

# A Network-centric Framework for the Evaluation of Mutual Exclusivity Tests on Cancer Drivers

Rafsan Ahmed <sup>1†</sup>, Cesim Erten <sup>2</sup>, Aissa Houdjedj <sup>2</sup>, Hilal Kazan <sup>2,\*</sup>, Cansu Yalcin, <sup>2</sup>

<sup>1</sup>*Electrical and Computer Engineering Graduate Program, Antalya, Antalya Bilim University, Turkey*

<sup>2</sup>*Department of Computer Engineering, Antalya, Antalya Bilim University, Turkey*

<sup>†</sup>*Current Address: Department of Experimental Medical Science, Lund, Lund University, Sweden*

Correspondence\*:

Dr. Hilal Kazan, Department of Computer Engineering, Antalya Bilim University, 07190, Antalya, Turkey  
hilal.kazan@antalya.edu.tr

## 2 ABSTRACT

3 One of the key concepts employed in cancer driver gene identification is that of mutual  
4 exclusivity (ME); a driver mutation is less likely to occur in case of an earlier mutation that  
5 has common functionality in the same molecular pathway. Several ME tests have been proposed  
6 recently, however the current protocols to evaluate ME tests have two main limitations. Firstly the  
7 evaluations are mostly with respect to simulated data and secondly the evaluation metrics lack a  
8 network-centric view. The latter is especially crucial as the notion of common functionality can be  
9 achieved through searching for interaction patterns in relevant networks. We propose a network-  
10 centric framework to evaluate the pairwise significances found by statistical ME tests. It has three  
11 main components. The first component consists of metrics employed in the network-centric ME  
12 evaluations. Such metrics are designed so that network knowledge and the reference set of known  
13 cancer genes are incorporated in ME evaluations under a careful definition of proper control  
14 groups. The other two components are designed as further mechanisms to avoid confounders  
15 inherent in ME detection on top of the network-centric view. To this end, our second objective  
16 is to dissect the side effects caused by mutation load artifacts where mutations driving tumor  
17 subtypes with low mutation load might be incorrectly diagnosed as mutually exclusive. Finally, as  
18 part of the third main component, the confounding issue stemming from the use of nonspecific  
19 interaction networks generated as combinations of interactions from different tissues is resolved  
20 through the creation and use of tissue-specific networks in the proposed framework. The data,  
21 the source code and useful scripts are available at: <https://github.com/abu-compbio/NetCentric>.

22 **Keywords:** mutual exclusivity, network-centric mutual exclusivity evaluation, cancer drivers, cancer genomics, tumor mutation load

## 1 INTRODUCTION

23 Cancer is a disease caused mostly due to a gradual accumulation of somatic alterations that give rise  
24 to pathway dysregulation through alterations in copy number, DNA methylation, gene expression, and  
25 molecular function. An important challenge in cancer genomics is to distinguish driver mutations from  
26 passenger mutations. The former are those determined to be causal for cancer progression, whereas the  
27 latter are characterized as those not leading to any selective advantage. Several computational methods  
28 have been proposed for the identification of cancer driver genes or driver modules of genes by integrating  
29 mutations data with various other types of genetic data; see Dimitrakopoulos and Beerenwinkel [2017],  
30 Zhang and Zhang [2018], Bailey et al. [2018], Tokheim et al. [2016] for recent comprehensive evaluations  
31 and surveys on the topic.

32 A phenomenon observed frequently in the data pertaining to the alterations that the tumors acquire is  
33 *mutual exclusivity (ME)*; a driver mutation is less likely to occur in case of an earlier mutation that has  
34 common functionality in the same molecular pathway [Leiserson et al., 2016, van de Haar et al., 2019,  
35 Thomas et al., 2007, Yeang and Levine, 2008]. Therefore several driver gene or module identification  
36 approaches employ ME detection as part of their problem definitions and optimization goals [Babur  
37 et al., 2015, Ciriello et al., 2012, Leiserson et al., 2013b, Kim et al., 2015, Ahmed et al., 2019, Baali  
38 et al., 2020]. Such a central role in driver gene and module identification has led to the design of many  
39 different approaches for defining and computing mutual exclusivity. Some of these approaches are based  
40 on combinatorial definitions of mutual exclusivity [Vandin et al., 2012, Leiserson et al., 2013a, Basso  
41 et al., 2019, Ahmed et al., 2019, Song et al., 2020, Baali et al., 2020]. In most cases the combinatorial  
42 definitions are incorporated and tested within a driver gene or module identification framework, rather  
43 than as stand-alone ME tests. On the other hand, the vast majority of the ME detection approaches are  
44 based on statistical tests [Ciriello et al., 2012, Szczurek and Beerenwinkel, 2014, Leiserson et al., 2015a,  
45 Constantinescu et al., 2015, Hua et al., 2016, Canisius et al., 2016, Leiserson et al., 2016, Kim et al., 2017,  
46 Liu et al., 2020, Zhang et al., 2020] and in most cases for such approaches the specific goal is to provide  
47 ME significance results. Therefore the focus of the proposed framework is the evaluation of the latter set of  
48 approaches consisting of the statistical ME tests.

49 Among such approaches, MEMo builds a graph based on gene similarities and extracts cliques from this  
50 graph. To determine whether each clique has significant mutual exclusivity, it then proposes a null model  
51 generated by randomly permuting the set of genomic events, while preserving the overall distribution of  
52 observed alterations across both genes and samples, and introduces a Markov Chain Monte Carlo (MCMC)  
53 permutation strategy based on random network generation models [Ciriello et al., 2012]. Szczurek and  
54 Beerenwinkel [2014] propose a probabilistic, generative model of mutual exclusivity, explicitly taking  
55 coverage, impurity, and error rates into account. Based on such a model, they provide a statistical test of  
56 mutual exclusivity by comparing its likelihood to the null model that assumes independent gene alterations.  
57 *Mutex* [Constantinescu et al., 2015] defines the alteration of two genes to be mutually exclusive if their  
58 overlap in samples is significantly less than expected by chance, where the statistical significance of the  
59 overlaps are calculated using a hypergeometric test with the assumption of a uniform alteration frequency  
60 among samples. This may not always be the case as in many data sources there are hyper-mutated samples.  
61 The problem is resolved partially by simply excluding such samples from the analysis. CoMEt [Leiserson  
62 et al., 2015a] on the other hand provides an exact statistical test for mutual exclusivity conditional on  
63 the observed frequency of each alteration with the goal of introducing less bias towards high frequency  
64 alterations. Based on this it provides a tail enumeration procedure to compute the exact test, as well as a  
65 binomial approximation. DISCOVER provides a statistical independence test that makes no assumption of

66 identical gene alteration probabilities across tumors [Canisius et al., 2016]. The alteration probabilities  
67 are estimated by solving a constrained optimization problem guaranteeing the probabilities are consistent  
68 with both the observed number of alterations per gene and the observed number of alterations per tumor.  
69 The tumor-specific gene alteration probabilities are then used to compute the probability of concurrent  
70 alterations which in turn are used to decide whether the number of tumors altered in both genes deviates  
71 from the expectation through an analytical test based on the Poisson-binomial distribution. WeXT provides  
72 a weighted exact test that conditions simultaneously on the number of samples with a mutation and the  
73 per-event, per-sample mutation probabilities [Leiserson et al., 2016]. A recursive formulation to compute  
74 P-values for this weighted test exactly and a saddle-point approximation of the test are proposed. WeSMe  
75 provides a permutation-based test and an approximation of significance through a weighted sampling  
76 technique that enables further improvements in running time spent for sampling and a way to obtain a  
77 better precision without increasing the computational time significantly [Kim et al., 2017]. Two recently  
78 suggested ME tests are FSME [Zhang et al., 2020] and MEScan [Liu et al., 2020]. The former proposes a  
79 *seed-and-extend* strategy to alleviate the computational cost of a permutation-based test. The seed pairs  
80 are constructed by a combinatorial formulation incorporating both ME and the coverage of the pair. The  
81 seeds are then grown with new genes by employing an independence test. MEScan provides a test statistic  
82 that incorporates a patient and gene-specific background mutation rate in the calculation to adjust for the  
83 background noise, and that includes a gene-specific weight to down-weight genes with high mutation rates.  
84 Such a statistic is then employed in an MCMC algorithm followed by a false discovery rate control.

