

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20

# **Identification of Pathways Associated with Chemosensitivity through Network Embedding**

Authors: Sheng Wang, Edward Huang, Junmei Cairns, Jian Peng, Liewei Wang, Saurabh Sinha

Running title: Network-Based Discovery of Chemosensitivity Pathways

Keywords: Network embedding, Pathway identification, Drug response, Molecular networks.

Corresponding author information:

Full name: Saurabh Sinha

Mailing address: 2122 Siebel Center, 201 N. Goodwin Ave, Urbana, IL 61801

Phone and fax numbers: 217-333-3233

Email address: [sinhas@illinois.edu](mailto:sinhas@illinois.edu)

Conflict of interest: none.

21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62

## ABSTRACT

Basal gene expression levels have been shown to be predictive of cellular response to cytotoxic treatments. However, such analyses do not fully reveal complex genotype-phenotype relationships, which are partly encoded in highly interconnected molecular networks. Biological pathways provide a complementary way of understanding drug response variation among individuals. In this study, we integrate chemosensitivity data from a recent pharmacogenomics study with basal gene expression data from the CCLE project and prior knowledge of molecular networks to identify specific pathways mediating chemical response. We first develop a computational method called PACER, which ranks pathways for enrichment in a given set of genes using a novel network embedding method. It examines known relationships among genes as encoded in a molecular network along with gene memberships of all pathways to determine a vector representation of each gene and pathway in the same low-dimensional vector space. The relevance of a pathway to the given gene set is then captured by the similarity between the pathway vector and gene vectors. To apply this approach to chemosensitivity data, we identify genes with basal expression levels in a panel of cell lines that are correlated with cytotoxic response to a compound, and then rank pathways for relevance to these response-correlated genes using PACER. Extensive evaluation of this approach on benchmarks constructed from databases of compound target genes, compound chemical structure, as well as large collections of drug response signatures demonstrates its advantages in identifying compound-pathway associations, compared to existing statistical methods of pathway enrichment analysis. The associations identified by PACER can serve as testable hypotheses about chemosensitivity pathways and help further study the mechanism of action of specific cytotoxic drugs. More broadly, PACER represents a novel technique of identifying enriched properties of any gene set of interest while also taking into account networks of known gene-gene relationships and interactions.

## 63 INTRODUCTION

64 Large-scale cancer genomics projects, such as the Cancer Genome Atlas<sup>1</sup>, the Cancer Genome  
65 project<sup>2</sup>, and the Cancer Cell Line Encyclopedia project<sup>3</sup>, and cancer pharmacology projects  
66 such as the Genomics of Drug Sensitivity in Cancer project<sup>4</sup> have generated a large volume of  
67 genomics and pharmacological profiling data. As a result, there is an unprecedented opportunity  
68 to link pharmacological and genomic data to identify therapeutic biomarkers<sup>5-7</sup>. In pursuit of this  
69 vision, significant efforts have been invested in identifying the genetic basis of drug response  
70 variation among individual patients. For instance, a recent study performed a comprehensive  
71 survey of genes with basal expression levels in cancer cell lines that correlate with drug  
72 sensitivity, revealing potential gene candidates for explaining mechanisms of action of various  
73 drugs<sup>8</sup>.

74  
75 While significant efforts have focused on specific genes that interact with compounds and confer  
76 observed cellular phenotypes, there has been relatively little progress in studying the synergistic  
77 effects of genes. These effects are key factors in comprehensively deciphering the mechanisms  
78 of action of compounds and understanding complex phenotypes<sup>9</sup>. Similarly, pathways, which  
79 comprise a set of interacting genes, have emerged as a useful construct for gaining insights into  
80 cellular responses to compounds. Analysis at the pathway level not only reduces the analytic  
81 complexity from tens of thousands of genes to just hundreds of pathways, but also contains  
82 more explanatory power than a simple list of differentially expressed genes<sup>10</sup>. Consequently, an  
83 important yet unsolved problem is the effective identification of pathways mediating drug  
84 response variation. Although the associated pathways for certain drugs have been studied  
85 experimentally<sup>11-13</sup>, *in vitro* pathway analysis is costly and inherently difficult, making it hard to  
86 scale to hundreds of compounds.

87  
88 Fortunately, a growing compendium of genomic, proteomic, and pharmacologic data allows us  
89 to develop scalable computational approaches to help solve this problem. Although statistical  
90 significance tests and enrichment analyses can be naturally applied to compound-pathway  
91 association identification (e.g., by testing the overlap between pathway members and  
92 differentially expressed genes), these approaches fail to leverage well-established biological  
93 relationships among genes<sup>14-17</sup>. Even when analyzing individual genes, molecular networks  
94 such as protein-protein interaction networks have been shown to play crucial roles in  
95 understanding the cellular drug response<sup>9, 18-21</sup>. Therefore, we propose to combine similar  
96 molecular networks with gene expression and drug response data for pathway identification.  
97 However, integrating these heterogeneous data sources is statistically challenging. Moreover,  
98 networks are high-dimensional, incomplete, and noisy. Thus, our algorithm needs to accurately  
99 and comprehensively identify pathway while exploiting suboptimal networks.

100  
101 In this work, we present PACER, a novel, network-assisted algorithm that identifies pathway  
102 associations for any gene set of interest. Additionally, we apply the algorithm to discover  
103 chemosensitivity-related pathways. PACER first constructs a heterogeneous network that  
104 includes pathways and genes, pathway membership information, and gene-gene relationships  
105 from a molecular network such as protein-protein interaction network. It then applies a novel  
106 dimensionality reduction algorithm to this heterogeneous network to obtain compact, low-

107 dimensional vectors for pathways and genes in the network. Pathways that are topologically  
108 close to drug response-related genes in the network are co-localized with those genes in this  
109 low-dimensional vector space. Hence, PACER ranks each pathway based on its proximity in the  
110 low-dimensional space to genes that have basal expressions highly correlated with drug  
111 response. We evaluated PACER's ability to identify compound-pathway associations with three  
112 'ground truth' sets built from compound target data<sup>8</sup>, compound structure data<sup>22</sup>, and LINCS  
113 differential expression data<sup>23</sup>. When comparing PACER to state-of-the-art methods that ignore  
114 prior knowledge of interactions among genes, we observed substantial improvement of the  
115 concordance with the chosen benchmarks. Even though we developed PACER and tested its  
116 ability to identify compound-pathway associations, the algorithm is applicable to any scenario in  
117 which one seeks to discover pathways related to a pre-specified gene set of interest, while  
118 utilizing a given gene network.

119  
120  
121  
122  
123  
124  
125  
126  
127  
128  
129  
130  
131  
132  
133  
134  
135  
136  
137  
138  
139  
140  
141  
142  
143  
144  
145  
146

## 147 MATERIALS AND METHODS

### 148 Compound response data and gene expression data

149 We obtained a large-scale compound response screening dataset from Rees *et al.*<sup>8</sup>, which  
150 spans 481 chemical compounds and 842 human cancer cell lines encompassing 25 lineages.  
151 These 481 compounds were collected from different sources including clinical candidates, FDA-  
152 approved drugs and previous chemosensitivity profiling experiments. Area under the drug  
153 response curve (AUC) was used by the authors of that study to measure cellular response to  
154 individual compound. We also obtained gene expression profiles for these cell lines from the  
155 Cancer Cell Line Encyclopedia (CCLE) project<sup>3</sup>, profiled using the GeneChip Human Genome  
156 U133 Plus 2.0 Array. Since these expression measurements were done in each cell line without  
157 any drug treatment, they are referred to as ‘basal’ expression levels. In contrast, the expression  
158 profiling of a cell line was performed *after* treatment with a drug in certain studies<sup>23</sup>. We  
159 obtained the SMILE specification of each drug from PubChem<sup>22</sup> and then calculated the  
160 Tanimoto similarity scores between all pairs of drugs based on their SMILE specifications.

### 161 STRING-based molecular network and NCI pathway collection

162 We obtained a collection of six human molecular networks from the STRING database v9.1<sup>24</sup>.  
163 These six networks include experimentally derived protein-protein interactions, manually  
164 curated protein-protein interactions, protein-protein interactions transferred from model  
165 organism based on orthology, and interactions computed from genomic features such as fusion-  
166 fusion events, functional similarity and co-expression data. There are 16,662 genes in the  
167 network. We used all the STRING channels except “text-mining” and used the Bayesian  
168 integration method provided by STRING. Since our approach can deal with different edge  
169 weights, we did not set a threshold to remove low confidence edges. We referred to this  
170 integrated network as the ‘STRING-based molecular network’. To test whether genes that are  
171 highly correlated with many compounds tend to have higher degrees in the network, we formed  
172 two groups of genes. One group contained genes that are correlated with over 100 compounds,  
173 and the other group contained the remaining genes. We then used the Wilcoxon signed-rank  
174 test to test whether the degrees of genes in these two groups were from the same distribution.  
175 We obtained a collection of 223 cancer-related pathways from the National Cancer Institute  
176 (NCI) pathway database. These manually curated pathways include human signaling and  
177 regulatory pathways as well as key cellular processes<sup>25</sup>.

