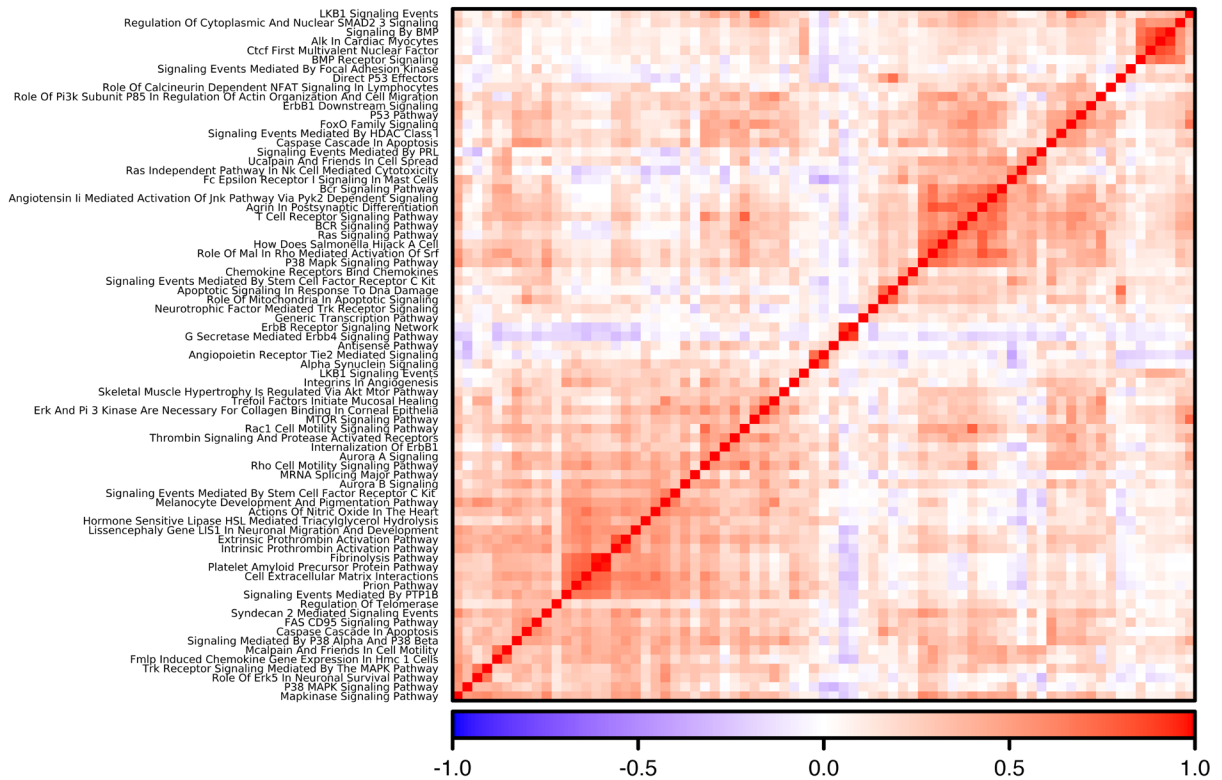


525

526 **Supplementary Figure 12**

527 **Co-expression of subnetwork risk scores in breast cancer.** Heatmap of
 528 correlation and cluster analysis of patient's risk score of top ranked 50
 529 subnetwork modules of breast cancer (validation datasets only). The plot
 530 displays activity of subnetworks as well as clusters of highly co-expressed
 531 modules as indicated in dark red clusters.

532

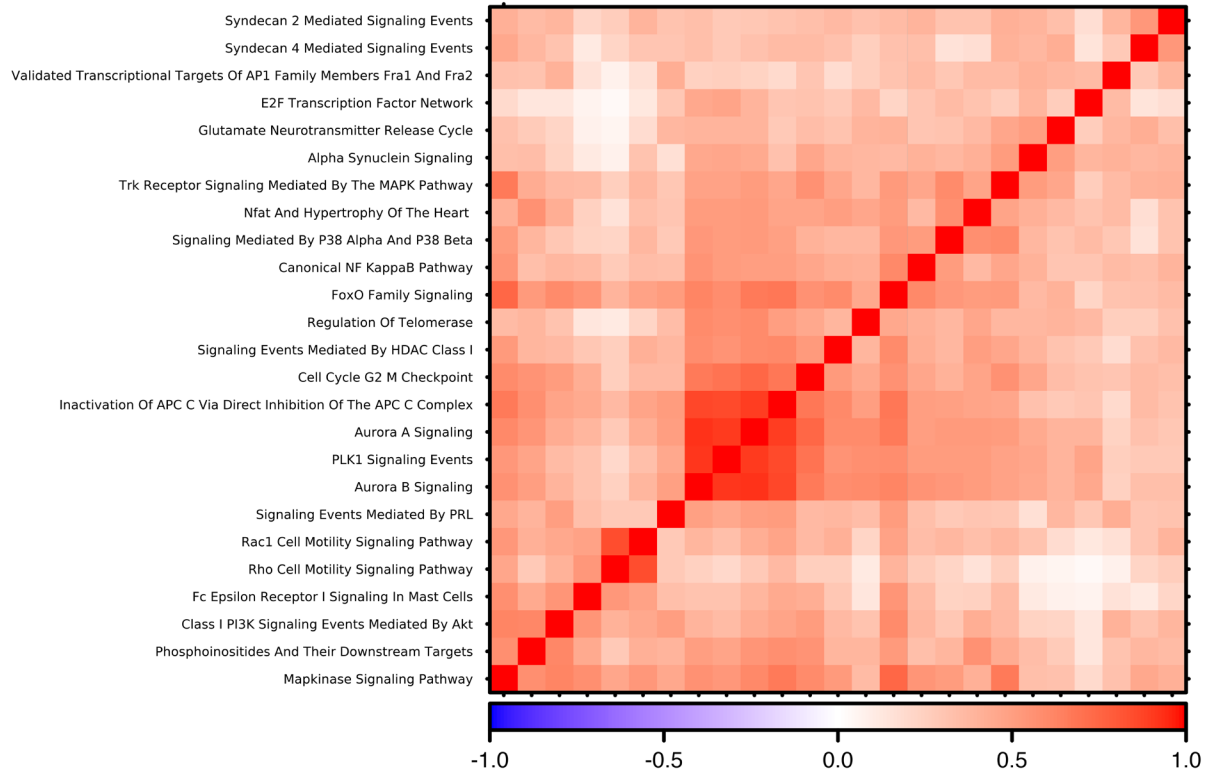


533

534 **Supplementary Figure 13**

535 **Co-expression of subnetwork risk scores in colon cancer.** Heatmap of
 536 correlation and cluster analysis of patients' risk score of top ranked 75
 537 subnetwork modules of colon cancer (validation datasets only). The plot displays
 538 biological activity of subnetworks as well as clusters of highly co-expressed
 539 modules as indicated in dark red clusters.