85 We propose a network-centric framework to evaluate the pairwise significances found by statistical  
86 ME tests. It is important to make a distinction between the network-centric view of the current study  
87 and that of the previous studies employing both network data and the concept of ME [Ciriello et al.,  
88 2012, Leiserson et al., 2013b, Kim et al., 2015, Ahmed et al., 2019, Baali et al., 2020]. The latter are  
89 network-centric in the sense that the proposed ME tests are applied on interacting pairs or subnetworks as  
90 part of a more general goal of identifying cancer driver genes/modules. Thus due to the nature of the set  
91 objectives their evaluations focus on the success of output genes/modules matching reference cancer-related  
92 drivers/pathways. The proposed study takes on an approach in the opposite direction; we assume the  
93 interaction network and the reference cancer-related drivers to be inputs to our framework which evaluates  
94 the success of various ME tests. The focus of the proposed framework is on pairwise significances since one  
95 of the major application areas where ME tests are commonly employed is knowledge-based cancer driver  
96 identification where pairwise ME significances are of major essence. In terms of the general objectives our  
97 work is most similar to that of Deng et al. [2017], where a framework for performance comparisons of  
98 statistical ME detection approaches is proposed and executed on six such tests. An important distinction is  
99 that the performance analysis of Deng *et al.* is based on experiments with simulated data and the framework  
100 does not suggest any mechanism to avoid confounders inherent in ME detection. One such confounder is  
101 due to the alterations specific to cancer subtypes [Deng et al., 2017, van de Haar et al., 2019]. Alterations  
102 in different subtypes may be incorrectly diagnosed with ME, although the alterations are not due to any  
103 natural root causes of ME such as redundant functionality. Our network-centric view aims to recognize  
104 such false positives by constructing reference sets based on known drivers gathered from neighborhoods  
105 of interaction networks. Furthermore, inspired by the *mutation load confounding* concept of van de Haar  
106 et al. [2019], we extend our network-centric framework to dissect side effects caused by mutation load  
107 artifacts; mutations that drive tumor subtypes with low mutation load might be incorrectly diagnosed as  
108 mutually exclusive. A possible drawback of the proposed network-centric evaluation framework would be  
109 due to the use of nonspecific interaction networks that are generated as combinations of interactions from  
110 different tissues and are thus suboptimal in resolving confounding issues of mutual exclusivity. In order

111 to detect whether there exists such discrepancies or to limit their effect if they do, we therefore refine the  
112 network-centric approach by designing further tests on *tissue-specific networks (TSN)* we construct based  
113 on gene co-expression.

## 2 METHODS

114 The overall network-centric ME evaluations framework has three main components. The first one consists  
115 of definitions of the metrics employed in the network-centric ME evaluations. Such metrics are designed so  
116 that network knowledge and the reference set of known cancer genes are incorporated in ME evaluations  
117 under a careful definition of proper control groups. The second component detects whether the use of the  
118 interactome information provides similar advantages in ME corrections of pairwise mutual exclusivity  
119 findings as the subtype-stratification idea suggested by van de Haar et al. [2019]. Finally, the third  
120 component extends our framework to incorporate tissue-specific networks with the aim of reducing the  
121 possible side effects of using nonspecific interaction networks.

### 2.1 Metrics for the Network-centric ME Evaluations

123 Assuming that cancer driver genes in the same pathway are more likely to show mutually exclusive  
124 mutation profiles, we utilize the interactome to devise a strategy for evaluating the ME methods and the  
125 effects of the interactome information on quantifying ME. Let  $\mathcal{G}, \mathcal{C}, \mathcal{T}, \mathcal{S}, p_t, c$  denote respectively the  
126 input Protein-Protein Interaction (PPI) network, the employed cohort, the statistical ME test undergoing  
127 the network-centric ME evaluations, the golden standard reference gene set of known cancer drivers, the  
128 p-value threshold for significance, and the type of the control group to be employed. Let  $N_{\mathcal{S}}(g_i)$  denote the  
129 set of genes from  $\mathcal{S}$  that are in the neighborhood of the node corresponding to gene  $g_i$  in the PPI network  
130  $\mathcal{G}$ . For a gene  $g_i \in \mathcal{S}$ , corresponding to each neighbor  $g_j \in N_{\mathcal{S}}(g_i)$ , we randomly select a gene  $g_r$  from a  
131 control group  $\mathcal{X}_c(g_i)$ , and compute  $TP^{cur}, FP^{cur}$ , based on the  $-\log$ -transformed p-values  $p_{i,j}$  and  $p_{i,r}$   
132 as computed by the ME test  $\mathcal{T}$ . Here  $p_{i,j}$  denotes the significance of the mutual exclusivity of the pair  $g_i, g_j$   
133 for  $g_i \in \mathcal{S}$  and  $g_j \in N_{\mathcal{S}}(g_i)$ , and  $p_{i,r}$  denotes the significance of the mutual exclusivity of the pair  $g_i, g_r$   
134 for a random gene  $g_r$  from the control group. Based on the premise that cancer driver genes interacting in the  
135 PPI network are likely to exhibit ME, a pair  $g_i, g_j$  belongs to the set of True Positives if  $p_{i,j}$  is significant  
136 and a pair  $g_i, g_r$  belongs to the set of False Positives if  $p_{i,r}$  is significant. To obtain robust results, the  
137 selection of the random genes from the control group is repeated 100 times. Finally the medians of these  
138 100 instances are summed over all genes  $g_i \in \mathcal{S}$  to provide the necessary statistics  $TP, FP$ . Thus precision,  
139 sensitivity, and the F1 scores are computed based on these 4 statistics. The number of condition positives  
140 used in the sensitivity calculations corresponds to the number of pairs  $g_i, g_j \in \mathcal{S}$  where  $g_i, g_j$  interact in  $\mathcal{G}$ .

141 We note that limiting our focus solely on these conventionally formed  $TP, FP$  classes may be misleading  
142 as each one considers the significance of  $p_{i,j}$  and  $p_{i,r}$  individually. A more detailed inspection with a  
143 simultaneous consideration of their values could prove more insightful in certain cases since they both  
144 involve a common gene  $g_i$ . Towards this aim we introduce the *strict* versions of these conventional classes.  
145 More specifically  $TP_{strict}$  consists of  $g_i, g_j$  pairs where  $p_{i,j}$  is significant not only with respect to the  
146 given threshold but also as compared to the p-value of the control pair  $g_i, g_r$ . Similarly  $FP_{strict}$  consists  
147 of the control pairs  $g_i, g_r$ , where  $p_{i,r}$  is more significant than both the threshold value and  $p_{i,j}$ . Based on  
148 these strict classes we can compute three metrics:  $\text{precision}_{strict}$ ,  $\text{sensitivity}_{strict}$ , and  $\text{F1}_{strict}$ . Such a  
149 consideration is especially convenient in reducing any potential bias inherent in genes like TP53 which  
150 have large mutation frequencies almost exclusively in tumors with small numbers of mutations; both  $p_{i,j}$   
151 and  $p_{i,r}$  are likely to be significant in such a scenario giving rise to vagueness in the conventional F1 score.

152 A comparison of  $F1_{strict}$  values based on the two statistics simultaneous by their nature,  $TP_{strict}$  and  
153  $FP_{strict}$  provides a more rigorous evaluation in such cases.

154 For the network-centric ME evaluations we employ two different definitions for the control groups. For  
155 the first one, the control group  $\mathcal{X}_1(g_i)$  consists of genes in  $\mathcal{S}$  that do not interact with  $g_i$  in the PPI network.  
156 For the second one,  $\mathcal{X}_2(g_i)$  consists of neighbors of  $g_i$  in the PPI network that are not in  $\mathcal{S}$ . In the latter case  
157 only the genes  $g_i \in \mathcal{S}$  for which the number of neighbors not in  $\mathcal{S}$  is larger than or equal to the number of  
158 neighbors in  $\mathcal{S}$  are taken into account.

## 159 2.2 Network-centric ME Corrections in Relation to MLA

160 Some statistical mutual exclusivity tests are based on the assumption that genes' alterations across tumors  
161 are identically distributed. Among the approaches considered in this study Fisher's Exact Test and MEGSA  
162 belong to this category. However, it has been observed that the number of alterations per tumor can vary  
163 quite considerably, even in tumors of the same type; colorectal tumors with microsatellite stability have  
164 a median of 66 non-synonymous mutations, but colorectal tumors with microsatellite instability have a  
165 median of 777 mutations [Vogelstein et al., 2013, Leiserson et al., 2016]. It has been shown that under  
166 such settings the mutual exclusivity tests relying on identical alteration probabilities across tumors may  
167 lead to reduced sensitivity for mutual exclusivity analysis [Canisius et al., 2016]. The effects of varying  
168 alteration probabilities on pairwise mutual exclusivity calculations have been formalized within the context  
169 of the so-called *mutation load confounding (MLC)* in a recent study by van de Haar et al. [2019]. MLC is a  
170 correlation between the number of statistically significant mutual exclusivity findings and the *mutation load*  
171 *association (MLA)* of a gene, where logistic regression is used to compute MLA as a standardized score of  
172 association between the mutation likelihood of each gene and the *mutation load*, that is the genome-wide  
173 number of somatic mutations observed in a tumor. Note that negative MLA values correspond to higher  
174 mutation frequencies in tumors with low mutation loads, whereas positive values correspond to higher  
175 mutation frequencies in tumors with high mutation loads. Strong negative correlations between the MLA  
176 of a gene and the number of statistically significant pairwise mutual exclusivities have been observed,  
177 implicating the finding that the more negative a gene's MLA, the higher the number of other genes that  
178 show mutual exclusivity with that particular gene [van de Haar et al., 2019]. However, such a negative  
179 correlation does not always imply true ME since a gene that exclusively shows large mutation frequency in  
180 tumors with low mutation loads, naturally has a better chance of forming mutually exclusive pairs with  
181 other genes. Thus extra sources of information are necessary to filter out the pairs with true ME relations  
182 among a set of statistically significant pairwise mutual exclusivities postulated by some exclusivity test.  
183 van de Haar et al. [2019] make use of the subtype information for such a purpose and show that MLC can  
184 be reduced by correcting via tumor subtype stratification. Such a correction greatly reduces the number of  
185 gene pairs reported to show mutual exclusivity, especially for pairs that include genes with low MLA. A  
186 major drawback is the absence of subtype information for many tumors. As part of our network-centric  
187 ME framework, we suggest that such a correction can be efficiently done with the interaction network data,  
188 rather than or better yet on top of the subtype information. For this purpose we calculate the correlation  
189 between the number of statistically significant pairwise ME findings and the MLA for two settings; one  
190 where pairwise mutual exclusivities are sought between a gene in  $\mathcal{S}$  and all other genes in  $\mathcal{S}$ , and the other  
191 where a gene in  $\mathcal{S}$  is checked against only its PPI neighbors that are in  $\mathcal{S}$ . The computations of the two  
192 settings are repeated with the subtype-stratified data as well, to see the added value of the network-centric  
193 ME corrections on top of the subtype-based corrections on statistically significant pairwise MEs.