### 178 179 The PACER Framework

180 PACER integrates pathway information with the STRING-based molecular network described  
181 above by constructing a heterogeneous network of genes and pathways. An edge exists  
182 between two genes if they are connected in the network. An edge exists between a pathway  
183 and a gene if the gene belongs to the pathway. There are no direct pathway-pathway edges in  
184 the heterogeneous network. PACER adopts diffusion component analysis (DCA), a recently  
185 developed network representation algorithm to learn a low-dimensional vector for each node in  
186 the network<sup>26</sup>. Because of its ability to handle noisy and missing edges in the biological network,  
187 DCA has achieved state-of-the-art results in different tasks<sup>26, 27</sup>. Since compounds are not

188 nodes in the constructed heterogeneous network, only genes and pathways are projected onto  
189 the low-dimensional space. After learning the low-dimensional representations of all nodes  
190 (genes and pathways), DCA ranks pathways based on the weighted cosine similarities between  
191 a pathway and the set of 250 genes most correlated with response to a compound. We  
192 henceforth refer to this set of genes as "response-correlated genes" (RCG) for the compound.  
193 These genes' expression values are most significantly correlated with chemosensitivity. We  
194 found that the performance of the PACER method is stable for different choices (200, 250, and  
195 300) for the number of RCGs considered in this step (**Suppl. Figure 7-9.**).

## 196 **LINCS drug perturbation profiles**

197 LINCS is a data repository of over 1.3 million genome-wide expression profiles of human cell  
198 lines subjected to a variety of perturbation conditions, which include treatments with more than  
199 20 thousand unique compounds at various concentrations. Each perturbation experiment is  
200 represented by a list of differentially expressed genes that are ranked based on z-scores of  
201 perturbation expression relative to basal expression. For each gene, we first took the difference  
202 between its expression in a perturbation condition and its expression in a control condition (i.e.,  
203 treatment with pure DMSO solvent). We then considered the differential expression of the gene  
204 in multiple perturbation experiments involving that compound (i.e., different concentrations, time  
205 points, and cell lines). We used the maximum differential expression to represent the  
206 compound's effect on that gene's expression. All genes were then ranked by their differential  
207 expression on treatment with the compound, and the top 250 genes were treated as  
208 differentially expressed genes (DEGs) of the compound, provided their z-score has an absolute  
209 value greater than 2.

210

## 211 **Comparison with method of Huang *et al.***

212 We implemented the method of Huang *et al.*<sup>16</sup> ourselves using the exact same input (i.e.,  
213 chemosensitivity and gene expression data) as PACER. We first computed a gene's correlation  
214 to a drug by calculating the Pearson correlation coefficient between the gene's expression  
215 values and the drug response values across cell lines. Let the set of genes in pathway  $p$  be  
216 denoted by  $G_p$ , and their correlation values to a drug  $d$  by  $C(G_p, d)$ . Conversely, the set of genes  
217 not in pathway  $p$  is denoted as  $\overline{G_p}$ , and their correlation values to  $d$  as  $C(\overline{G_p}, d)$ . We then  
218 performed the Kruskal-Wallis H-test, following Huang *et al.*, to test if the medians of  $C(G_p, d)$   
219 and  $C(\overline{G_p}, d)$  were significantly different. We used the resulting  $p$ -value to rank pathways for  
220 each drug.

221

## 222 **Software availability:**

223 The PACER software is currently available at <https://github.com/KnowEnG/PACER>. It will also  
224 be available through the cloud-based analysis framework KnowEnG ([knoweng.org](http://knoweng.org)) upon  
225 publication.

226

## 227 RESULTS

### 228 Global analysis of correlations between basal gene expression and compound response

229 Following the work of Rees *et al.*<sup>8</sup>, we first examined correlations between the compound  
230 sensitivity and basal gene expression profiles across hundreds of cell lines. We calculated  
231 Pearson correlation coefficients between each gene's expression and the cellular response  
232 (measured as the area under the curve or AUC) to each compound, across different cell lines  
233 (**Figure 1A**). Compared to the IC50 and EC50 scores, AUC simultaneously captures the  
234 efficacy and potency of a drug. Of the ~8.7 million pairs of genes and compounds tested, we  
235 found 294,789 to be significantly correlated ( $p$ -value < 0.0001 after Bonferroni correction,  
236 corresponding to a Pearson correlation coefficient of 0.215). Within these significantly correlated  
237 pairs, 1,749 genes were correlated with over 100 compounds (**Figure 1B, Suppl. Table 1**). We  
238 note that these key genes tend to be high-degree nodes in STRING-based molecular network  
239 (Wilcoxon rank sum test  $p$ -value < 9.6e-14, see Methods). We also found that some (10 of 481)  
240 compounds were significantly correlated (Pearson correlation  $p$ -value < 0.0001 after Bonferroni  
241 correction) with more than 3,200 genes (**Figure 1C**). Five of these ten compounds are  
242 chemotherapeutic agents (**Suppl. Table 2**). In contrast, about 100 compounds were not  
243 significantly correlated with any genes; these compounds are mostly probes that either lack  
244 FDA approval or are not clinically used. The large disparity among the examined compounds in  
245 terms of the number of correlated genes reflects the diversity of these 481 small molecules.  
246 While many of them are chemotherapeutic, which can affect the expression of a large number of  
247 genes, some compounds may be targeting specific mutations, post-translational modifications,  
248 or protein expression. A closer examination revealed that the compounds with the highest  
249 cytotoxicity had the fewest gene correlations (i.e., fewest genes whose expression correlates  
250 with cytotoxic response were mostly those with low cytotoxicity) (**Suppl. Figure 1**). This  
251 suggests that the strategy of identifying compound-associated genes by correlating basal gene  
252 expression profiles with cytotoxicity is likely to be more effective for more potent compounds, for  
253 which average response is stronger. Note that the gene expression profiles used here are basal  
254 and not in response to treatment with compound, hence it was not clear *a priori* that more  
255 effective compounds would have larger numbers of gene correlates. In summary, examination  
256 of individual genes' correlations with chemical response confirmed previous reports<sup>4, 8, 28</sup> that  
257 basal gene expression significantly correlates with cytotoxicity across cell lines, especially for  
258 effective cytotoxic drugs.

259

### 260 Identifying compound-specific pathways via enrichment tests

261 The above evidence for correlations between basal gene expression and chemical response  
262 raised the possibility that one might discover important biological pathways associated with the  
263 response by a systems-level analysis of gene expression data. To explore this, we considered a  
264 collection of 223 cancer-related pathways from the National Cancer Institute (NCI) pathway  
265 database<sup>25</sup> and used Fisher's exact test to quantify the overlap between the set of genes in a  
266 pathway  $p$  and (RCGs). A significantly large overlap between the two sets indicates an  
267 association between the pathway and the compound. We performed a multiple hypothesis

268 correction on all pathway association tests for each compound, using  $FDR \leq 0.05$ . The results  
269 of this baseline method for predicting pathway associations are shown in **Figure 1D** (distribution  
270 of the number of compounds that are significantly associated with each pathway) and **Figure 1E**  
271 (distribution of the number of pathways significantly associated with each compound). Both  
272 distributions revealed a long tail. For instance, while each pathway was associated with an  
273 average of 18 compounds (of the 481 tested), there were 10 pathways that were associated  
274 with over 150 compounds (**Suppl. Table 3**). Likewise, while each compound was associated  
275 with an average of eight pathways, there were 12 compounds associated with over 25 pathways  
276 (**Suppl. Table 4**). We show the details of these long tails in **Suppl. Figure 2**.