540

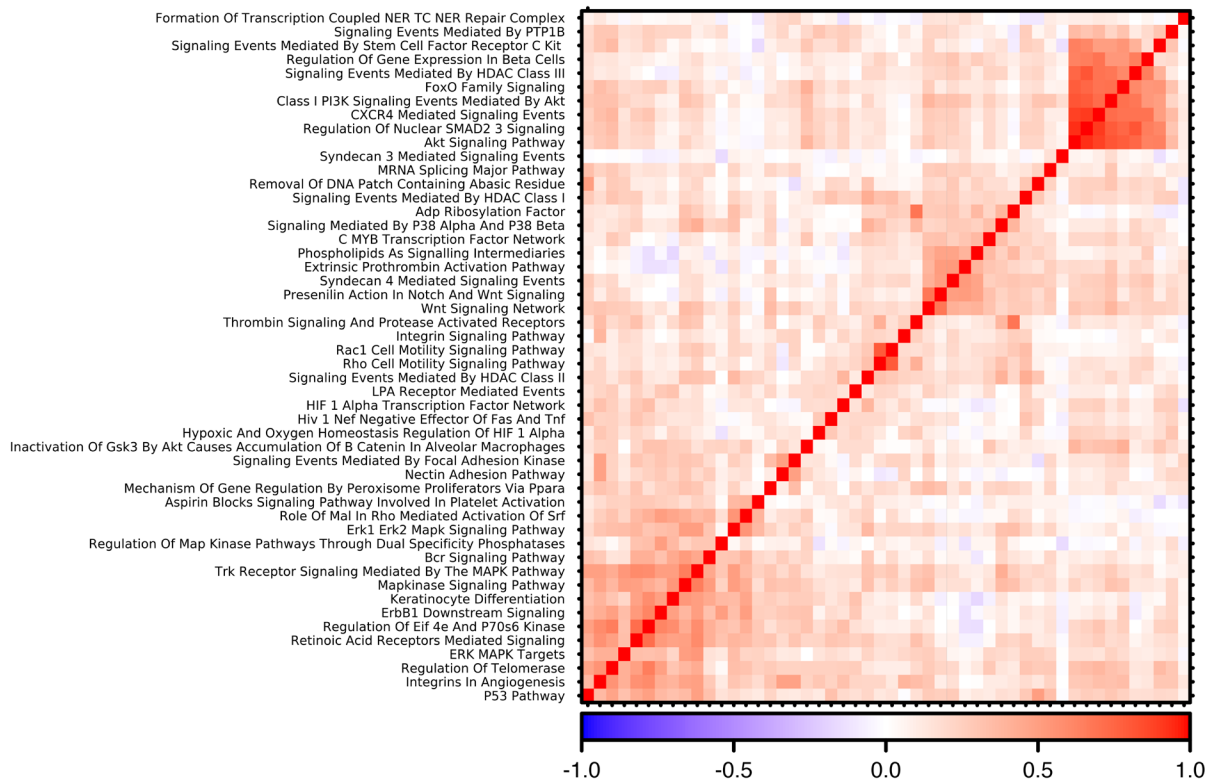


541

542 **Supplementary Figure 14**

543 **Co-expression of subnetwork risk scores in NSCLC.** Heatmap of correlation
 544 and cluster analysis of patients' risk score of top ranked 25 subnetwork modules
 545 of NSCLC (validation datasets only). The plot displays biological activity of
 546 subnetworks as well as clusters of highly co-expressed modules as indicated in
 547 dark red clusters.