## 194 2.3 Network-centric ME Evaluations in Relation to TSN

195 Rather than using a common nonspecific network for all the cancer types, in this component of our  
196 evaluation framework we employ TSN based on the tissue in which the tumor develops. To construct the  
197 TSN for a particular tissue, we start with the original PPI network and remove the edges between the pairs  
198 of genes that are not co-expressed in the corresponding tissue. For this purpose, we download RNA-seq  
199 datasets from GTEX portal [GTExConsortium, 2020]. See Supplementary Table 49 for the total number of  
200 available samples for each tissue. To determine the co-expressed genes, we follow the procedure described  
201 in Luck et al. [2020]. For each pair of genes that have an edge in the original PPI network, we identify the  
202 number of samples where both genes have Transcripts Per Kilobase Million (TPM) values  $\geq 1$ . We then  
203 divide this number with the total number of samples where either gene has a TPM value  $\geq 1$ . The resulting  
204 value is called the *co-expression ratio*. Gene pairs interacting in the original network are included in the  
205  $TSN_{cor}$  if the *co-expression ratio* is  $\geq cor$ , for a given threshold *cor*.

206 In addition to applying the network-centric metrics introduced in Section 2.1 on the constructed TSNs,  
207 we also propose a more detailed evaluation in terms of ROC analysis based on tissue-specificity. For this  
208 purpose, we define the gene pairs with co-expression ratio value of 1 as *tissue-specific gene pairs*. Similarly,  
209 the gene pairs with co-expression ratio values  $\leq 0.5$  are called *non-tissue-specific gene pairs*. To test  
210 whether a specific ME test identifies stronger mutual exclusivities for the tissue-specific gene pairs in  $\mathcal{S}$ ,  
211 we rank the gene pairs in  $\mathcal{S}$  in increasing order of p-values. To construct the control group, we rank the  
212 same number of random samples of gene pairs not in  $\mathcal{S}$  with respect to the p-values making sure that the  
213 sizes of the positive (or negative) sets of gene pairs not in  $\mathcal{S}$  are exactly the same as those that are found  
214 for the gene pairs in  $\mathcal{S}$ . For both gene pairs in  $\mathcal{S}$  and gene pairs not in  $\mathcal{S}$ , the set of positives consists of  
215 the tissue-specific gene pairs, whereas non-tissue-specific gene pairs are labelled as negatives. We then  
216 compute the True Positive Rate (TPR) and the False Positive Rate (FPR) for each case. Note that for  
217 robustness considerations the control group computations are repeated 100 times and the median TPR and  
218 FPR values are reported.

## 3 RESULTS

### 219 3.1 Input Data and Parameter Settings

220 The somatic mutation data from TCGA was preprocessed and provided by van de Haar et al., 2019. The  
221 8 different cancer types and their corresponding tumor samples within the dataset is as follows: BLCA  
222 (411), BRCA (1026), COADREAD (498), LUAD (568), LUSC (485), SKCM (468), STAD (438) and  
223 UCEC (531). The preprocessing step involves the removal of all mutations with ‘variant\_classification’  
224 of ‘Silent’, ‘3’UTR’, ‘Intron’, ‘5’UTR’, ‘RNA’, ‘3’Flank’ and ‘5’Flank’ from the TCGA data. The input  
225 data is then further filtered by mutation frequency threshold,  $t$ , to include genes with  $> t$  mutations across  
226 the cohort. More specifically, with  $t = 20$  we include the genes that are mutated in more than 20 samples  
227 within the cancer type under study. Regarding subtypes, we download subtype information for BRCA  
228 from the cBioPortal [Cerami et al., 2012, Gao et al., 2013] and the CMS stratification for COADREAD  
229 from [Guinney et al., 2015]. We use the COSMIC Cancer Gene Census database to compile the set of  
230 known cancer genes [Sondka et al., 2018]. For the results presented in the main document we employ the  
231 IntAct PPI network as it is a comprehensive and well-characterized database [Orchard et al., 2014]. As a  
232 preprocessing step, we remove duplicate edges and edges below the confidence threshold of 0.35 from  
233 the network. The final network contains 15,079 nodes and 103,520 edges. For the gene expression data

234 employed in the construction of TSNs, we download RNA-Seq data from the Genotype-Tissue Expression  
235 (GTEx) portal [GTExConsortium, 2020] (05-06-2017).

236 For the comparative evaluations of our network-centric framework described in the previous section, we  
237 choose six popular statistical mutual exclusivity methods: DISCOVER [Canisius et al., 2016], DISCOVER  
238 Strat [Canisius et al., 2016, van de Haar et al., 2019], Fisher's Exact Test, WeXT [Leiserson et al.,  
239 2016], MEMo [Ciriello et al., 2012] and MEGSA [Hua et al., 2016]. Among these, MEMo and MEGSA are  
240 originally designed to output p-values for a set of genes with size  $> 2$ . For MEMo, we re-implement the first  
241 part of the algorithm where pairwise ME p-values are estimated. We use  $Q = 100$  and  $N = 10,000$  as  
242 suggested by the original paper [Ciriello et al., 2012]. For MEGSA, pairwise ME p-values are calculated by  
243 applying chi-square cumulative probability less than or equal to the value of the log likelihood calculated by  
244 the *funestimate* function. We should note that all the employed methods output multiple testing corrected  
245 p-values. With regards to the parameter settings of our proposed framework, we employ the values of 5 and  
246 20 for  $t$ .

### 247 3.2 ME Evaluations Based on Defined Metrics

248 Table 1 and 2 show the results of evaluating the six ME detection methods on COADREAD data where  
249  $t = 20$  and we use the data from 498 patients for which subtype information is available. We use  $\mathcal{X}_1$  and  
250  $\mathcal{X}_2$  as the control group in Table 1 and 2, respectively. We first discuss the results of  $\mathcal{X}_1$ . We observe that  
251 DISCOVER Strat gives the highest precision and precision<sub>strict</sub> values. The ranking of the other methods  
252 from best to worst in terms of precision or precision<sub>strict</sub> is as follows: WeXT, DISCOVER, MEMo,  
253 MEGSA and Fisher's Exact Test. A comparison of the precision and precision<sub>strict</sub> values distinguishes  
254 two groups of ME methods; for DISCOVER, DISCOVER Strat, Fisher's Exact Test, and WeXT the  
255 precision<sub>strict</sub> values are greater than or equal to the precision values, whereas the exact opposite is  
256 observed for MEGSA and MEMo. This suggests that the performance of the methods in the latter group  
257 gets worse when random control gene pair is considered simultaneously in the precision calculation, that is  
258 precision<sub>strict</sub>. Compared to the precision, we observe much larger differences among the sensitivity or  
259 the sensitivity<sub>strict</sub> values output by the employed methods. We can group the methods into two where  
260 the first group contains WeXT, MEMo and DISCOVER, and the second group contains the rest of the  
261 methods. The first group of methods give much larger sensitivity or sensitivity<sub>strict</sub> values than the second.  
262 For instance, the sensitivity value obtained with WeXT is an order of magnitude larger than that of Fisher's  
263 Exact Test. This also shows that the second group of methods are more conservative than the first group of  
264 methods. WeXT is the least conservative approach based on its high sensitivity value. Even though WeXT  
265 predicts many significant p-values, it still has a competitive precision value which is slightly lower than  
266 the maximum observed value (0.725 vs 0.727). Accordingly, WeXT obtains the best F1 score and F1<sub>strict</sub>  
267 score which is followed by MEMo and DISCOVER. The remaining three methods give much smaller F1  
268 scores and they rank as follows from highest to lowest: MEGSA, DISCOVER Strat and Fisher's Exact Test.  
269 Comparing the conventional F1 score with the F1<sub>strict</sub> score of each ME method, the largest difference is  
270 observed for MEMo indicating that the consideration of the random pair as a control affects its performance  
271 dramatically. Another interesting observation is the lower performance of DISCOVER Strat compared to  
272 DISCOVER which suggests that the use of subtype information is not useful for COADREAD. Table 2  
273 shows the results where  $\mathcal{X}_2$  is used as the control group. Since  $\mathcal{X}_2(g_i)$  is defined as the non-CGC neighbors  
274 of  $g_i$  in the PPI network, we can only consider the CGC genes that have more non-CGC neighbors than  
275 CGC neighbors. As such, the number of pairs included in this analysis is much smaller than that of Table 1  
276 (107 vs 196). The ranking of the methods in Table 2 with respect to F1 score and sensitivity remain the  
277 same as Table 1. However, there are differences in the ranking with respect to other metrics. For instance,