## 277 **A new method for identifying pathways associated with chemical response, based on** 278 **network embedding**

279 We observed above that key RCGs – those correlated with many compounds – tend to be  
280 enriched in high degree nodes of the STRING-based molecular network. This suggests that an  
281 analysis combining this network with pathway enrichment tests might provide additional insights.  
282 We therefore developed a novel network-based method, called PACER, for scoring compound-  
283 pathway associations. PACER (**Figure 2A**) first constructs a heterogeneous network consisting  
284 of genes and pathways as nodes. In this network, gene-pathway edges denote pathway  
285 memberships based on a compendium of pathways and gene-gene edges from the STRING-  
286 based molecular network introduced above (also see Methods). PACER then creates a low-  
287 dimensional vector representation for each gene and pathway node in the heterogeneous  
288 network, reflecting the node's position in this heterogeneous network. This is done by the  
289 Diffusion Component Analysis (DCA) approach reported in previous work<sup>26, 27</sup>. Nodes (i.e.,  
290 pathways or genes) will have similar vector representations if they are near each other in the  
291 network. For instance, two pathway nodes will have similar vector representations if the  
292 pathways share genes and/or their genes are related by the STRING-based molecular network.  
293 In a similar vein, two genes will have similar representations if they belong to the same  
294 pathway(s) and/or exhibit the same network neighbors. A gene and a pathway can also be  
295 compared in the low-dimensional space, and will be deemed similar if the gene is in the  
296 pathway and/or the gene is related by network to other genes of the pathway. DCA performs  
297 network-based embedding without utilizing gene expression or chemical response data. Next,  
298 PACER identifies RCGs as a fixed number of genes the expression of which shows the greatest  
299 correlation with chemical response to a specific compound. Finally, it scores a pathway based  
300 on the average cosine similarity between the vector representation of the pathway and those of  
301 the RCGs. A pathway can thus be found to be associated with a compound if, in the network,  
302 the pathway genes are closely related to the compound's RCGs; this association can be  
303 discovered even if the pathway does not actually include the RCGs. We note that scores  
304 assigned by PACER are not statistical significance scores and are meant only to rank pathways  
305 for association with a given compound. Also, a negative score assigned to a compound-  
306 pathway pair does not imply a negative correlation between expression levels of pathway genes  
307 and chemosensitivity. Rather, it only implies a lack of evidence for an association between the  
308 compound-pathway pair.

309  
310 The PACER association scores for all combinations of 481 compounds and 223 NCI signaling

311 pathways are shown in **Figure 2B**. The pathways cluster into many distinct groups, each with  
312 different compound association profiles. One group (cyan branches in the row dendrogram),  
313 associated with more than half of the compounds, consists of pathways describing various  
314 integrin cell surface interactions (e.g., ‘Integrins in angiogenesis’ pathway, ‘Alpha 4 Beta1  
315 integrin cell surface interactions’ pathway). These pathways are known to play crucial roles in  
316 communications among cells in response to small molecules<sup>29</sup>. Notably, integrins are a major  
317 family of cell surface adhesion receptors, and are involved in major pathways that contribute to  
318 cancer cell survival and resistance to chemotherapy<sup>30</sup>. PACER found 329 of the 481  
319 compounds to be associated with the “integrin family cell surface interactions” pathway. Since  
320 PACER scores are not easily assigned statistical significance levels, we chose for each  
321 compound  $n$  pathways with the highest PACER scores, where  $n$  is the number of statistically  
322 significant pathway associations ( $FDR \leq 0.05$ ) found by the baseline method above for the same  
323 compound. We found literature support for some of these associations. For example, ruxolitinib,  
324 a JAK/STAT inhibitor, is associated with integrin pathways by PACER analysis. In a previous  
325 study, it was shown that beta 4 integrin enhances activation of the transcription factor STAT3,  
326 which is a target of ruxolitinib<sup>31</sup>. Furthermore, vorinostat, a member of a larger class of  
327 compounds that inhibit histone deacetylases (HDAC), can induce integrin  $\alpha 5\beta 1$  expression and  
328 activate MET, leading to resistance<sup>32</sup>. **Figure 2C** reveals another example of functionally related  
329 pathways being grouped together. The pathways ‘VEGFR1 specific signals’, ‘ErbB4 signaling  
330 events’, ‘EGFR-dependent Endothelin signaling events’, ‘ErbB receptor signaling network’, and  
331 ‘PDGF receptor signaling network’ form one group (red branches in row dendrogram), and are  
332 associated with masitinib (PDGFRB inhibitor) and RAF265 (VEGFR2 inhibitor), among other  
333 compounds (see **Suppl. Table 5**). Vascular endothelial growth factor receptor (VEGFR inhibitor)  
334 and platelet-derived growth factor receptor (PDGFR inhibitor) are both members of the family of  
335 58 known tyrosine kinase receptors in humans<sup>33</sup>. Tyrosine kinases have various modulatory  
336 functions in growth factor signaling and several of their inhibitors are known for their anti-tumor  
337 activity<sup>34</sup>.

338  
339 **Figure 2B** also shows compounds clustered into different groups based on their associations  
340 with pathways. We found that many compounds with similar structure were grouped together.  
341 For example, teniposide and etoposide had a Tanimoto similarity score of 0.94 between their  
342 SMILE specifications, which was substantially higher than the average Tanimoto similarity score  
343 of 0.3716 for all pairs of drugs. They were clustered together in the same group (**Figure 2D**,  
344 also marked as a rectangle in **Figure 2B**), which had seven compounds. This group is  
345 associated with a set of similar pathways, including ‘p53 pathway’, ‘Direct p53 effectors’,  
346 ‘Signaling mediated by p38-alpha and p38-beta’, and ‘Signaling mediated by p38-gamma and  
347 p38-delta’. We found support in the literature in favor of some of these associations. For  
348 example, a previous study reported that etoposide activates p38MAPK and can be used as a  
349 new combined treatment approach when used with p38MAPK inhibitor SB203580<sup>35</sup>. To take  
350 another example, temsirolimus and tacrolimus, which are both epipodophyllotoxins and inhibit  
351 topoisomerase II, have a Tanimoto similarity score of 0.82, and are grouped closely in **Figure**  
352 **2B**.

353

354 **PACER improves pathway identification**

355  
356 We noted a substantial degree of complementarity between the top predictions of PACER and  
357 those of the baseline method that uses Fisher's exact test between RCGs and pathway genes  
358 (see **Suppl. Table 5**). For instance, PACER found that PD153035, an ErbB2 inhibitor, is  
359 associated with the 'C-MYC pathway', reflecting the fact that PD153035 is able to reduce c-Myc  
360 protein levels in breast tumor cells<sup>36</sup>. The baseline approach did not find this association to be  
361 significant. Similarly, PACER reported that the 'EGFR-dependent Endothelin signaling events'  
362 pathway is associated with EGFR inhibitor gefitinib<sup>37</sup>, while the baseline method did not.

363  
364 For a more systematic comparison between the two methods, we evaluated PACER based on a  
365 database of known compound targets. We performed the evaluation under the assumption that  
366 a pathway containing at least one known target is an associated pathway. Huang *et al.* used  
367 and suggested this approach<sup>16</sup>. We used it here to evaluate PACER, the baseline method, as  
368 well as a third method presented by Huang *et al.*<sup>16</sup> Although this third method was proposed to  
369 detect association between pathways and drug clades, it can directly detect pathway-compound  
370 associations. We implemented the method ourselves (see Methods) and included it in our  
371 evaluations. We obtained the known targets for 246 compounds in our compound set from Rees  
372 *et al.*<sup>8</sup> We then computed the AUROC of pathway predictions made by PACER for each  
373 compound, and plotted this information alongside analogous information for the baseline  
374 method and the method of Huang *et al.*<sup>16</sup> As shown in **Figure 2E**, PACER identified pathways  
375 with higher AUROC compared to the other two methods. For example, PACER identified  
376 pathways with an AUROC greater than 0.75 for 22 different compounds, while the baseline  
377 method achieved this level of AUROC for only 7 compounds. **Table 1** shows the 10 compounds  
378 for which PACER achieved highest AUROC. We present a closer examination of PACER  
379 predictions for two of these compounds: 'cyclophosphamide' (0.80 AUROC) and 'nsc23766'  
380 (0.84 AUROC) in **Table 2** and **Table 3**, respectively. These tables also show the  
381 complementarity between PACER predictions and those of the baseline method.

382  
383 We note that the AUROC values reported here are likely to be underestimates, as there is  
384 literature evidence for some of the reported pathways being associated with the compound,  
385 even though the pathway does not include a known target (and is thus considered a false  
386 positive in our AUROC estimate). For example, our method identified the 'Fanconi anemia  
387 pathway' as being associated with compound 'cyclophosphamide' (**Table 2**). The Fanconi  
388 anemia pathway is one of the major DNA damage response pathways disrupted in breast  
389 cancer<sup>38</sup>, and is known to play an important role in cancer treatment by DNA crosslinking agents  
390 such as cyclophosphamide<sup>39, 40</sup>. PACER also identified 'EPO signaling pathway' to be  
391 associated with cyclophosphamide. Cyclophosphamide treatment has been reported to affect  
392 EPO receptor expression in murine erythropoiesis<sup>41</sup>. Although we found no existing study to  
393 corroborate the PACER-predicted association between cyclophosphamide and the 'ATR  
394 signaling pathway', the cyclophosphamide analogue, mafosfamide, has been reported to  
395 activate the ATM/ATR-Chk1/Chk2 pathway<sup>42</sup>. In addition, the predicted association of  
396 cyclophosphamide with 'Class I PI3K signaling events' is supported by reports of the compound  
397 activating the PI3K/Akt/mTOR signaling pathway in the ovary<sup>43</sup>.