548

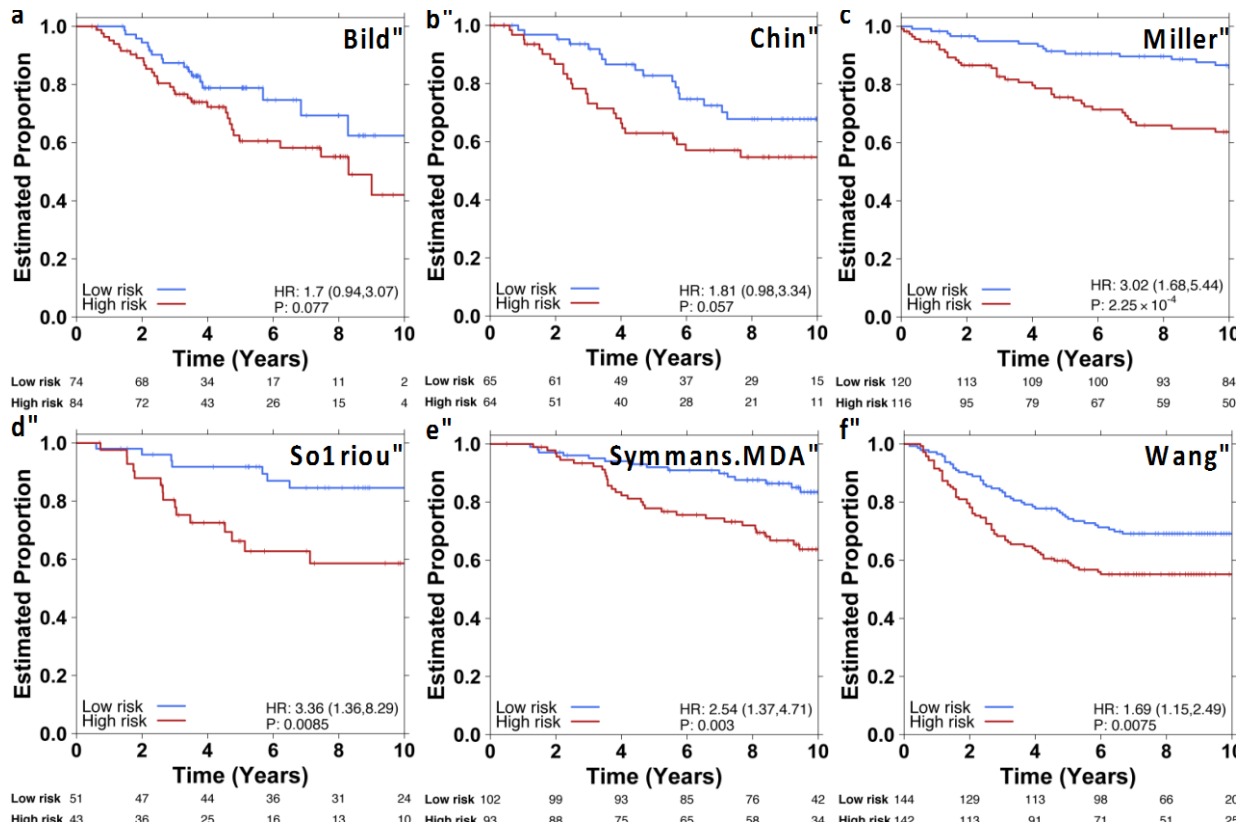


549

550 **Supplementary Figure 15**

551 **Co-expression of subnetwork risk scores in ovarian cancer.** Heatmap of
 552 correlation and cluster analysis of patients' risk score of top ranked 50
 553 subnetwork modules of ovarian cancer (validation datasets only). The plot
 554 displays biological activity of subnetworks as well as clusters of highly co-
 555 expressed modules as indicated in dark red clusters.

556

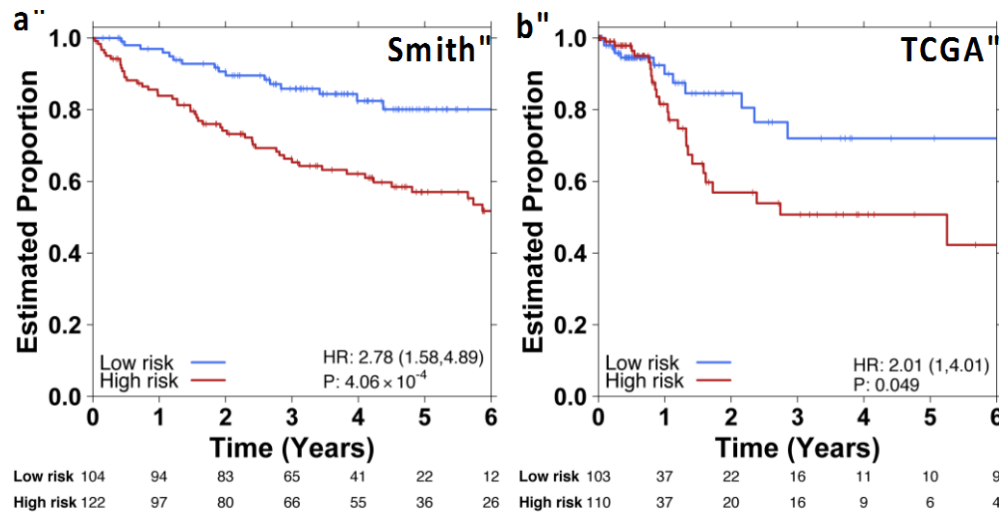


557

558 **Supplementary Figure 16**

559 **Independent validation in breast cancer cohorts.** Kaplan-Meier survival plots
 560 using SIMMS' Model N on 6 breast cancer validation sets (**Supplementary**
 561 **Table 2**) (10-year survival truncation) with subnetwork module selection
 562 performed through generalized linear models with *L1*-regularization (10-fold
 563 cross validation on training set). Model was initialised with the top ranked 50
 564 subnetwork modules.