278 WeXT ranks best in terms of precision whereas the best ranking method in Table 1, DISCOVER Strat, ranks  
279 the fifth. Compared to Table 1, the precision values of all the methods are smaller in Table 2. We see the  
280 opposite trend for sensitivity values. These changes are in parallel with the increase in percent significant  
281 p-values output by the methods. For instance, the percentage of significant p-values output by DISCOVER  
282 is 12% in Table 1 and 18% in Table 2. We also observe differences between the conventional and the  
283 strict versions of the employed metrics. WeXT and DISCOVER have increased  $\text{precision}_{\text{strict}}$  values  
284 compared to precision whereas we observe the opposite trend for the rest of the methods. Additionally, the  
285 ranking of the methods with respect to F1 score and  $\text{F1}_{\text{strict}}$  score is different. Namely, MEMo's ranking  
286 decreases from second highest to third highest when we switch from F1 score to  $\text{F1}_{\text{strict}}$  score. Accordingly,  
287 DISCOVER's ranking improves from third highest to second highest based on F1 score. This increases the  
288 confidence of DISCOVER results as  $\text{F1}_{\text{strict}}$  requires a stricter definition of true and false positives.

289 Table S1 show the results with  $\mathcal{X}_1$  control group and  $t = 20$  filtering for the other cancer types. We  
290 observe that the methods report a small number of significant p-values for BLCA data (Table S1-a). In line  
291 with this, we observe smaller sensitivity values for BLCA data compared to the results on COADREAD  
292 data. For BRCA dataset, we observe that the top F1 score is obtained by DISCOVER Strat indicating  
293 the benefit of considering subtype information for BRCA (Table S1-b). When we compare the values of  
294 the conventional and strict versions of the metrics, we observe that DISCOVER Strat's  $\text{precision}_{\text{strict}}$  is  
295 significantly larger than its precision value whereas we observe negligible changes for the other methods.  
296 Accordingly, DISCOVER Strat's  $\text{F1}_{\text{strict}}$  score is slightly lower than its F1 score whereas we observe  
297 significant drops for the other methods. We also observe that MEGSA and Fisher's Exact Test perform  
298 significantly worse than the other methods. Interestingly MEGSA ranks the best in terms of F1 score  
299 for LUSC dataset (Table S1-e). However, we should point out that the methods report a small number  
300 of significant p-values similar to the BLCA dataset. WeXT shows a dramatically better performance on  
301 SKCM dataset (Table S1-f). For UCEC dataset, we observe that Fisher's Exact Test and MEGSA perform  
302 poorly due to their conservative calculation of p-values when compared to DISCOVER and WeXT (Table  
303 S1-h). Table S13 show the analogous results when the control group is defined as  $\mathcal{X}_2$ . We observe that  
304 the analysis includes less than 50 pairs for BLCA, BRCA and LUSC. Again, we observe that Fisher's  
305 Exact Test and MEGSA report zero or very few number of significant p-values across the majority of the  
306 the cancer types. For LUAD, SKCM, STAD and UCEC, WeXT gives the largest F1-values. Similar to  
307 the results obtained with  $\mathcal{X}_1$ , MEGSA ranks the top in terms of F1-score on LUSC dataset. MEGSA and  
308 MEMo results are not available for some cancer types since we are unable to run these methods due to  
309 memory issues.

310 Table 3 and 4 show the COADREAD results of  $t = 5$  setting with  $c = X_1$  and  $c = X_2$ , respectively.  
311 Using a lower value for  $t$  increases the number of gene pairs tested in our analysis. When we compare these  
312 results with the results we obtained when  $t = 20$ , we observe few differences. Though the number of tested  
313 gene pairs is larger, the percentage of significant p-values obtained by the methods decreases. For instance,  
314 the percentage of significant p-values output by WeXT for COADREAD data decreases from 42% to 14%  
315 when  $t$  is changed from 20 to 5. This is likely related to the larger inclusion of low mutation frequency  
316 genes when  $t = 5$ . An interesting observation for  $t = 5$  results is the decrease in DISCOVER Strat's  
317 performance. For COADREAD, DISCOVER Strat's precision and  $\text{precision}_{\text{strict}}$  value is the highest for  
318  $t = 20$  when  $\mathcal{X}_1$  is used as the control group. However, when  $t = 5$ , we observe that it ranks after WeXT  
319 and DISCOVER in terms of precision/ $\text{precision}_{\text{strict}}$  value. Similarly, for BRCA dataset, DISCOVER Strat  
320 ranks after WeXT for both control groups  $\mathcal{X}_1$  and  $\mathcal{X}_2$  (Table S25-b, Table S37-b).



### 321 3.3 Robustness Analysis of Evaluations Based on Defined Metrics

322 We also investigate the robustness of our results with respect to *robustness.iterations* value, the p-value  
323 significance threshold value, the reference gene set and the employed PPI network. For *robustness.iterations*,  
324 we try the values 300 and 500. For p-value significance threshold, we try a more stringent significance  
325 threshold of 0.01. Since the employed ME methods output multiple testing corrected p-values, we also try  
326 a less stringent significance threshold: 0.1. Regarding the reference gene set, we try using a subset of CGC  
327 genes to include only those which have SNV type of mutations in cancer (378 out of 723 genes). To this end,  
328 we filter out the genes where the *mutation type* column consists of only A (amplification), D (large deletion)  
329 or T (translocation). Additionally, we use an alternative source named *IntoGen* [Martínez-Jiménez et al.,  
330 2020] to compile reference cancer genes. We download *Unfiltered\_driver\_results\_05.tsv* file (2020-02-  
331 02 release) and include the genes where FILTER column is *PASS*, which results in 503 genes. For the PPI  
332 network, we try different confidence threshold values for filtering IntAct. 0.45 threshold value is commonly  
333 used in the literature to filter out the interactions with low confidence [Sügis et al., 2019, Porras et al.,  
334 2020]. Additionally, to observe the effect of using a lower threshold than the current one, we also tried  
335 filtering IntAct with a threshold value of 0.25. Lastly, we utilize two alternative networks in our analyses:  
336 *HINT + HI2012* [Das and Yu, 2012, Yu et al., 2011, Leiserson et al., 2015b] and STRING network  
337 [Szkarczyk D et al., 2018]. For the latter, we download the file *9606.protein.physical.links.v11.0.txt* and  
338 only use physical interactions with scores greater than 700. This results in 9,524 nodes and 146,120 edges.

339 The results with these different settings are available in Table S2-S9 for  $t = 20$  and  $c = X_1$ . Table S2 and  
340 S3 show the results with *robustness.iterations* value 300 and 500, respectively. When *robustness.iterations*  
341 value is increased from 100 to 300, the ranking of the ME methods based on F1 score remains the same  
342 for all the cancer types except for BRCA (Table S2). Similarly, when the *robustness.iterations* value is  
343 increased from 100 to 500, the ranking of the ME methods based on F1 score changes only for BLCA  
344 and BRCA (Table S3). Table S4 shows the results where p-value significance threshold is decreased from  
345 0.05 to 0.01. We observe that the 0.01 threshold is too strict for most of the methods as evident from  
346 low sensitivity values for most cancer types. In fact, Fisher's Exact Test predicts no mutually exclusive  
347 cases for six out of eight cancer types. When we use the p-value threshold of 0.1, we observe an overall  
348 improvement in F1 scores for all cancer types and for all the methods except for MEGSA which shows  
349 decreased F1 scores for four cancer types (Table S5). The observed increase in F1 scores for the majority  
350 of the cases is due to the large increase in sensitivity values. Additionally, for some cases precision values  
351 also increase with this less stringent p-value threshold. In particular, we observe a dramatic increase in  
352  $F1/F1_{strict}$  scores for LUAD cancer type. In terms of ranking of the methods, we observe a difference in  
353 BLCA and BRCA types. For BLCA, WeXT's ranking improves from fourth place to first place whereas  
354 MEMO's ranking decreases from first place to third place. We observe a similar switch in ranking of the  
355 methods for BRCA. When  $CGC_{SNV}$  is used as the reference gene set, F1 score-based ranking of the ME  
356 methods changes for three cancer types: BLCA, BRCA and LUSC (Table S6). Though, among these three  
357 cancer types, the top ranking method remains the same for BRCA and LUSC. When *IntoGen* is used as  
358 the reference set, we only observe a change in ranking of the ME methods for BLCA and BRCA types  
359 (Table S7). Next, for the experiments where we try different confidence thresholds for IntAct or switch  
360 to the *HINT + HI2012* network, the only cancer types where the ranking of the ME methods change  
361 are BLCA and BRCA (Table S8-S10). When we use the STRING network, we observe a decrease in  
362  $F1/F1_{strict}$  scores of all the methods except for WeXT for BRCA and UCEC cancer types as compared to  
363 the results obtained with the IntAct network (Table S11). For COADREAD, LUAD and STAD we observe  
364 both increases and decreases in  $F1/F1_{strict}$  scores. Lastly, all the methods show decreased  $F1/F1_{strict}$

365 scores in LUSC and SKCM. The ranking of the methods based on F1 score and  $F1_{strict}$  score differs  
366 for BRCA and COADREAD datasets. For BRCA, WeXT and DISCOVER Strat switch positions in the  
367 ranking based on  $F1/F1_{strict}$  scores. For COADREAD, MEMO's ranking decreases from first position to  
368 third position showing that its top performance is not preserved with stricter versions of evaluation metrics.  
369 Additionally, the ranking of the methods change in five out of eight cancer types as compared to the results  
370 obtained with IntAct network suggesting that the input PPI plays a role in the performance of ME methods.