398

399 The Rac1-specific inhibitor nsc23766 was also identified by PACER as being associated with  
400 several pathways that include a known target (**Table 3**), as well as some pathways that do not  
401 but whose association is supported by literature evidence. For example, the pathway 'Beta5  
402 beta6 beta7 and beta8 integrin cell surface interactions' was predicted as being associated with  
403 this compound, and even though the latter's known target is not in this pathway, various lines of  
404 evidence support the association. A study of cholangiocarcinoma found beta6 integrin to  
405 promote invasiveness by activating Rac1, and that the compound nsc23766 is able to suppress  
406 this invasiveness<sup>44</sup> and can be used to identify poor prognostic HER2 amplified breast cancer  
407 patients<sup>45</sup>. Beta8 integrin influences Rac1 levels to promote cell invasiveness in glioblastoma<sup>46</sup>.  
408 Also, Beta8 integrin is reported to activate Rac1 signaling in endometrial epithelial cells<sup>47</sup>. In  
409 addition, ncs23766 was predicted to be associated with 'a4b7 Integrin signaling'. While this  
410 pathway does not directly include the compound's target, Rac1 was previously shown to induce  
411 a4b7-mediated T cell adhesion to MAdCAM-1<sup>48</sup>, supporting the predicted association.  
412 Furthermore, 'Plexin-D1 Signaling pathway' was found by PACER to be associated with  
413 ncs23766. Plexin-D1 is a receptor for SEMA4A which inhibits Rac activation<sup>49</sup>. Thus, the  
414 anecdotal observations made above suggest that compound-pathway predictions made by  
415 PACER may sometimes be worth pursuing even if the compound's target is not included in the  
416 pathway.

417  
418 We further evaluated PACER based on the identified pathways for compounds with similar  
419 chemical structures. We performed this evaluation under the assumption that compounds with  
420 similar chemical structures tend to be associated with the same pathways. To this end, we  
421 calculated the Tanimoto similarity between each pair of compounds according to their SMILE  
422 specifications. We regarded two compounds as having similar chemical structures if their  
423 Tanimoto similarity is larger than 0.8. This threshold was used in previous work to indicate high  
424 Tanimoto similarity<sup>50</sup>. For a given compound, we used PACER to rank other compounds  
425 according to the Spearman correlation coefficient between their rankings of pathways, and  
426 asked if this ranking was predictive of chemical structure similarity. For each compound, an  
427 AUROC score was computed to measure this predictive ability. We repeated this evaluation  
428 with the three different methods of ranking pathways for a compound: PACER, the baseline  
429 method, and the method of Huang *et al.* **Figure 2F** shows the number of compounds for which  
430 the AUROC is above a specified threshold, for each of the three methods. We found that  
431 PACER achieves better performance in identifying compounds with similar chemical structure  
432 based on similarity of pathway ranking. For example, 21 (of 42) compounds yielded an AUROC  
433 greater than 0.8 when using PACER, compared to 10 compounds meeting the same criterion  
434 when using the baseline method. As compounds with similar chemical structure tend to be  
435 functionally similar, our results demonstrate that PACER can be used to identify similar  
436 compounds by integrating prior network information into chemosensitivity data.

437  
438 We also compared the associations predicted by the three methods to those identified from an  
439 external data set. To this end, we mined the Library of Integrated Network-Based Cellular  
440 Signatures (LINCS) L1000 data<sup>23</sup>, which reports genes differentially expressed upon treatment  
441 of various cell lines with a compound. For each compound in our analysis that is also included  
442 as a pertubagen in the L1000 compendium, we established a LINCS-based benchmark of

443 significantly associated pathways. This was based on a Fisher's exact test ( $p$ -value  $\leq 0.05$ )  
444 between pathway genes and the most differentially expressed genes from treatments with the  
445 same compound (see Methods). We required this criterion to be met in at least one of the cell  
446 lines for which data was available from LINCS. We then assessed the concordance between  
447 this set of LINCS-based compound-pathway associations and those predicted by either method  
448 presented above. We recognize that this is not an ideal benchmark: LINCS data points to genes  
449 (and, indirectly, to pathways) that are differentially expressed in response to treatment, while  
450 PACER and the compared methods base their pathway predictions on genes that have basal  
451 expression levels across cell lines that correlate with chemical response. At the same time, we  
452 expect the pathways affected by chemical treatment to also be, to an extent, involved in  
453 interpersonal variation of chemosensitivity, making this a suitable evaluation procedure. This  
454 was inspired by similar observations in cancer biology: genes and pathways disrupted in cancer  
455 tissues overlap with genes and pathways whose mutation status in germline non-tumor samples  
456 is informative about disease susceptibility and progression.

457  
458 To test whether the significant pathways identified from LINCS data agree with the pathways  
459 predicted by one of the methods being evaluated (based on chemical response variation in  
460 CCLE cell lines), we counted the compounds for which the two sets of predicted pathways  
461 overlapped significantly (Fisher's exact test  $p$ -value  $\leq 0.05$ ). As shown in **Figure 2G**, the  
462 PACER approach predicts pathways concordant with the corresponding LINCS-based  
463 benchmark for more compounds, compared to the baseline method and that of Huang *et al.*<sup>16</sup>  
464 For instance, when the baseline method used an FDR threshold of 10% to designate significant  
465 pathway associations for each drug, and the PACER method predicted the same number of  
466 pathways, the latter's predictions were concordant with the LINCS-based benchmark for 110 of  
467 the 481 compounds, a nearly two-fold improvement over the baseline method's predictions. Our  
468 evaluations actually provide evidence for the above-mentioned possibility that pathways  
469 predictive of drug sensitivity overlap with genes that mediate drug response. In fact, we found  
470 86 compounds for which the pathways identified from basal expression correlations and the  
471 pathways identified from LINCS signatures overlap with FDR  $< 5\%$ .

472  
473 After observing the substantial improvement of PACER, we then investigated whether the  
474 performance of PACER is stable when only using experimental derived protein-protein  
475 interactions as input. We found that the performance of PACER, as per the three evaluations  
476 presented above, was stable when only using experimental derived protein-protein interactions  
477 as input (**Suppl. Figure 4-6**). We further demonstrated that the result of our method is robust to  
478 different numbers of top response-correlated genes used in PACER, as shown in **Suppl.**  
479 **Figures 7-9**. We compared different values for  $k$  in the top  $k$  genes chosen by PACER. We  
480 found that for two of the three evaluation schemes, results were comparable when using  $k=100$ ,  
481 150, 200, 250 and 300. For the third evaluation scheme, results were comparable when using  
482  $k=200$ , 250, and 300. This demonstrates the stability of the algorithm's performance to different  
483 but reasonable values of  $k$  in its choice of top  $k$  response-correlated genes.

484  
485  
486

## 487 DISCUSSION AND CONCLUSION

488 We have shown that embedding prior knowledge in a gene network can more accurately identify  
489 compound-pathway. Our new method, called PACER, identified many compound-pathway  
490 associations that are supported by known compound targets as well as literature evidence. Due  
491 to its unique ability to incorporate any suitable compendium of gene interactions, our approach  
492 may provide complementary insights into drug mechanisms of action.

493  
494 Historically, pathways associated with a particular gene set are identified by using popular  
495 statistical methods such as Gene Set Enrichment Analysis<sup>51</sup>, Fisher's exact test (DAVID<sup>52</sup>) or  
496 the Binomial test (Reactome<sup>53</sup>). These tools test the overlap between differentially expressed  
497 genes and pathway members. They may also be applied to the set of drug-response-correlated  
498 genes (RCGs) analyzed here. Ingenuity Pathway Analysis<sup>54</sup> is another related tool, which  
499 utilizes information about causal interactions between pathway members. Our study is similar to  
500 the above tools in that PACER also seeks to find pathways implicated by a gene set. However,  
501 our approach differs from these existing tools in that known molecular interactions (e.g., PPI)  
502 among different genes are taken into consideration. Thus, a gene set, be it the RCGs of a  
503 compound or the members of a pathway, is not treated merely as the sum of its parts, but also  
504 includes the relationships among those parts. Since the dominant theme in existing approaches  
505 is assessment of overlaps between two gene sets (MSigDB, DAVID, and Reactome adopt  
506 variations on this theme), our extensive comparisons between PACER and the baseline method  
507 of Fisher's exact test shed light on the relative merits of the new approach. A related line of work  
508 aims to identify differentially expressed subnetworks in a given interaction network, e.g.,  
509 KeyPathwayMiner<sup>55</sup>, but these studies are only superficially relevant to our work since we aim to  
510 prioritize existing pathways instead of finding new pathways.