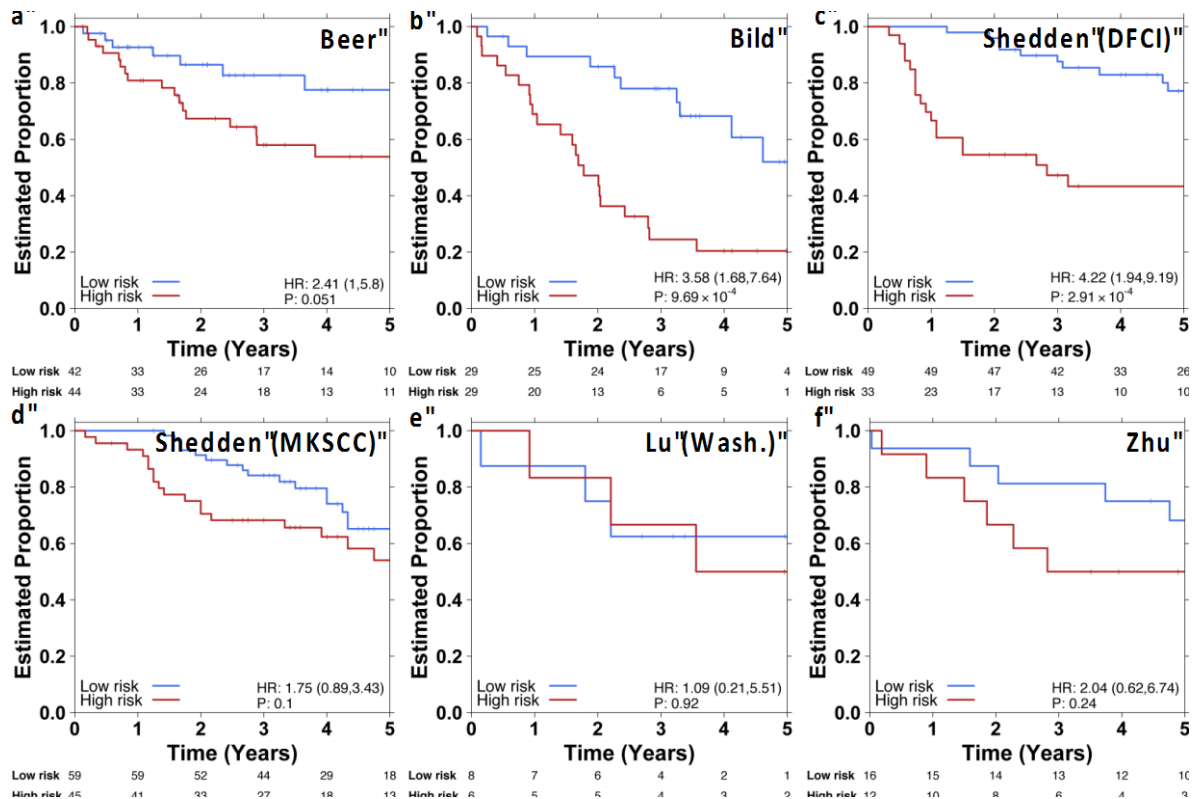
565



566

567 **Supplementary Figure 17**

568 **Independent validation in colon cancer cohorts.** Kaplan-Meier survival plots
 569 using SIMMS' Model N on 2 colon cancer validation sets (**Supplementary Table**
 570 **3**) (6-year survival truncation) with subnetwork module selection performed
 571 through generalized linear models with $L1$ -regularization (10-fold cross validation
 572 on training set). Model was initialised with the top ranked 75 subnetwork
 573 modules.



574

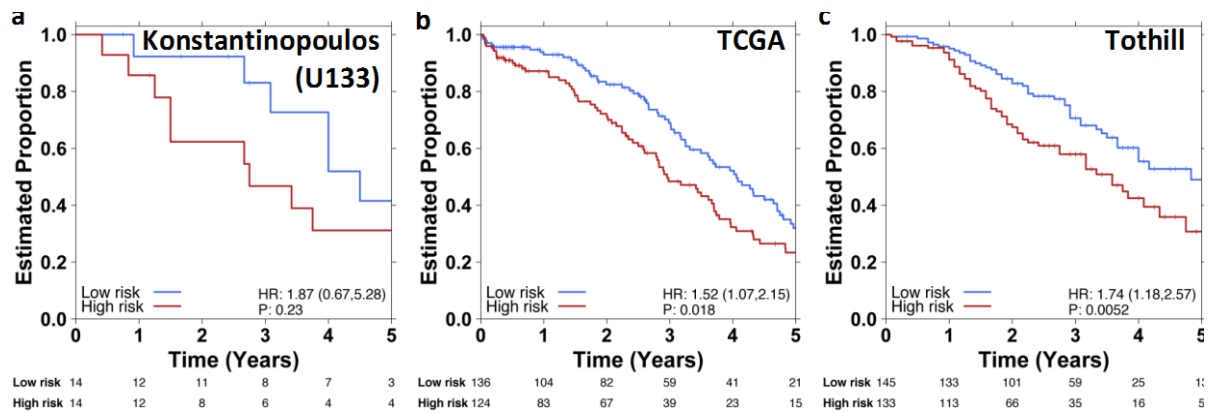
575 **Supplementary Figure 18**

576 **Independent validation in NSCLC cohorts.** Kaplan-Meier survival plots using
 577 SIMMS' Model N on 6 NSCLC validation sets (**Supplementary Table 4**) (5-year
 578 survival truncation) with subnetwork module selection performed through
 579 generalized linear models with *L1*-regularization (10-fold cross validation on
 580 training set). Model was initialised with the top ranked 25 subnetwork modules.

581

582

583



584

585 **Supplementary Figure 19**

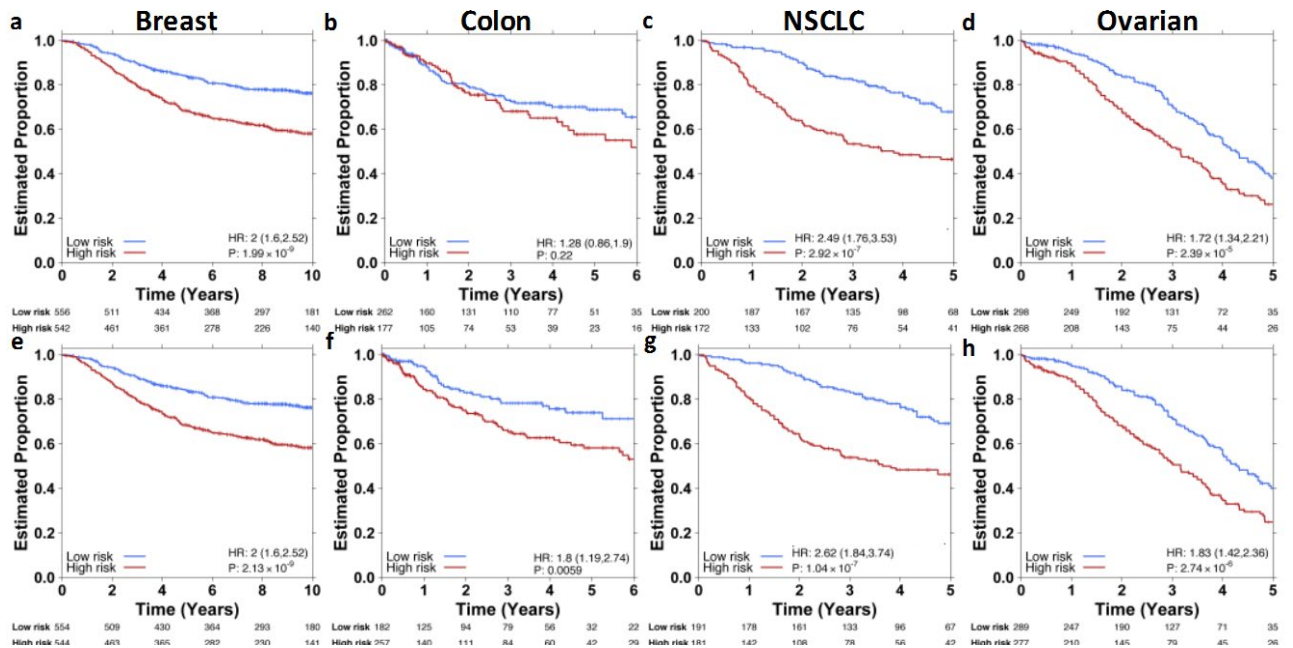
586 **Independent validation in ovarian cancer cohorts.** Kaplan-Meier survival plots
 587 using SIMMS' Model N on 3 ovarian cancer validation sets (**Supplementary**
 588 **Table 5**) (5-year survival truncation) with subnetwork module selection performed
 589 through generalized linear models with L_1 -regularization (10-fold cross validation
 590 on training set). Model was initialised with the top ranked 50 subnetwork
 591 modules.