371 Results of robustness analyses for  $t = 20$  and  $c = X_2$  setting is available in Table S13-S20. When  
372 *robustness.iterations* is increased to 300 or to 500, the ranking of the ME methods remains the same  
373 for all the cancer types. Using a stricter p-value significance threshold of 0.01 results in too few mutual  
374 exclusivities predicted for many of the cancer types (Table S15). For instance, for LUAD, the sensitivity  
375 values are zero for all the methods. Similarly, for BLCA, SKCM and STAD cancer types the maximum  
376 sensitivity value across the ME methods is 0.1. When p-value threshold is switched to 0.1, we observe  
377 changes in both directions for different cancer types (Table S16). For BRCA,  $F1/F1_{strict}$  scores of all the  
378 methods either remain the same or decrease. On the other hand, all the methods have increased  $F1/F1_{strict}$   
379 scores for COADREAD. We also observe a change in the ranking of the methods when we change the  
380 p-value threshold from 0.05 to 0.1. For LUAD, SKCM and UCEC, we observe significant increase in  
381  $F1/F1_{strict}$  scores of all the methods. For LUSC and STAD, we see both increases and decreases in  
382  $F1/F1_{strict}$  scores. In particular, WeXT and DISCOVER show dramatic increases in  $F1/F1_{strict}$  scores  
383 for STAD type. Overall, we can conclude that the p-value threshold of 0.1 leads to a better performance  
384 across the employed methods. When the reference gene set is changed to  $CGC_{SNV}$  or to *IntoGen*, the F1  
385 score-based ranking of the methods only change for BLCA (Table S18-S19). When we use a version of  
386 IntAct filtered with confidence value 0.25, the ranking remains the same for all the cancer types except  
387 for BLCA, BRCA and LUSC (Table S20). When the confidence value threshold is increased to 0.45, the  
388 number of considered CGC-CGC gene pairs decreases to values that are  $< 20$  for all the cancer types  
389 except for SKCM and UCEC (Table S21). For SKCM and UCEC, we observe the same ranking as in  
390 our original parameter settings. Switching to *HINT + HI2012* network also leads to the inclusion of  
391 too few gene pairs for many cancer types (Table S22). For the rest, the ranking of the ME methods is in  
392 accordance with our original results. When we use the STRING network as the input PPI, we observe very  
393 few CGC-CGC pairs for BLCA, BRCA and LUSC. Comparing these results with those obtained with the  
394 IntAct network, for COADREAD, MEMO and WeXT are the only methods that show decreased  $F1/F1_{strict}$   
395 scores where the magnitude of change is much larger for MEMO. However, these changes do not lead to a  
396 difference in the ranking of the methods. For LUAD, we observe large improvements in  $F1/F1_{strict}$  scores  
397 of all the methods. Lastly, for SKCM, STAD and UCEC we observe smaller  $F1/F1_{strict}$  scores in almost  
398 all the cases (Table S23). These results reveal that the change of the network leads to distinct changes in  
399 different cancer types.

400 Robustness analyses for  $t = 5$  setting is available in Table S25-S35 and in Table S37-S47 for  $c = X_1$  and  
401  $c = X_2$ , respectively. We observe similar patterns compared to the  $t = 20$  setting. For instance, changing  
402 the p-value threshold to 0.01 decreases the  $F1/F1_{strict}$  scores overall whereas increasing the threshold to  
403 0.1 also increases the  $F1/F1_{strict}$  scores. Overall, changing the different settings of the analysis do not lead  
404 a change in the ranking of the methods.

405 We also assess whether the F1 scores improve or worsen when different confidence thresholds are used  
406 to filter the IntAct network. Increasing the interaction confidence threshold increases the support of both  
407 the reference pair and the control pair with respect to  $X_2$ , since they are both interacting pairs. In parallel  
408 with this, for the  $t = 5$  setting where we have adequate number of gene pairs for all the cancer types, we

409 observe an increased F1 score for all the methods for all the cancer types except for BRCA and LUAD.  
410 For these two cancer types, we observe both increases and decreases in F1 scores. On the other hand, it  
411 is difficult to propose a similar argument for  $c = X_1$  setting, since increasing the interaction confidence  
412 threshold increases the support of the reference pair but decreases that of the control pair. This is due to the  
413 possibility of considering a random reference cancer gene as a non-neighbor with the 0.45 threshold even  
414 though it would be considered as a neighbor with a lower confidence threshold e.g. 0.35.

### 415 3.4 ME Evaluations Based on Corrections via MLA

416 Having compared the ME tests with respect to our novel network-centric evaluation framework, we now  
417 assess whether including network knowledge reduces the mutation load confounding (MLC) problem  
418 introduced by van de Haar *et al.* [2019]. van de Haar *et al.* identified a strong negative correlation between  
419 the MLAs of genes and their percent significant findings in mutual exclusivity tests. In van de Haar *et al.*  
420 [2019], these statistics are computed for a set of 341 genes from an established cancer gene panel [Cheng  
421 *et al.*, 2015] where, for each gene, mutual exclusivity tests are performed with all the other genes in the  
422 panel. Here, we first perform a similar analysis where we use the COSMIC CGC database [Forbes *et al.*,  
423 2017] to define the reference cancer gene set as it is more comprehensive and up to date.

424 Figure 1-A shows the MLA of the reference cancer genes vs the percent significant findings in mutual  
425 exclusivity tests performed with DISCOVER for the TCGA COADREAD cohort (498 tumors). We observe  
426 a strong negative correlation between MLA values and percent significant findings in mutual exclusivity  
427 tests (Pearson correlation  $-0.88$ , p-value  $3.0e - 25$ ) similar to van de Haar *et al.* [2019]. In Figure 1-B,  
428 we take into account the PPI information to calculate percent significant findings. Namely, for each CGC  
429 gene, we perform mutual exclusivity tests only with its PPI neighbors that are also in CGC. Note that CGC  
430 genes which do not have any CGC neighbors are excluded from this analysis. To make a fair comparison  
431 between Figure 1-A and 1-B, only the CGC genes that have CGC neighbors are shown in Figure 1-A. We  
432 also ensure that the mutual exclusivity of a gene of interest is checked with same sized group of genes in  
433 both Figure 1-A and 1-B. To achieve this in Figure 1-A, for each gene, we compute mutual exclusivity with  
434 a random subsample of the CGC reference set, the same size as the set of CGC neighbors of that gene. We  
435 repeat this random sampling 100 times and plot the mean percent significant findings value. For reference,  
436 Supplementary Figure 3-A and 3-D contains versions of Figure 1-A and 1-C, where all CGC genes (i.e.,  
437 with and without CGC neighbors) are plotted and mutual exclusivities are checked between all CGC pairs,  
438 as it was done in van de Haar *et al.* [2019].

439 In Figure 1-B, we observe a reduced correlation when network information is included (Pearson  
440 correlation  $-0.4$ , p-value  $4.91e - 4$ ). We also run DISCOVER Strat where stratification is based on  
441 CMS subtypes [Guinney *et al.*, 2015]. We plot these results in Fig 1-C where we again ensure comparability  
442 with Fig 1-D where both subtype and network information are considered. Comparing Figure 1-A and  
443 Figure 1-C, we verify the findings of van de Haar *et al.*, although with less significance in correlation  
444 difference (Pearson correlation  $-0.73$ , p-value  $1.1e - 13$ ). It should be noted that the subtype stratification  
445 inherently causes an overall decrease in percent significant findings, not specific to genes with low MLA.  
446 On the contrary the idea of ME corrections through network incorporation, materialized in the comparison  
447 of Figure 1-A and Figure 1-B, inherently leads to an increase in percent significant findings. Most of the  
448 decreases occur in genes with small number of CGC neighbors. When we compare Fig 1-D to Fig 1-B, the  
449 decrease in correlation from  $-0.4$  to  $-0.36$  indicates that including subtype information is still useful when  
450 used on top of network-based corrections we propose.