511  
512 We consider two potential reasons for the strong performance of PACER. First, it is widely  
513 appreciated that a chemical compound not only affects individual genes, but also combinations  
514 of genes in molecular networks corresponding to core processes, such as cell proliferation and  
515 apoptosis. Our method postulates that even if the RCGs and a pathway may only have a few  
516 genes in common, they may be close to each other in the network. Although current compound  
517 pathway maps are incomplete, much relevant information is available in public databases of  
518 human molecular networks. While traditional pathway enrichment analysis methods like Fisher's  
519 exact test identify pathways according to the number of shared genes, PACER prioritizes  
520 pathways based on their proximities to RCGs in molecular networks. Second, manually curated  
521 pathways may have arbitrary boundaries due to the need to capture knowledge at different  
522 levels of detail. Consequently, identifying drug related pathways might be hindered by pathway  
523 boundaries. By leveraging the prior knowledge in molecular networks, PACER is more robust to  
524 pathways with different boundaries, thus improving the sensitivity of detecting compound-  
525 pathway associations.

526  
527 We see many opportunities to improve upon the basic concept of PACER in future work. First,  
528 although the current PACER framework was developed in an unsupervised fashion, the scores  
529 assigned to each pathway for the given gene set can be used as the feature and plugged into  
530 off-the-shelf machine learning classifiers for compound-pathway association identification.

531 Second, although this study focused on chemosensitivity response, the PACER method is  
532 broadly applicable to testing the association between two sets of genes according to their  
533 proximity in the network. Finally, although we use gene expression data as the molecular profile  
534 of each cell line, it might be interesting to test our method based on other molecular data such  
535 as somatic mutations and copy number alterations.

536

537

538

539 **FIGURE LEGENDS**

540

541 **Figure 1. Global analysis of correlations between basal gene expression and compound**  
542 **response.** (A) Heatmap of the Pearson correlation coefficient between genes (expression) and  
543 compounds (chemosensitivity, measured by AUC values). (B) Histogram of the number of  
544 compounds associated with each gene. The *y*-axis shows the number of genes associated with  
545 *k* compounds, where *k* is shown on the *x*-axis. (C) Histogram of the number of genes associated  
546 with each compound. The *y*-axis shows the number of compounds associated with *k* genes,  
547 where *k* is shown on the *x*-axis. (D) Histogram of the number of compounds significantly  
548 associated with each pathway (Fisher's exact test  $FDR \leq 0.05$ ). (E) Histogram of the number of  
549 pathways significantly associated with each compound (Fisher's exact test  $FDR \leq 0.05$ ).

550

551 **Figure 2. Identifying pathways associated with chemical response using PACER.**

552 (A) Schematic description of PACER. (B) Heat map of associations between compounds and  
553 pathways (PACER scores). Rows are pathways and columns are compounds. (C) Detailed view  
554 of a subset of the red branch cluster of pathways, marked by 'C' in Figure 2B. (D) Detailed view  
555 of a subnet of purple branch cluster of compounds, marked by a rectangle in Figure 2B. (E)  
556 Comparative evaluation of different methods for predicting compound-pathway associations.  
557 The ground truth used here is the pathways that contain any known target gene of the  
558 compound. (F) Comparison of PACER, Fisher's exact test and Huang *et al.* on predicting  
559 compounds with similar chemical structure. The *y*-axis shows the number of compounds with an  
560 AUROC larger than *k*, where *k* is shown on the *x*-axis. Compound structure similarity is  
561 determined by the Tanimoto similarity score calculated based on their SMILE specifications.  
562 Prediction is made according to the Spearman correlation between the pathway rankings of two  
563 compounds. (G) Number of compounds with significant overlap ( $p < 0.05$ ) between pathways  
564 from LINCS and pathways from PACER, from Huang *et al.* 2005 and from the baseline method  
565 (Fisher's exact test) respectively, at different levels of stringency in pathway prediction.  
566 Stringency refers to the FDR control used by the baseline method in determining significant  
567 pathways. Both PACER and the Huang *et al.* 2005 method were used to predict the same  
568 number of (highest scoring) pathways as the baseline method, for a fair comparison.

569 **TABLES**

570 **Table 1.** Compounds for which PACER predicted pathways with greatest precision. Evaluation was  
571 performed with known targets.

Compound	AUROC
bms-536924	0.936652
ml239	0.869446
kx2-391	0.846667
skepinone-l	0.844854
mgcd-265	0.841279
nsc23766	0.835607
vorapaxar	0.834101
pf-3758309	0.824359
cyclophosphamide	0.798150
pf-573228	0.792370

572  
573  
574

575 **Table 2.** Top 10 pathways predicted by PACER for the compound 'cyclophosphamide'

576

Pathway	Contains known target?	P-value of pathway (baseline method)
Fanconi anemia pathway	NO	0.2670
CXCR4-mediated signaling events	YES	1
TCR signaling in naive CD4+ T cells	YES	1
IL2 signaling events mediated by PI3K	YES	1
ATR signaling pathway	NO	1
TCR signaling in naive CD8+ T cells	YES	1
Class I PI3K signaling events	NO	1
p53 pathway	YES	0.3185
BARD1 signaling events	NO	1
EPO signaling pathway	NO	1

577  
578

579

580 **Table 3.** Top 10 pathways predicted by PACER for the compound 'nsc23766'

581

Pathway	Contains known target?	P-value of pathway (baseline method)
Signaling events mediated by focal adhesion kinase	YES	0.0001
Integrins in angiogenesis	YES	0.0716
Alpha4 beta1 integrin signaling events	YES	0.0084
a4b7 Integrin signaling	NO	0.1006
Nectin adhesion pathway	YES	0.0521
Integrin-linked kinase signaling	YES	0.4421
Neurotrophic factor-mediated Trk receptor signaling	YES	0.0080
Integrin family cell surface interactions	NO	0.0003
Beta5 beta6 beta7 and beta8 integrin cell surface interactions	NO	0.0013
Plexin-D1 Signaling	NO	0.0364

582

583

584 **REFERENCE**

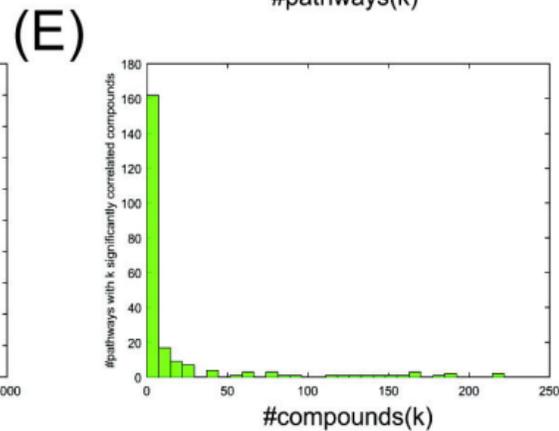
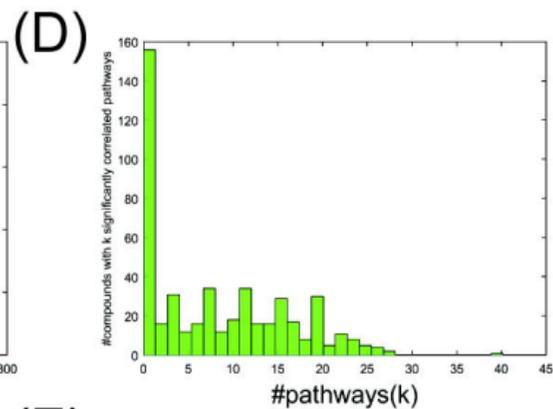
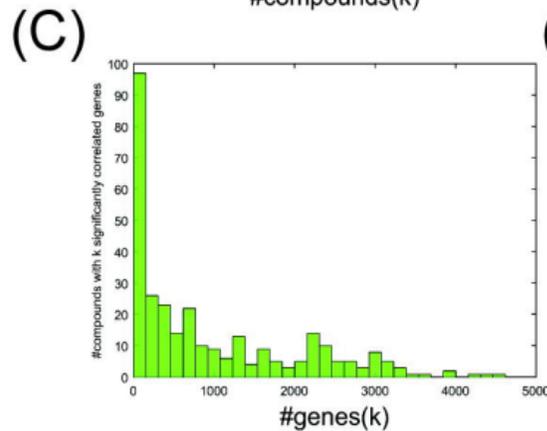
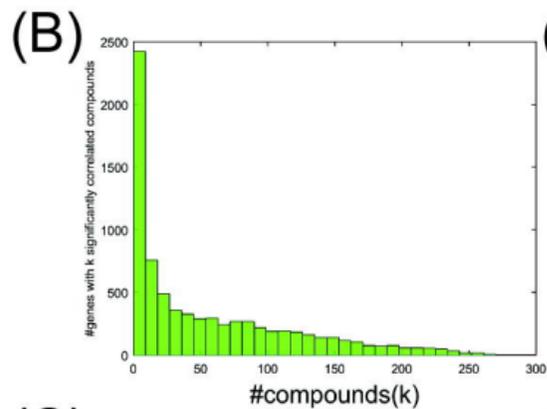
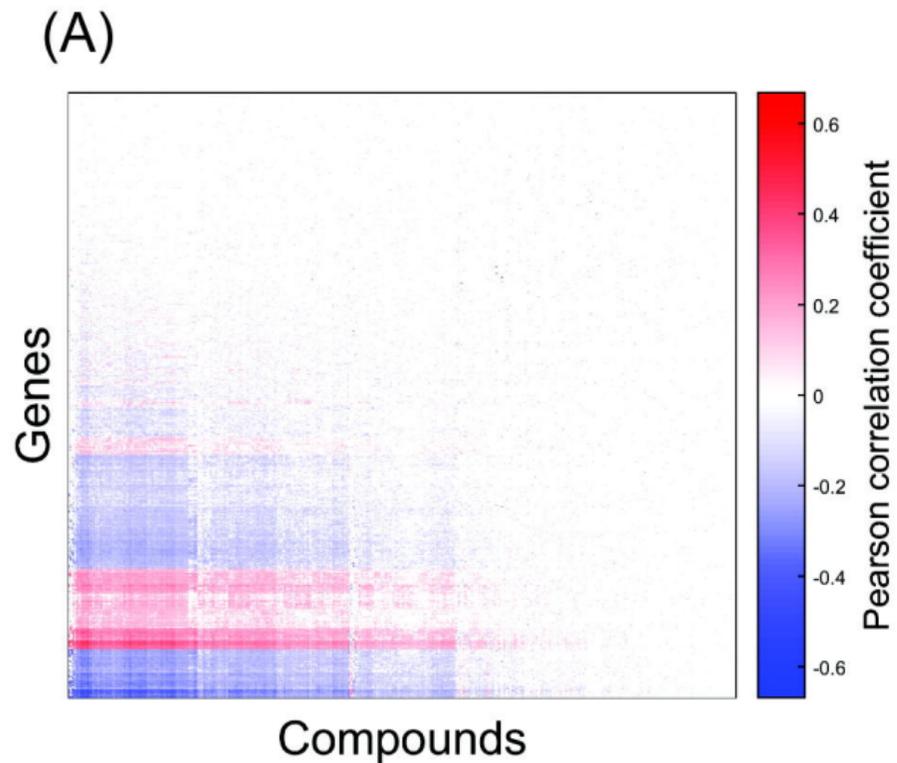
585