592

593

594

595

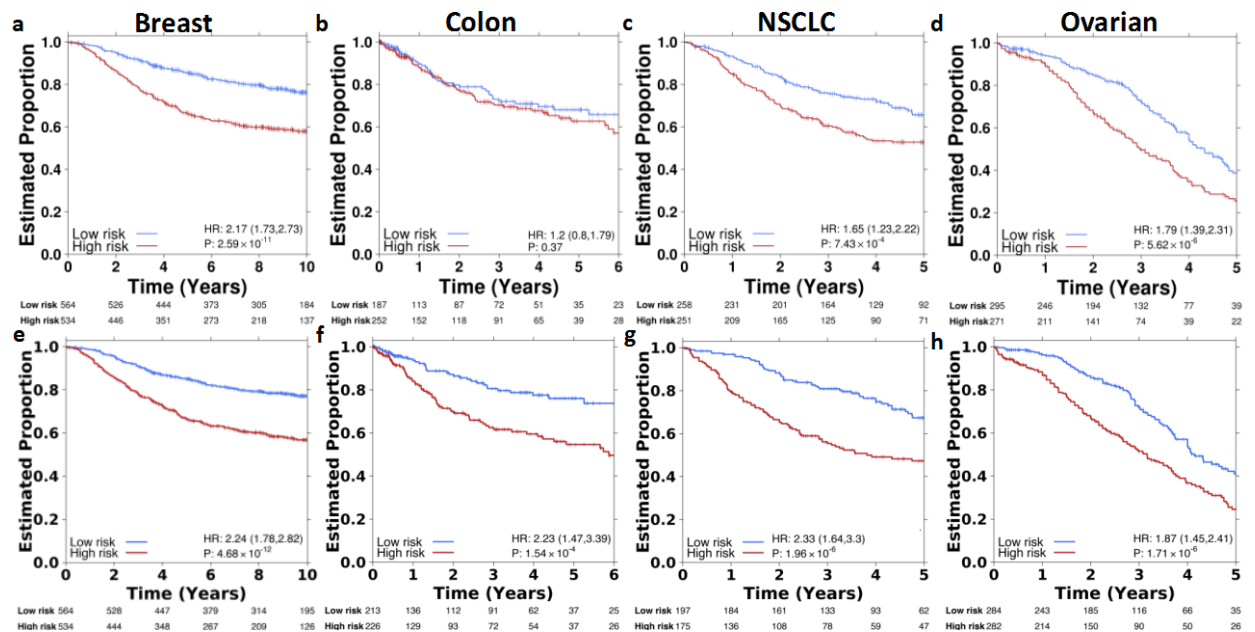


596

597 **Supplementary Figure 20**

598 **Assessment of alternative machine learning algorithms.** Kaplan-Meier
 599 survival plots of SIMMS' Model N in validation cohorts of various tumour types
 600 using alternative training algorithms; backwards elimination (**a-d**) and forward
 601 selection (**e-h**).

602



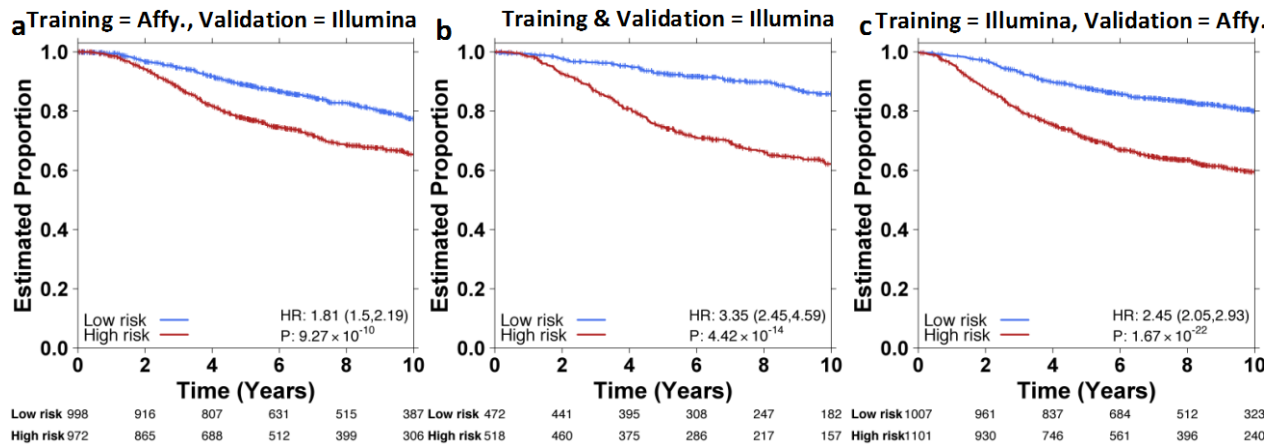
603

604 **Supplementary Figure 21**

605 **Prognostic assessment of naïve and SIMMS model with all the genes in the**
 606 **subnetwork database.** Kaplan-Meier survival plots of validation sets in each
 607 tumour type (a-d) for a Cox proportional hazard model using LASSO ($L1$ -
 608 regularization) with all genes contained in any subnetwork as model variables.
 609 (e-h) Kaplan-Meier survival plots of validation sets in each tumour type for a Cox
 610 proportional hazard model fitted using risk scores estimated by SIMMS on a
 611 single module containing all the genes across all subnetworks.