451 Supplementary Figure 2 shows the results of the same analysis repeated for BRCA. Similar to the results  
452 that we obtain for COADREAD data, including network information reduces the correlation between  
453 MLA and ME detection rate (Supp Fig 2-A vs 2-B). The magnitude of reduction is even more significant  
454 than what we observe for COADREAD data (Pearson correlation -0.93 vs -0.27). Interestingly, including  
455 subtype information results in a very slight decrease in correlation coefficient (-0.93 to -0.92)(Supp Fig 2-A  
456 vs Supp Fig 2-C) as opposed to what we observe for COADREAD. We again observe that including subtype  
457 information on top of network information results in a negligible decrease in correlation. (Supp Fig 2-B vs  
458 2-D). This difference in the effect of including subtype information for BRCA and COADREAD datasets  
459 could be related to the average tumor mutation load of subtypes. BRCA subtypes have comparable average  
460 TML values (Her2: 146, LumA:65, LumB: 71, Normal: 55) whereas the CMS1 subtype in COADREAD  
461 has a dramatically larger average TML value compared to the other subtypes of COADREAD (CMS1:  
462 1387, CMS2:93, CMS3: 272, CMS4: 212) We repeat the same analysis with the other ME detection  
463 methods as well as for other cancer types when  $t$  is set to 20 (Supplementary Figures 1-8). We observe  
464 that the percent significant finding values can vary remarkably across the tumor types. Compared to other  
465 cancer types, we observe smaller percent significant findings for LUSC (Supplementary Figs 5 A,D,G).  
466 Similarly, very few pairs have percent significance value  $\geq 20$  when we consider network information in  
467 LUSC (Supplementary Figs 5 C,F,I). On the contrary, we observe many pairs with large percent significant  
468 values for CGC-CGC neighbors in UCEC data. This is particularly true for DISCOVER and WeXT results  
469 (Supplementary Figs 8 C-I).

470 When we consider the correlation between MLA and percent significant values, we observe that adding  
471 network information decreases the correlation coefficient values for all cancer types and for all ME detection  
472 methods except for Fisher's Exact Test. Fisher's Exact Test results show an increased correlation with the  
473 addition of network information for LUSC and SKCM (Supplementary Fig. 5-6 D-F). Also, the correlation  
474 coefficient can not be computed for LUAD and STAD since Fisher's Exact Test gives a value of 0 for the  
475 percent significant findings of all considered genes (Supplementary Fig. 4-7 D-F). Another interesting  
476 observation is the variance in magnitude of decrease in correlation values across different tumor types.  
477 In particular, we observe a smaller decrease in correlation values for LUAD compared to other cancer  
478 types. The analogous results are also available for  $t = 5$  setting (Supplementary Figs 9-15). For all the  
479 cancer types, the correlation between MLA values and percent significant findings decreases and becomes  
480 non-significant for most cases.

481 We should also note that the majority of CGC genes have only one neighbor within the data setting of  
482 the cancer type under consideration. This leads to percentage significant findings of either 0 or 1 in many  
483 cases simply because these are the only possible values; for COADREAD see 1-B and 1-D where 41 out  
484 of 74 genes under study have only one CGC neighbor in the COADREAD data settings. To avoid any  
485 such possible biases, we repeat the same evaluations after filtering out those CGC genes with only one  
486 neighbor. The evaluations still provide significant decreases in correlation coefficient values analogous to  
487 the decreases observed in 1-B as compared to 1-A and 1-D as compared to 1-C. For detailed results, see  
488 Supplementary Figures 17-24 for  $t = 20$  and Supplementary Figures 25-32 for  $t = 5$ .

489 Individual genes of interest are those that have increased percent significant findings when network  
490 neighborhood information is incorporated while at the same have significant number of CGC neighbors.  
491 More specifically, for the former constraint, we identify the CGC genes with at least 0.1 increase in  
492 percentage of significant findings value of WeXT, DISCOVER and MEMo when the network information  
493 is included as opposed to the scenario when it is not (e.g. for COADREAD, Figure 1-A vs Figure 1-  
494 B). We choose these three ME methods since they are top performers based on the defined metrics in

495 Section 2.1 . For STAD, SKCM and UCEC, since MEMo results are unavailable, we only consider  
496 WeXT and DISCOVER results. For the second constraint, we include the CGC genes with at least 3 CGC  
497 neighbors. For COADREAD, this selection procedure results in four genes: EP300, CREBBP, NCOA2  
498 and NCOR2. Among these, EP300 is a well-known tumor suppressor in epithelial cancer types including  
499 COADREAD [Gayther et al., 2000]. For BRCA, the only identified gene is PIK3R1. PIK3R1 is found to  
500 be significantly mutually exclusive with PIK3CA and SPEN based on both WeXT, DISCOVER and MEMo  
501 results. PIK3R1 and PIK3CA are members of the PI3K pathway and their mutual exclusivity has been  
502 previously established in the literature [Chen et al., 2018]. For LUAD, PTPRB is the only identified gene  
503 and is found to be mutually exclusive with EGFR, a well-known oncogene in non-small cell lung cancer  
504 [Bethune et al., 2010]. The set of identified genes for STAD are NCOA2, NCOR2 and CREBBP; all of  
505 which are found to be mutually exclusive with TP53. For SKCM, we identify ERBB4, RAC1, EP300 and  
506 ITK. ERBB4 is a well-known oncogene in skin cancer and found to be mutually exclusive with ERBB2  
507 [Prickett et al., 2009, Nielsen et al., 2014]. ERBB2 and ERBB4 indeed belong to the same family (i.e.  
508 ErbB family of receptor tyrosine kinases) and form a heterodimer receptor for Heparin-binding EGF-like  
509 growth factor (HB-EGF) [Iwamoto et al., 2017]. RAC1 mutation P29S is an established driver in melanoma  
510 [Jiang et al., 2018]. RAC1 is found to be mutually exclusive with MYH9, a tumor suppressor in melanoma  
511 [Singh et al., 2020]. Lastly, ITK has been shown to be an oncogene in melanoma [Carson et al., 2015]. For  
512 UCEC, we identify 33 genes in total. Among these, KIT and PTEN have established roles in UCEC cancer  
513 development [Chang et al., 2015, Wang et al., 2020]. Moreover, PTEN is found to be strongly mutually  
514 exclusive with SPOP, whose mutations are also associated with endometrial cancer [Clark and Burleson,  
515 2020]. Lastly, for BLCA and LUSC, no gene satisfies the abovementioned criteria. Overall these results  
516 suggest that the CGC genes that show increased ME with network incorporation as well as their mutually  
517 exclusive partner genes often have established roles in the development of the particular cancer type.

### 518 3.5 ME Evaluations Based on Corrections via TSN

519 We first provide our ME evaluations with respect to the metrics defined in Section 2.1 by replacing the  
520 non-specific networks with TSNs. We provide two types of comparisons; one where we compare  $TSN_{0.5}$   
521 with the original non-tissue specific Intact network and one where results of  $TSN_{0.5}$  are compared against  
522  $TSN_0$ . We do the latter to avoid artifacts that may be introduced due to the fact that some genes in the  
523 original Intact network might be simply missing from even  $TSN_0$  since they may be nonexistent in the  
524 GTEX database. For the BLCA dataset, comparing the F1 scores of the ME methods under  $TSN_0$  and  
525  $TSN_{0.5}$  settings, we observe that the scores of all methods are higher for the latter network. The largest  
526 percent increase of 10% is observed for WeXT when the control group is  $\mathcal{X}_1$ . Similarly, the largest percent  
527 increase of 12% is observed for MEMo when the control group is  $\mathcal{X}_2$ . On the other hand, when we compare  
528 the scores of  $TSN_0$  against the original network, the differences are negligible. The next largest difference  
529 between the F1 scores obtained under under  $TSN_{0.5}$  as compared to  $TSN_0$  is observed in STAD where we  
530 see a 7% increase in DISCOVER's score for  $\mathcal{X}_1$ , and a 10% increase in WeXT's score for  $\mathcal{X}_2$ . For the rest  
531 of the cancer types under study, for LUSC and UCEC we observe slight increase in performances of all the  
532 ME methods comparing the metrics under  $TSN_{0.5}$  against  $TSN_0$ . For COADREAD, BRCA and SKCM  
533 we observe both increases and decreases in performances but the differences are almost negligible; see  
534 Supplementary Tables 50-81 for detailed results.

535 Figure 2 compares the ROC curves of CGC gene pairs and non-CGC gene pairs for COADREAD  
536 data where mutual exclusivities are estimated with DISCOVER, DISCOVER Strat, Fisher's Exact Test,  
537 MEGSA, MEMo and WeXT with  $t = 20$ . We observe that all the ME methods estimate stronger mutual  
538 exclusivities for tissue-specific CGC gene pairs compared to non-tissue-specific CGC gene pairs since

539 AUROCs are greater than 0.5. Additionally, we observe much smaller AUROCs for the control group  
540 where we repeat the same analysis with non-CGC gene pairs. Analogous results are available for the other  
541 cancer types where both the positive and negative set contains at least 10 number of pairs when  $t$  is set to  
542 20. (Supplementary Figs. 33-35). We observe a similar result for SKCM where CGC pairs result in larger  
543 AUROCs compared to non-CGC pairs for all ME methods (Supplementary Fig. 34). We observe a steep  
544 increase in the ROC curves plotted for MEGSA results. This is due to the utilized likelihood ratio test  
545 that results in a p-value of 0.5 when the likelihood values are equal to each other. For UCEC, we see a  
546 significant difference between the ROC curves of CGC-pairs vs non-CGC pairs for Fisher's Exact Test and  
547 MEGSA; whereas the corresponding difference is negligible for DISCOVER and WeXT.