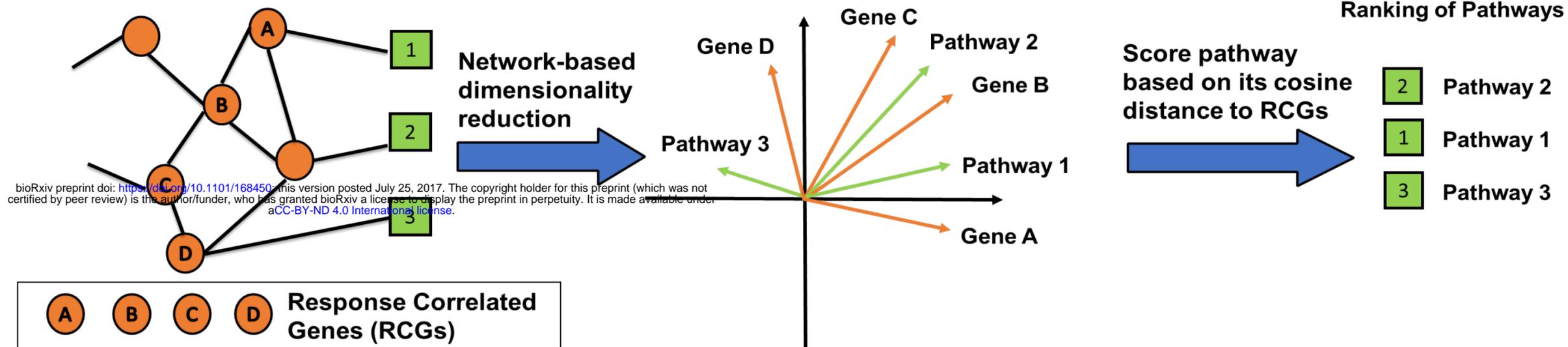
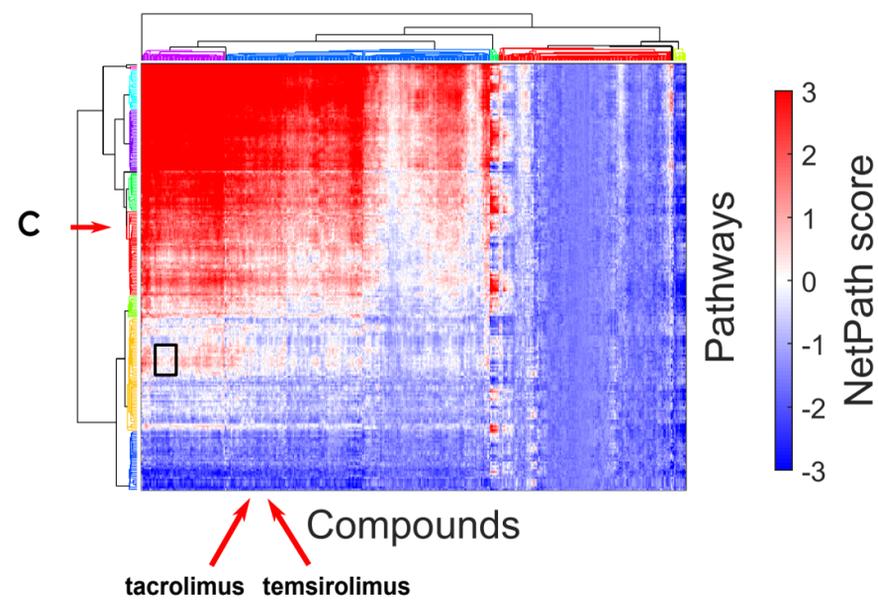
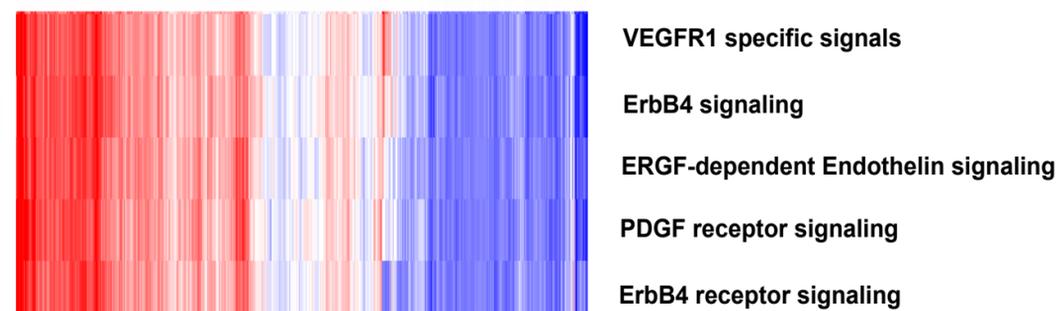
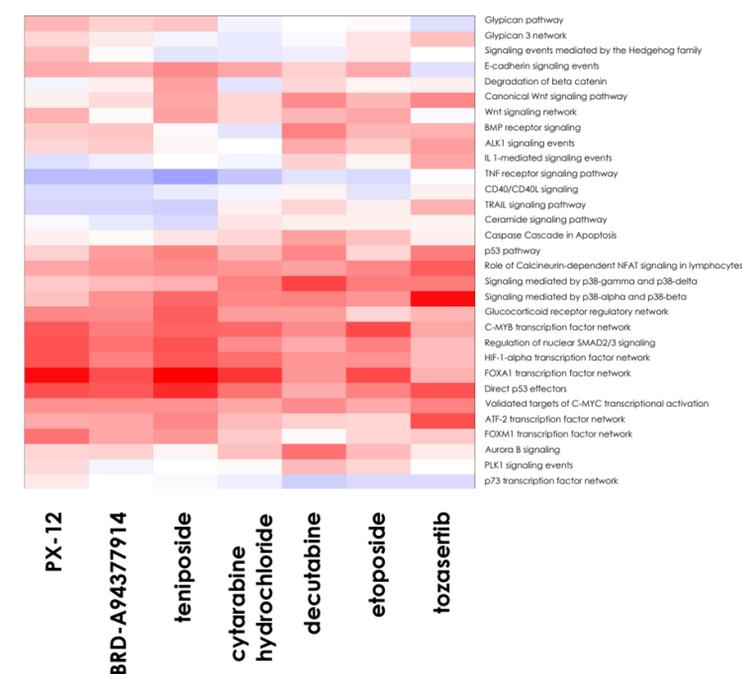
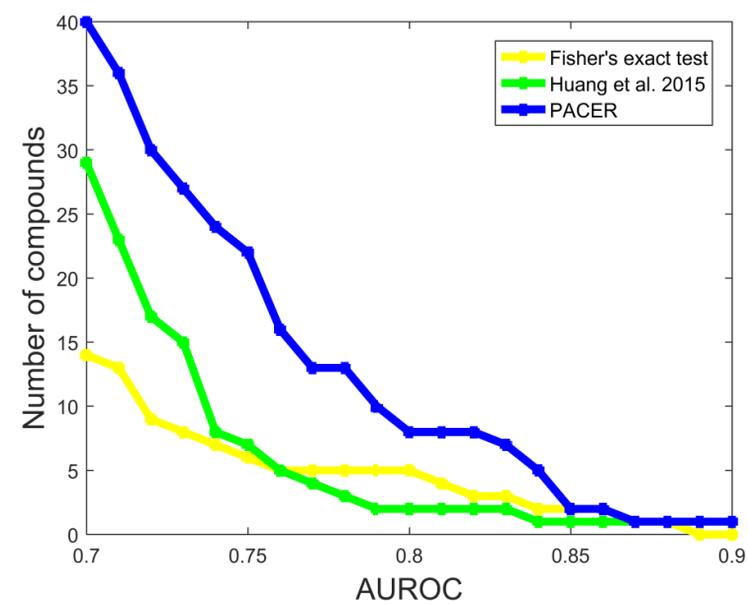
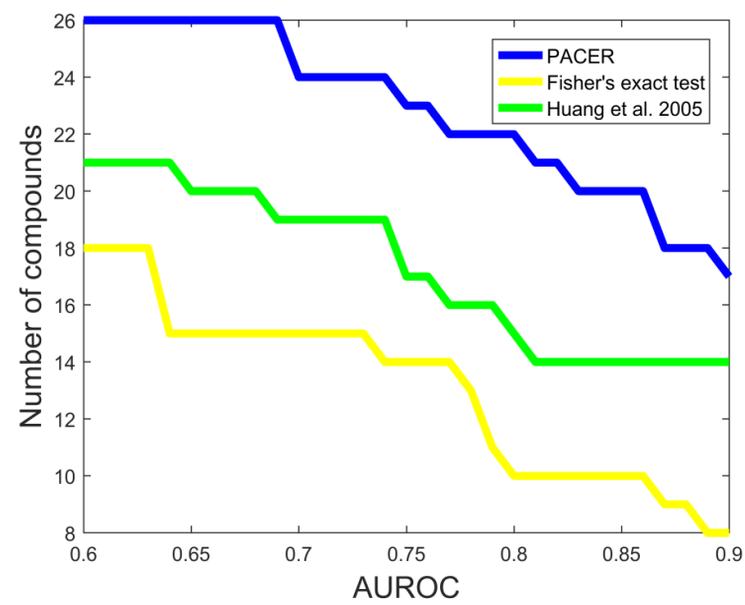
- 586 1. Cancer Genome Atlas Research N, Weinstein JN, Collisson EA, Mills GB, Shaw KR,  
587 Ozenberger BA, *et al.* The Cancer Genome Atlas Pan-Cancer analysis project. *Nat*  
588 *Genet* 2013; **45**(10): 1113-1120.
- 589  
590 2. Stratton MR, Campbell PJ, Futreal PA. The cancer genome. *Nature* 2009; **458**(7239):  
591 719-724.
- 592  
593 3. Barretina J, Caponigro G, Stransky N, Venkatesan K, Margolin AA, Kim S, *et al.* The  
594 Cancer Cell Line Encyclopedia enables predictive modelling of anticancer drug  
595 sensitivity. *Nature* 2012; **483**(7391): 603-607.
- 596  
597 4. Yang W, Soares J, Greninger P, Edelman EJ, Lightfoot H, Forbes S, *et al.* Genomics of  
598 Drug Sensitivity in Cancer (GDSC): a resource for therapeutic biomarker discovery in  
599 cancer cells. *Nucleic Acids Res* 2013; **41**(Database issue): D955-961.
- 600  
601 5. Wang L, McLeod HL, Weinshilboum RM. Genomics and drug response. *N Engl J Med*  
602 2011; **364**(12): 1144-1153.
- 603  
604 6. Wang L, Weinshilboum RM. Pharmacogenomics: candidate gene identification,  
605 functional validation and mechanisms. *Hum Mol Genet* 2008; **17**(R2): R174-179.
- 606  
607 7. Xie Y, Xiao G, Coombes KR, Behrens C, Solis LM, Raso G, *et al.* Robust gene  
608 expression signature from formalin-fixed paraffin-embedded samples predicts prognosis  
609 of non-small-cell lung cancer patients. *Clin Cancer Res* 2011; **17**(17): 5705-5714.
- 610  
611 8. Rees MG, Seashore-Ludlow B, Cheah JH, Adams DJ, Price EV, Gill S, *et al.* Correlating  
612 chemical sensitivity and basal gene expression reveals mechanism of action. *Nat Chem*  
613 *Biol* 2016; **12**(2): 109-116.
- 614  
615 9. Castoreno AB, Eggert US. Small molecule probes of cellular pathways and networks.  
616 *ACS Chem Biol* 2011; **6**(1): 86-94.
- 617  
618 10. Khatri P, Sirota M, Butte AJ. Ten years of pathway analysis: current approaches and  
619 outstanding challenges. *PLoS Comput Biol* 2012; **8**(2): e1002375.
- 620  
621 11. Clevers H. Wnt/beta-catenin signaling in development and disease. *Cell* 2006; **127**(3):  
622 469-480.
- 623  
624 12. Mikkelsen TS, Thorn CF, Yang JJ, Ulrich CM, French D, Zaza G, *et al.* PharmGKB  
625 summary: methotrexate pathway. *Pharmacogenet Genomics* 2011; **21**(10): 679-686.
- 626  
627 13. Thorn CF, Oshiro C, Marsh S, Hernandez-Boussard T, McLeod H, Klein TE, *et al.*  
628 Doxorubicin pathways: pharmacodynamics and adverse effects. *Pharmacogenet*  
629 *Genomics* 2011; **21**(7): 440-446.
- 630  
631 14. Braun R, Cope L, Parmigiani G. Identifying differential correlation in gene/pathway  
632 combinations. *BMC Bioinformatics* 2008; **9**: 488.
- 633