612

613



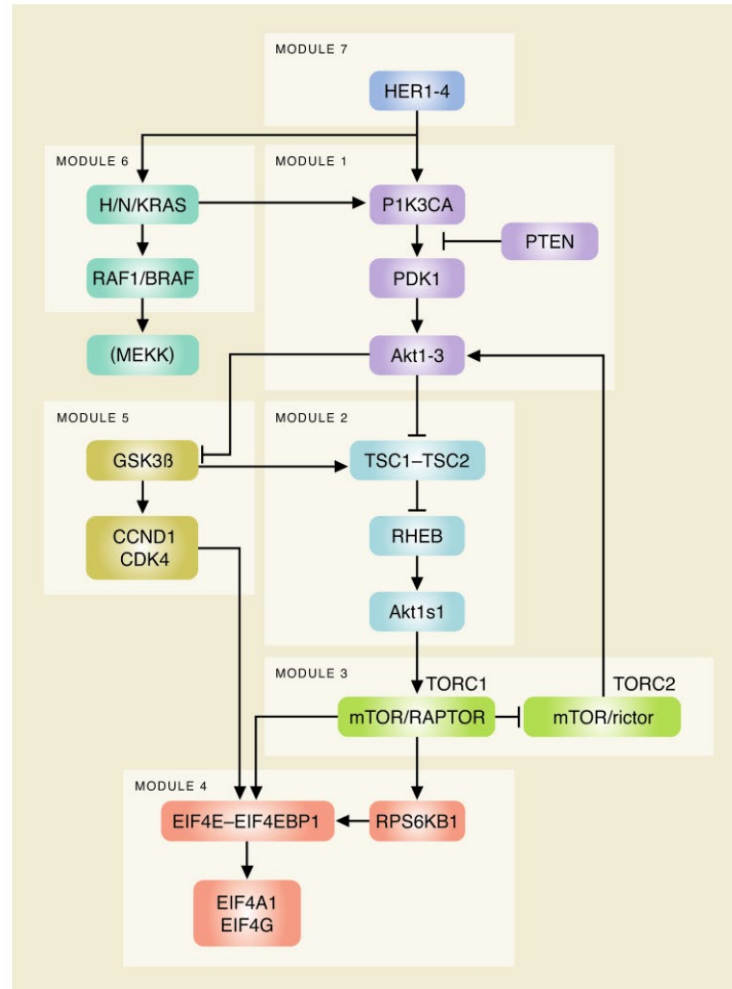
614

615 **Supplementary Figure 22**

616 **Reproducibility of SIMMS' models across mRNA quantification platforms.**

617 Kaplan-Meier survival plots of SIMMS' Model N based predictions on the
 618 Metabric validation cohort. Separate classifiers were created using the Affymetrix
 619 based breast cancer training cohorts (**Supplementary Table 2**) and Illumina
 620 based breast cancer cohort (Metabric training set). These two classifiers were
 621 validated on Illumina based breast cancer cohort (Metabric validation set) (**a,b**)
 622 and Affymetrix based breast cancer validation cohorts, respectively (**c**). All
 623 models were trained in 10-fold cross validation setting.

a''



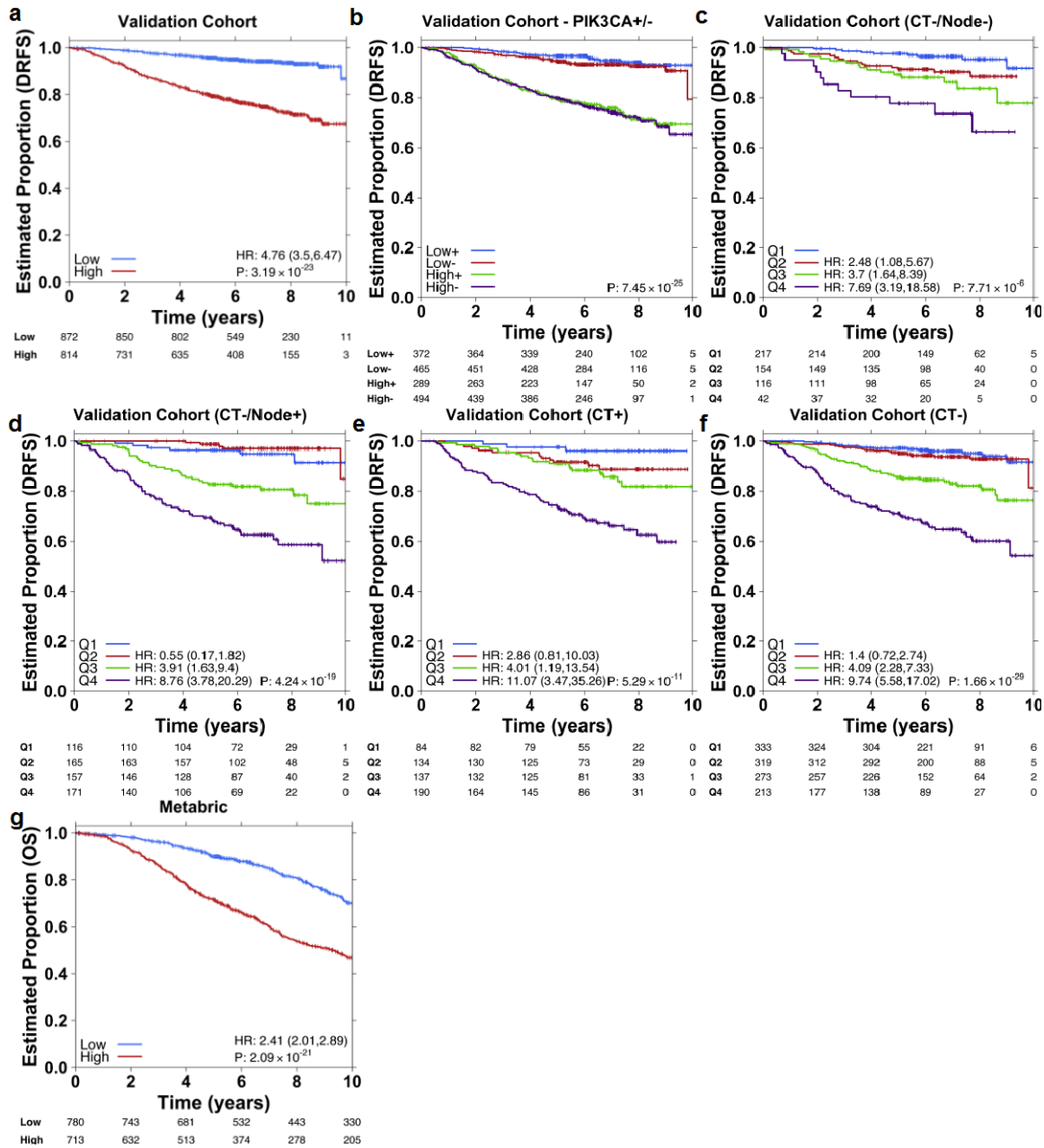
624

625

626 **Supplementary Figure 23**

627 **Schematic overview of the PI3K signalling pathway.** Figure illustrating key
 628 relationships between modules assessed in the current study. Modules 1-7 are
 629 highlighted with key signalling inter-relationships between the member genes.