## 4 CASE STUDY

548 Apart from the defined network-centric ME evaluation framework, we discuss a case study where we  
549 assess whether mutual exclusivities estimated by the considered ME methods improve the performance  
550 of driver identification methods that utilize mutual exclusivity information. To this end, we compare the  
551 original version of MEXCOWalk with its alternatives where mutual exclusivity estimates are provided  
552 by the employed ME methods. Assuming that  $g_i$  and  $g_j$  genes are mutated in patient sets  $S_i$  and  $S_j$ ,  
553 respectively; MEXCOWalk simply computes the mutual exclusivity between these two genes with the  
554 following formula:  $|S_i \cup S_j| / (|S_i| + |S_j|)$ . MEXCOWalk uses the estimated mutual exclusivity values as  
555 part of edge weights. As such, to utilize the p-values output by ME detection methods in MEXCOWalk, we  
556 first compute  $-\log(\text{p-value})$  and then convert the resulting values between 0 and 1. To this end, we replace  
557 all  $-\log(\text{p-value})$ 's larger than 10 with 1. We then find the maximum  $-\log(\text{p-value})$  less than 10 and divide  
558 all other  $-\log(\text{p-value})$ 's with this value. The reason why we set a threshold for finding the maximum is  
559 the large differences across the smallest p-values output by different ME methods. For instance, WeXT  
560 outputs a very large range of p-values and if we use the smallest p-value to scale, all other  $-\log(\text{p-value})$ s  
561 will be converted to values that are very close to 0. In the original MEXCOWalk study, a threshold of  
562 0.7 is applied to ME values such that all values  $\leq 0.7$  are clamped to 0. This conversion is equivalent  
563 to removing those edges from the network since the edge weights include a multiplicative term for ME  
564 values. We find that the removal of these edges correspond to a 0.035 percent reduction in graph density.  
565 For the current analysis, we determine the threshold value for each ME detection method to achieve the  
566 same percent density reduction in the graph. Figure 3 shows the number of recovered CGC genes for fixed  
567 output gene sizes from 100 to 2500 as a ROC curve for original MEXCOWalk as well as for versions of  
568 MEXCOWalk where mutual exclusivity values are estimated with DISCOVER, Fisher's Exact Test and  
569 WeXT, respectively. We observe that MEXCOWalk with WeXT's ME values results in the best AUROC  
570 value for COADREAD. Supplementary Figure 36 shows the analogous results for the other cancer types.  
571 For, LUSC, STAD and UCEC, MEXCOWalk with DISCOVER gives the best AUROC whereas for BLCA,  
572 LUAD and SKCM MEXCOWalk with Fisher's Exact Test performs the best. An important observation is  
573 the worse performance of MEXCOWalk with Fisher's Exact Test compared to the original MEXCOWalk  
574 for COADREAD, STAD and UCEC. As such, using Fisher's Exact Test in place of MEXCOWalk's original  
575 ME values does have the potential to decrease the performance whereas for the other ME methods we  
576 do not observe such a risk. Note that for these analysis we employ  $t = 5$  since  $t = 20$  filtering does not  
577 provide enough number of genes to be evaluated.

## 5 DISCUSSION

578 It is important to investigate whether the employment of an interaction network within our ME evaluation  
579 framework causes any ascertainment bias in the findings and to elaborate on how any such potential  
580 bias is mediated within the framework. It is established that known cancer genes have larger number of  
581 interactions compared to other genes in the network [Hou and Ma, 2014]. This implies a potential bias  
582 that needs to be resolved in cancer driver gene identification methods employing interaction network data.  
583 Such a bias is less of a problem for the current study, since our aim is not to identify novel cancer driver  
584 genes but to utilize the interaction network and known cancer genes to form a ground truth of mutually  
585 exclusive interactions for evaluating existing ME methods. On the contrary, the fact that most known cancer  
586 genes have well-characterized interactions in the network provides a benefit for our work as it supports  
587 the confidence of our true positive examples. Additionally, our framework makes use of not only genes  
588 from the reference set  $\mathcal{S}$  but also genes not in  $\mathcal{S}$  to create random controls. Nevertheless, the fact that some  
589 known cancer genes have significantly larger number of interactions compared to other known cancer  
590 genes could lead to a bias. For instance, for our analysis of the COADREAD data ( $t = 20$ ,  $\mathcal{S} = CGC$ ),  
591 there are 74 CGC genes among which five CGC genes have more than ten CGC neighbors whereas 41  
592 have exactly one CGC neighbor. This could lead to a bias as CGC genes with large number of CGC  
593 neighbors contribute to the aggregate statistics and metrics much more than those CGC genes with small  
594 number of CGC neighbors. To mediate for this bias, our framework includes additional results where all  
595 the statistics and the traditional measures such as the F1 score are calculated in a degree-normalized way  
596 for each gene and the gene-level results are then aggregated by taking an average across the genes. These  
597 results are available in the Supplementary Document; see Table S12, S24, S36, and S48. To summarize,  
598 the degree-normalized results are in agreement with those of the previous settings in almost all the cases in  
599 terms of ranking based on F1 score.

600 Another important point worth emphasizing is that apart from the aggregate statistics provided in the  
601 previous sections as part of the metrics for the network-centric ME evaluations, our proposed framework  
602 also provides analogous statistics at the gene-level as well. Such statistics may in fact be of more interest to  
603 cancer biologists than the aggregate statistics in certain cases. Several interesting observations can be made  
604 through an inspection of these gene-level evaluations, especially for the settings where the conventionally  
605 defined F1 score fails in quantifying ME. Genes with low MLA comprise an example setting, where TP53  
606 is a leading member. Consider the case of TP53 in COADREAD evaluations for instance. With respect to  
607 the degree-normalized setting, the values of precision, sensitivity,  $\text{precision}_{strict}$  and  $\text{sensitivity}_{strict}$  for  
608 WeXT are respectively 0.5, 1, 0.25, 0.25 which gives rise to an F1 score of 0.66 and  $F1_{strict}$  score of 0.25.  
609 On the other hand, MEMo provides the same precision, sensitivity and F1 scores as WeXT whereas its  
610  $\text{precision}_{strict}$ ,  $\text{sensitivity}_{strict}$  and  $F1_{strict}$  scores are all 0. To summarize, although the inspection of the F1  
611 scores does not provide a distinction between the two results, an inspection of the  $F1_{strict}$  scores establishes  
612 that MEMo is worse than WeXT in this setting. We note that the advantages of inspections based on the  
613 strict definitions of the metrics rather than the conventional ones are also apparent in the aggregate analysis  
614 as well. In addition to the COADREAD evaluations shown in Table2, BRCA also contains an example  
615 instance where the conventional and the strict versions of the metrics provide different conclusions; see  
616 Table S3-b. In terms of the F1 scores WeXT and DISCOVER Strat obtain very close values with WeXT  
617 providing slightly better results. However comparing  $F1_{strict}$  scores reveals DISCOVER Strat a better ME  
618 test candidate for this instance. Also, overall we observe that MEMo's performance gets severely affected  
619 when the strict versions of the metrics are employed.

620 Lastly, our robustness analysis results reveal some suggestions for potential users of our framework.  
621 We recommend using a p-value threshold greater than 0.05 as lower threshold values are too stringent  
622 and lead to too few predicted positives. Regarding *robustness\_iterations*, we tested values both smaller  
623 than and higher than the default value of 100 for COADREAD evaluations: 5, 50, 100, 300 and 500.  
624 We repeated each experiment 20 times and calculated the standard deviation of the obtained set of F1  
625 and  $F1_{strict}$  scores. For the majority of the cases, we observe a large decrease in the standard deviation  
626 values when *robustness\_iterations* is increased from 5 to 50. (Table S82). This analysis suggests that the  
627 *robustness\_iterations* should be set to a at least 50. Lastly, we observe that different PPI networks can  
628 lead to large differences in both the F1/ $F1_{strict}$  scores and the ranking of the methods. As such, exploring  
629 different PPI sources would be beneficial.