- 634 15. Hoehndorf R, Dumontier M, Gkoutos GV. Identifying aberrant pathways through  
635 integrated analysis of knowledge in pharmacogenomics. *Bioinformatics* 2012; **28**(16):  
636 2169-2175.  
637
- 638 16. Huang R, Wallqvist A, Thanki N, Covell DG. Linking pathway gene expressions to the  
639 growth inhibition response from the National Cancer Institute's anticancer screen and  
640 drug mechanism of action. *Pharmacogenomics J* 2005; **5**(6): 381-399.  
641
- 642 17. Song M, Meiyue S, Yan Y, Zhenran J. Drug–pathway interaction prediction via multiple  
643 feature fusion. *Mol Biosyst* 2014; **10**(11): 2907-2913.  
644
- 645 18. Guo H, Dong J, Hu S, Cai X, Tang G, Dou J, *et al.* Biased random walk model for the  
646 prioritization of drug resistance associated proteins. *Sci Rep* 2015; **5**: 10857.  
647
- 648 19. Hopkins AL. Network pharmacology: the next paradigm in drug discovery. *Nat Chem*  
649 *Biol* 2008; **4**(11): 682-690.  
650
- 651 20. Kotlyar M, Fortney K, Jurisica I. Network-based characterization of drug-regulated  
652 genes, drug targets, and toxicity. *Methods* 2012; **57**(4): 499-507.  
653
- 654 21. Schenone M, Monica S, Vlado D, Wagner BK, Clemons PA. Target identification and  
655 mechanism of action in chemical biology and drug discovery. *Nat Chem Biol* 2013; **9**(4):  
656 232-240.  
657
- 658 22. Wang Y, Xiao J, Suzek TO, Zhang J, Wang J, Bryant SH. PubChem: a public  
659 information system for analyzing bioactivities of small molecules. *Nucleic Acids Res*  
660 2009; **37**(Web Server issue): W623-633.  
661
- 662 23. Duan Q, Flynn C, Niepel M, Hafner M, Muhlich JL, Fernandez NF, *et al.* LINCS Canvas  
663 Browser: interactive web app to query, browse and interrogate LINCS L1000 gene  
664 expression signatures. *Nucleic Acids Res* 2014; **42**(Web Server issue): W449-460.  
665
- 666 24. Franceschini A, Szklarczyk D, Frankild S, Kuhn M, Simonovic M, Roth A, *et al.* STRING  
667 v9.1: protein-protein interaction networks, with increased coverage and integration.  
668 *Nucleic Acids Res* 2013; **41**(Database issue): D808-815.  
669
- 670 25. Schaefer CF, Anthony K, Krupa S, Buchoff J, Day M, Hannay T, *et al.* PID: the Pathway  
671 Interaction Database. *Nucleic Acids Res* 2009; **37**(Database issue): D674-679.  
672
- 673 26. Cho H, Berger B, Peng J. Diffusion Component Analysis: Unraveling Functional  
674 Topology in Biological Networks. *Lecture Notes in Computer Science*, 2015, pp 62-64.  
675
- 676 27. Wang S, Cho H, Zhai C, Berger B, Peng J. Exploiting ontology graph for predicting  
677 sparsely annotated gene function. *Bioinformatics* 2015; **31**(12): i357-364.  
678
- 679 28. Hanson C, Cairns J, Wang L, Sinha S. Computational discovery of transcription factors  
680 associated with drug response. *Pharmacogenomics J* 2015.  
681
- 682 29. Hayes JS, Czekanska EM, Richards RG. The Cell–Surface Interaction. *Advances in*  
683 *Biochemical Engineering/Biotechnology*, 2011, pp 1-31.  
684