630



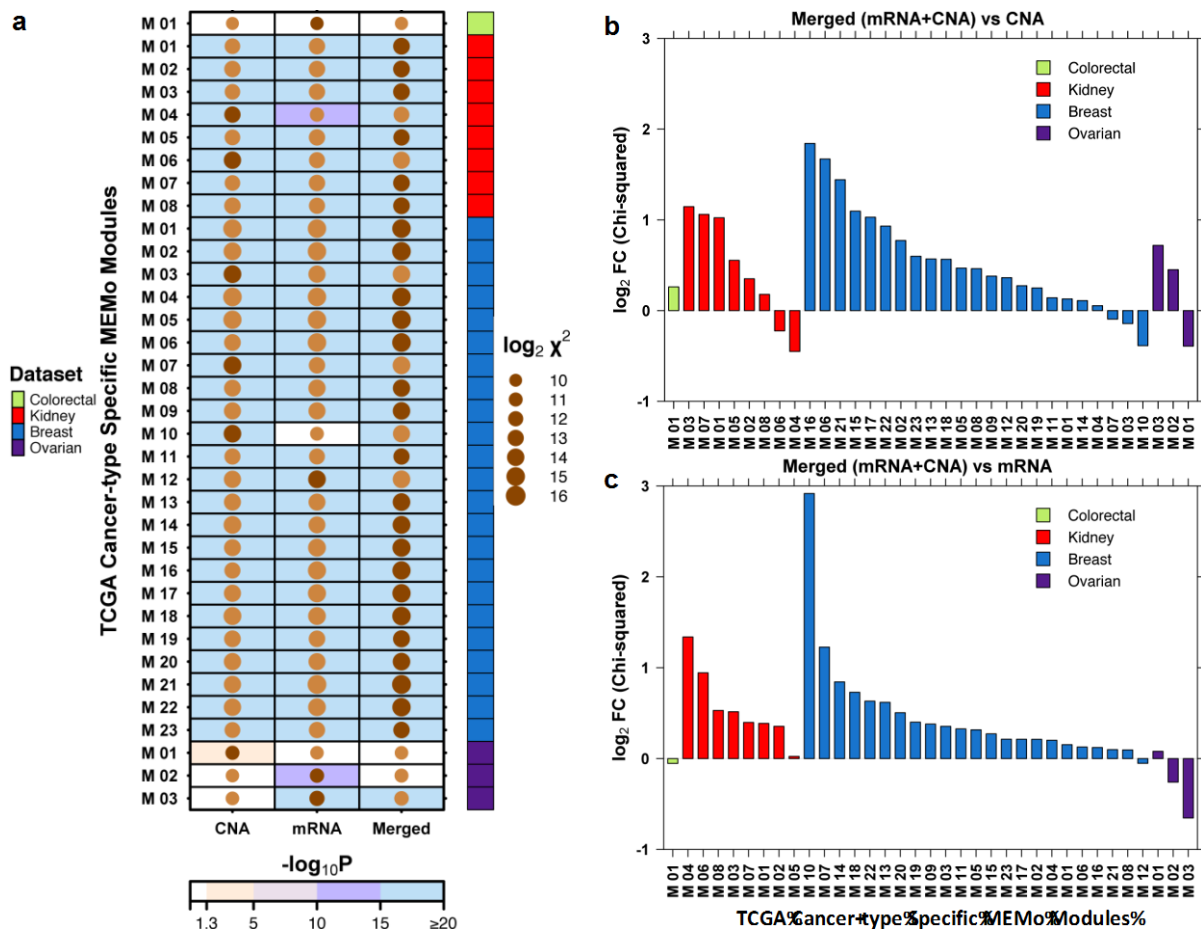
631

632 **Supplementary Figure 24**

633 **Validation of SIMMS' PI3K risk predictor.** (a) Prognostic assessment of
 634 SIMMS' PI3K risk predictor by median-dichotomizing predicted risk scores into
 635 low- and high-risk groups. (b) Prognostic assessment of model in (a) stratified by
 636 PIK3CA mutations. Patients were classified into low- and high-risk groups, and

637 each was further divided by PIK3CA mutant (+) and wild-type (-) status. **(c-d)**
638 Prognostic assessment of PI3K predictor on patients which were not treated with
639 chemotherapy and were further stratified into node -ve and node +ve groups. **(e,**
640 **f)** Prognostic performance assessment in patients with- and without
641 chemotherapy arms of the validation cohort. Within each subgroup, risk score
642 quartiles Q2-Q4 were compared against Q1 using Cox proportional hazards
643 modelling and the log-rank test. **(g)** Validation of SIMMS' PI3K risk predictor
644 (FFPE samples trained model) on ER+ subset of Metabric cohort (fresh frozen
645 samples). Risk scores of Metabric samples were dichotomised using median risk
646 score derived from TEAM cohort.

647



648

649 **Supplementary Figure 25**

650 **Multi-modal assessment of SIMMS.** Multi-modal prognostic biomarkers for
 651 breast, colon, kidney and ovarian cancers. **(a)** Dot plot of summarised (Fisher's
 652 combined probability test) chi-square estimates and P values for each of the
 653 MEMo derived cancer-type specific subnetwork modules (Mx) (**Supplementary**
 654 **Methods section 5, Supplementary Table 21**). Covariates represent colours of
 655 each cancer type. Size of the dot represents $\log(\text{chi-square})$ estimate resulting
 656 from the meta-analysis of Cox P values (1000 random subsets for each profile in
 657 each cancer type). A Cox proportional hazards model was fitted to dichotomous
 658 risk scores across the entire validation cohort to assess survival association of
 659 predicted risk groups. Crosses represent absence of a module from a particular
 660 cancer type. **(b, c)** Performance comparison of multi-modal prognostic models

661 (Merged mRNA+CNA) against CNA models (b) and mRNA models (c) in each
662 cancer type using MEMo modules of that particular cancer. Within each cancer
663 type, modules are sorted by the largest fold-change in chi-squared values; with
664 positive values indicating improved prognostication by the multi-modal model
665 over CNA or mRNA models.

666

667

668 **Supplementary References**

669

- 670 1. Breitling, R., Armengaud, P., Amtmann, A. & Herzyk, P. Rank products: a
671 simple, yet powerful, new method to detect differentially regulated genes in
672 replicated microarray experiments. *FEBS Lett* **573**, 83-92 (2004).
- 673 2. Sotiriou, C. et al. Gene expression profiling in breast cancer: understanding
674 the molecular basis of histologic grade to improve prognosis. *J Natl Cancer*
675 *Inst* **98**, 262-272 (2006).
- 676 3. Symmans, W.F. et al. Genomic index of sensitivity to endocrine therapy for
677 breast cancer. *J Clin Oncol* **28**, 4111-4119 (2010).
- 678 4. Lee, E., Chuang, H.Y., Kim, J.W., Ideker, T. & Lee, D. Inferring pathway activity
679 toward precise disease classification. *PLoS Comput Biol* **4**, e1000217 (2008).
- 680 5. Guo, Z. et al. Towards precise classification of cancers based on robust gene
681 functional expression profiles. *BMC Bioinformatics* **6**, 58 (2005).
- 682 6. Bild, A.H. et al. Oncogenic pathway signatures in human cancers as a guide to
683 targeted therapies. *Nature* **439**, 353-357 (2006).
- 684 7. Bueno, R. et al. Comprehensive genomic analysis of malignant pleural
685 mesothelioma identifies recurrent mutations, gene fusions and splicing
686 alterations. *Nat Genet* **48**, 407-416 (2016).
- 687 8. Zhao, X. et al. Systematic assessment of prognostic gene signatures for breast
688 cancer shows distinct influence of time and ER status. *BMC Cancer* **14**, 211
689 (2014).
- 690 9. Oh, S.C. et al. Prognostic gene expression signature associated with two
691 molecularly distinct subtypes of colorectal cancer. *Gut* **61**, 1291-1298 (2012).
- 692 10. Smith, J.J. et al. Experimentally derived metastasis gene expression profile
693 predicts recurrence and death in patients with colon cancer.
694 *Gastroenterology* **138**, 958-968 (2010).
- 695 11. Chen, H.Y. et al. A five-gene signature and clinical outcome in non-small-cell
696 lung cancer. *The New England journal of medicine* **356**, 11-20 (2007).
- 697 12. Lau, S.K. et al. Three-gene prognostic classifier for early-stage non small-cell
698 lung cancer. *Journal of clinical oncology : official journal of the American*
699 *Society of Clinical Oncology* **25**, 5562-5569 (2007).
- 700 13. Shedden, K. et al. Gene expression-based survival prediction in lung
701 adenocarcinoma: a multi-site, blinded validation study. *Nature medicine* **14**,
702 822-827 (2008).
- 703 14. Boutros, P.C. et al. Prognostic gene signatures for non-small-cell lung cancer.
704 *Proceedings of the National Academy of Sciences of the United States of*
705 *America* **106**, 2824-2828 (2009).
- 706 15. Starmans, M.H. et al. Exploiting the noise: improving biomarkers with
707 ensembles of data analysis methodologies. *Genome Med* **4**, 84 (2012).
- 708 16. Yoshihara, K. et al. High-risk ovarian cancer based on 126-gene expression
709 signature is uniquely characterized by downregulation of antigen

- 710 presentation pathway. *Clinical cancer research : an official journal of the*
711 *American Association for Cancer Research* **18**, 1374-1385 (2012).
- 712 17. The Cancer Genome Atlas Research Network Integrated genomic analyses of
713 ovarian carcinoma. *Nature* **474**, 609-615 (2011).
- 714 18. Mankoo, P.K., Shen, R., Schultz, N., Levine, D.A. & Sander, C. Time to
715 recurrence and survival in serous ovarian tumors predicted from integrated
716 genomic profiles. *PLoS One* **6**, e24709 (2011).
- 717 19. Wu, G. & Stein, L. A network module-based method for identifying cancer
718 prognostic signatures. *Genome biology* **13**, R112 (2012).
- 719 20. Waggott, D. et al. NanoStringNorm: an extensible R package for the pre-
720 processing of NanoString mRNA and miRNA data. *Bioinformatics* **28**, 1546-
721 1548 (2012).
- 722 21. Sabine, V.S. et al. Mutational analysis of PI3K/AKT signaling pathway in
723 tamoxifen exemestane adjuvant multinational pathology study. *J Clin Oncol*
724 **32**, 2951-2958 (2014).
- 725 22. Cuzick, J. et al. Prognostic value of a combined estrogen receptor,
726 progesterone receptor, Ki-67, and human epidermal growth factor receptor 2
727 immunohistochemical score and comparison with the Genomic Health
728 recurrence score in early breast cancer. *J Clin Oncol* **29**, 4273-4278 (2011).
- 729 23. The Cancer Genome Atlas Research Network Comprehensive genomic
730 characterization defines human glioblastoma genes and core pathways.
731 *Nature* **455**, 1061-1068 (2008).
- 732 24. Ciriello, G., Cerami, E., Aksoy, B.A., Sander, C. & Schultz, N. Using MEMo to
733 discover mutual exclusivity modules in cancer. *Curr Protoc Bioinformatics*
734 **Chapter 8**, Unit 8 17 (2013).
- 735 25. Network, T.C.G.A. Comprehensive molecular portraits of human breast
736 tumours. *Nature* **490**, 61-70 (2012).
- 737 26. Cancer Genome Atlas Research, N. Comprehensive molecular
738 characterization of clear cell renal cell carcinoma. *Nature* **499**, 43-49 (2013).
- 739 27. The Cancer Genome Atlas Research Network Comprehensive molecular
740 characterization of human colon and rectal cancer. *Nature* **487**, 330-337
741 (2012).
- 742