- 685 30. Aoudjit F, Vuori K. Integrin signaling in cancer cell survival and chemoresistance.  
686 *Chemother Res Pract* 2012; **2012**: 283181.  
687
- 688 31. Guo W, Pylayeva Y, Pepe A, Yoshioka T, Muller WJ, Inghirami G, *et al.* Beta 4 integrin  
689 amplifies ErbB2 signaling to promote mammary tumorigenesis. *Cell* 2006; **126**(3): 489-  
690 502.  
691
- 692 32. Ding L, Zhang Z, Liang G, Yao Z, Wu H, Wang B, *et al.* SAHA triggered MET activation  
693 contributes to SAHA tolerance in solid cancer cells. *Cancer Lett* 2015; **356**(2 Pt B): 828-  
694 836.  
695
- 696 33. Robinson DR, Yi-Mi W, Su-Fang L. The protein tyrosine kinase family of the human  
697 genome. *Oncogene* 2000; **19**(49): 5548-5557.  
698
- 699 34. Arora A, Scholar EM. Role of tyrosine kinase inhibitors in cancer therapy. *J Pharmacol*  
700 *Exp Ther* 2005; **315**(3): 971-979.  
701
- 702 35. Marengo B, De Ciucis CG, Ricciarelli R, Furfaro AL, Colla R, Canepa E, *et al.* p38MAPK  
703 inhibition: a new combined approach to reduce neuroblastoma resistance under  
704 etoposide treatment. *Cell Death Dis* 2013; **4**: e589.  
705
- 706 36. Neve RM, Sutterluty H, Pullen N, Lane HA, Daly JM, Krek W, *et al.* Effects of oncogenic  
707 ErbB2 on G1 cell cycle regulators in breast tumour cells. *Oncogene* 2000; **19**(13): 1647-  
708 1656.  
709
- 710 37. Kitazaki T, Oka M, Nakamura Y, Tsurutani J, Doi S, Yasunaga M, *et al.* Gefitinib, an  
711 EGFR tyrosine kinase inhibitor, directly inhibits the function of P-glycoprotein in  
712 multidrug resistant cancer cells. *Lung Cancer* 2005; **49**(3): 337-343.  
713
- 714 38. Wang W. A major switch for the Fanconi anemia DNA damage-response pathway. *Nat*  
715 *Struct Mol Biol* 2008; **15**(11): 1128-1130.  
716
- 717 39. Kennedy RD, D'Andrea AD. DNA repair pathways in clinical practice: lessons from  
718 pediatric cancer susceptibility syndromes. *J Clin Oncol* 2006; **24**(23): 3799-3808.  
719
- 720 40. Mulligan JM, Hill LA, Deharo S, Irwin G, Boyle D, Keating KE, *et al.* Identification and  
721 validation of an anthracycline/cyclophosphamide-based chemotherapy response assay  
722 in breast cancer. *J Natl Cancer Inst* 2014; **106**(1): djt335.  
723
- 724 41. Juaristi JA, Aguirre MV, Todaro JS, Alvarez MA, Brandan NC. EPO receptor, Bax and  
725 Bcl-x(L) expressions in murine erythropoiesis after cyclophosphamide treatment.  
726 *Toxicology* 2007; **231**(2-3): 188-199.  
727
- 728 42. Goldstein M, Roos WP, Kaina B. Apoptotic death induced by the cyclophosphamide  
729 analogue mafosfamide in human lymphoblastoid cells: contribution of DNA replication,  
730 transcription inhibition and Chk/p53 signaling. *Toxicol Appl Pharmacol* 2008; **229**(1): 20-  
731 32.  
732
- 733 43. Chen XY, Xia HX, Guan HY, Li B, Zhang W. Follicle Loss and Apoptosis in  
734 Cyclophosphamide-Treated Mice: What's the Matter? *Int J Mol Sci* 2016; **17**(6).  
735

- 736 44. Li Z, Biswas S, Liang B, Zou X, Shan L, Li Y, *et al.* Integrin beta6 serves as an  
737 immunohistochemical marker for lymph node metastasis and promotes cell invasiveness  
738 in cholangiocarcinoma. *Sci Rep* 2016; **6**: 30081.  
739
- 740 45. Desai K, Nair MG, Prabhu JS, Vinod A, Korlimarla A, Rajarajan S, *et al.* High expression  
741 of integrin beta6 in association with the Rho-Rac pathway identifies a poor prognostic  
742 subgroup within HER2 amplified breast cancers. *Cancer Med* 2016; **5**(8): 2000-2011.  
743
- 744 46. Reyes SB, Narayanan AS, Lee HS, Tchaicha JH, Aldape KD, Lang FF, *et al.*  
745 alphavbeta8 integrin interacts with RhoGDI1 to regulate Rac1 and Cdc42 activation and  
746 drive glioblastoma cell invasion. *Mol Biol Cell* 2013; **24**(4): 474-482.  
747
- 748 47. Kumar V, Soni UK, Maurya VK, Singh K, Jha RK. Integrin beta8 (ITGB8) activates VAV-  
749 RAC1 signaling via FAK in the acquisition of endometrial epithelial cell receptivity for  
750 blastocyst implantation. *Sci Rep* 2017; **7**(1): 1885.  
751
- 752 48. Zhang WV, Yang Y, Berg RW, Leung E, Krissansen GW. The small GTP-binding  
753 proteins Rho and Rac induce T cell adhesion to the mucosal addressin MAdCAM-1 in a  
754 hierarchical fashion. *Eur J Immunol* 1999; **29**(9): 2875-2885.  
755
- 756 49. Toyofuku T, Yabuki M, Kamei J, Kamei M, Makino N, Kumanogoh A, *et al.* Semaphorin-  
757 4A, an activator for T-cell-mediated immunity, suppresses angiogenesis via Plexin-D1.  
758 *EMBO J* 2007; **26**(5): 1373-1384.  
759
- 760 50. Cao Y, Jiang T, Girke T. Accelerated similarity searching and clustering of large  
761 compound sets by geometric embedding and locality sensitive hashing. *Bioinformatics*  
762 2010; **26**(7): 953-959.  
763
- 764 51. Subramanian A, Tamayo P, Mootha VK, Mukherjee S, Ebert BL, Gillette MA, *et al.* Gene  
765 set enrichment analysis: a knowledge-based approach for interpreting genome-wide  
766 expression profiles. *Proc Natl Acad Sci U S A* 2005; **102**(43): 15545-15550.  
767
- 768 52. Huang da W, Sherman BT, Lempicki RA. Bioinformatics enrichment tools: paths toward  
769 the comprehensive functional analysis of large gene lists. *Nucleic Acids Res* 2009;  
770 **37**(1): 1-13.  
771
- 772 53. Fabregat A, Sidiropoulos K, Viteri G, Forner O, Marin-Garcia P, Arnau V, *et al.*  
773 Reactome pathway analysis: a high-performance in-memory approach. *BMC*  
774 *Bioinformatics* 2017; **18**(1): 142.  
775
- 776 54. Kramer A, Green J, Pollard J, Jr., Tugendreich S. Causal analysis approaches in  
777 Ingenuity Pathway Analysis. *Bioinformatics* 2014; **30**(4): 523-530.  
778
- 779 55. Alcaraz N, Pauling J, Batra R, Barbosa E, Junge A, Christensen AG, *et al.*  
780 KeyPathwayMiner 4.0: condition-specific pathway analysis by combining multiple omics  
781 studies and networks with Cytoscape. *BMC Syst Biol* 2014; **8**: 99.  
782  
783



**(A)****Heterogeneous network of genes and pathways****(B)****(C)****(D)****(E)****(F)****(G)**