## 6 CONCLUSION

630 We propose a network-centric framework to evaluate pairwise mutual exclusivity findings reported by  
631 different ME algorithms. The first component of our framework consists of useful definitions of statistics  
632 employed in the network-centric ME evaluations. We observe that for the majority of the cancer types under  
633 study WeXT outperforms the other methods in terms of F1 score measured with respect to appropriately  
634 defined control groups. In half of the cancer types DISCOVER and in the other half MEMo perform as the  
635 second best methods. When comparing different cancer types we observe that BRCA and COADREAD  
636 are among the top two types leading to maximum F1 scores with at least one of the ME methods providing  
637 a score greater than 0.5. We note that DISCOVER Strat is only applicable in two cancer types among a  
638 total of eight since these are the only cancer types with well-defined subtypes. Furthermore, among these  
639 two cancer types, DISCOVER Strat outperforms original DISCOVER algorithm in BRCA, whereas it  
640 is the second worst method after Fisher's Exact Test in COADREAD. This is noteworthy since van de  
641 Haar *et al.* propose subtype stratification as employed by DISCOVER Strat as a way to emphasize true  
642 mutual exclusivity by reducing mutation load confounding [van de Haar et al., 2019]. We also observe  
643 that Fisher's exact test and MEGSA are more conservative compared to DISCOVER and WeXT, where  
644 from the latter group, WeXT outputs notably larger number of significant p-values. The second component  
645 of our framework evaluates ME tests by comparing two types of measures obtained with and without  
646 network information. First measure is with respect to the percent significant findings of mutually exclusive  
647 gene pairs, whereas the second is based on MLC values. In most of the cancer types and for most of the  
648 genes we observe an increase with respect to the former whereas a decrease with respect to the latter  
649 measure. Finally, we repeat the same analysis by considering TSNs in the network-centric framework.  
650 Considerable improvements achieved due to the use of TSNs as opposed tissue nonspecific interaction  
651 network are only observed for BLCA and STAD datasets. A more detailed analysis in terms of comparing  
652 ROCs of CGC gene pairs and non-CGC gene pairs on cancer types with considerable number of tissue-  
653 specific gene pairs indicate the advantages of employing tissue specificity in detecting mutual exclusivity in  
654 COADREAD, SKCM, and UCEC. Finally we extend out network-centric evaluation framework to assess  
655 whether including network knowledge reduces the mutation load confounding problem.

656 As noted earlier the proposed framework is intended for the network-centric evaluations of mutual  
657 exclusivities of pairs of genes rather than groups of genes. Such a choice stems from the fact that the mutual  
658 exclusivities are commonly made use of in driver gene/module identification algorithms which mostly  
659 employ pairwise mutual exclusivities. Furthermore the extensive evaluation settings proposed, the number  
660 of ME methods under study and their own computational requirements, and the potentially exponential  
661 computational complexity inherent in handling groups of genes limits the scope of the current study to



662 evaluations of pairwise ME scorings. Nonetheless most statistical ME methods are capable of providing  
663 ME results for groups of genes as well. Regarding the ME tests considered in this study, the main ME test  
664 provided by DISCOVER is based on a pairwise test definition but it also extends the definition for possible  
665 use in quantifying the ME of a group of genes, although the experiments involving the latter are based  
666 only on simulation data. The remaining tests MEGSA, MEMo, and WeXT are all ME tests specifically  
667 designed for groups of genes. An important direction for future work is to design a suitable extension of the  
668 proposed network-centric framework to evaluate the results of ME tests on groups of genes. Design choices  
669 relevant for such an extension would involve an appropriate and computationally efficient definition of the  
670 reference groups of genes analogous to a pair of interacting genes from the set  $\mathcal{S}$  in the current setting and  
671 the definitions of control groups analogous to  $\mathcal{X}_1$  and  $\mathcal{X}_2$ .

## 7 ADDITIONAL REQUIREMENTS

672 For additional requirements for specific article types and further information please refer to Author  
673 Guidelines.

## CONFLICT OF INTEREST STATEMENT

674 The authors declare that the research was conducted in the absence of any commercial or financial  
675 relationships that could be construed as a potential conflict of interest.

## AUTHOR CONTRIBUTIONS

676 Authors names are written in alphabetical order. C. E. and H.K. conceived the idea and supervised the study.  
677 R. A implemented the code and performed the initial experiments. C. Y. and A.H. repeated the experiments  
678 for other cancer types. All authors contributed to the preparation of the manuscript. All authors read and  
679 approved the manuscript.

## FUNDING

680 This work was supported by the Scientific and Technological Research Council of Turkey [grant number  
681 117E879 to H.K. and C.E.].

## ACKNOWLEDGMENTS

682 We thank Joris van de Haar for sharing us the data used in van de Haar et al. [2019].

## SUPPLEMENTAL DATA

683 Supplementary Material should be uploaded separately on submission, if there are Supplementary Figures,  
684 please include the caption in the same file as the figure. LaTeX Supplementary Material templates can be  
685 found in the Frontiers LaTeX folder.

## DATA AVAILABILITY STATEMENT

686 The datasets analyzed for this study can be found on the following repository: [https://github.com/abu-](https://github.com/abu-compbio/NetCentric)  
687 [compbio/NetCentric](https://github.com/abu-compbio/NetCentric).

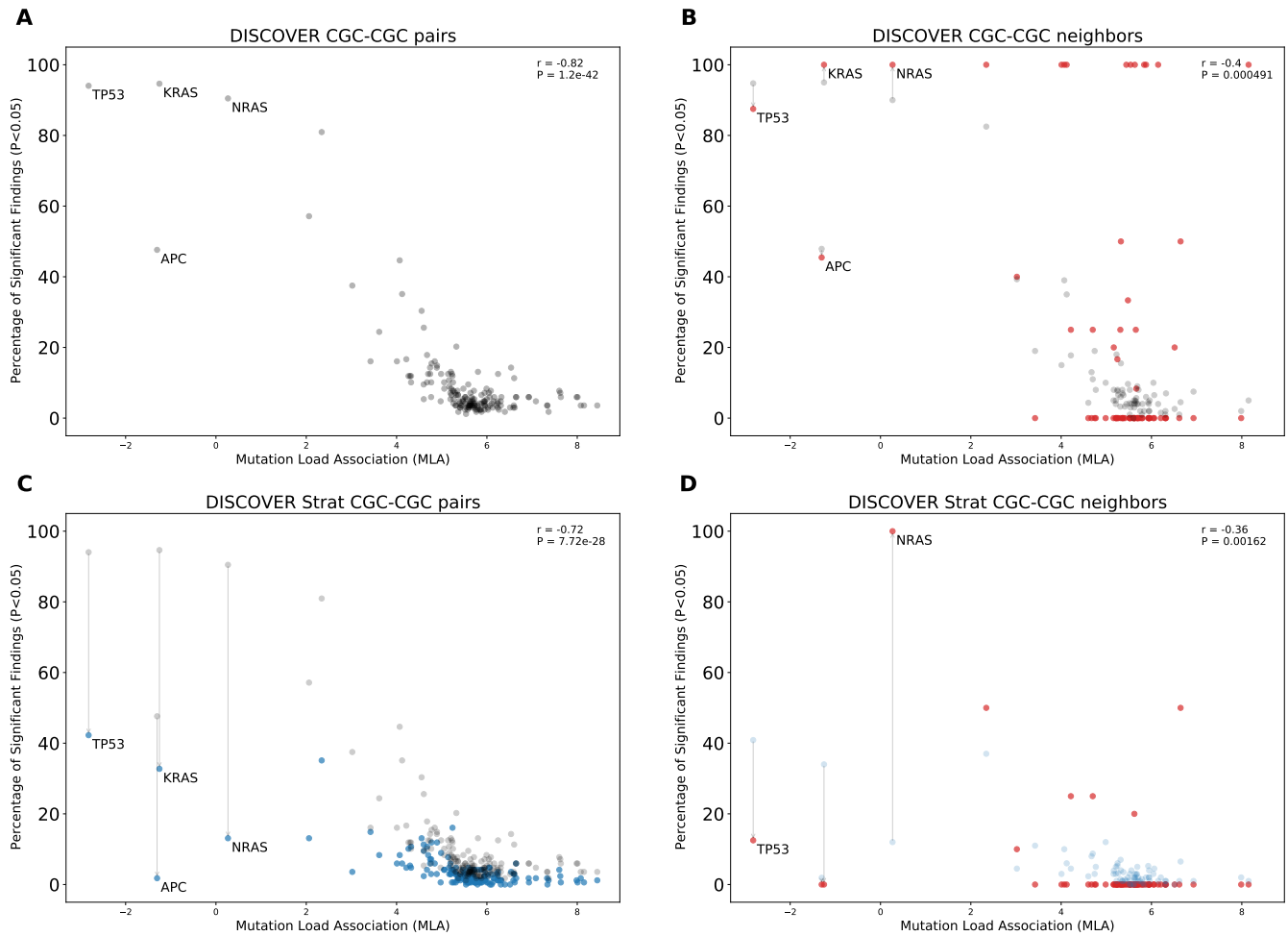
## REFERENCES

- 688 Rafsan Ahmed, Ilyes Baali, Cesim Erten, et al. MEXCOWalk: mutual exclusion and coverage based  
689 random walk to identify cancer modules. *Bioinformatics*, 36(3):872–879, 08 2019. ISSN 1367-4803.  
690 doi: 10.1093/bioinformatics/btz655.
- 691 Ilyes Baali, Cesim Erten, and Hilal Kazan. Driveways: A method for identifying possibly overlapping  
692 driver pathways in cancer. *Sci Rep*, 10, 2020. doi: 10.1101/2020.04.01.015388.
- 693 Özgün Babur, Mithat Gönen, Bülent Arman Aksoy, et al. Systematic identification of cancer driving  
694 signaling pathways based on mutual exclusivity of genomic alterations. *Genome Biology*, 16(1):45, Feb  
695 2015.
- 696 Matthew Bailey, Collin Tokheim, Eduard Porta, et al. Comprehensive characterization of cancer driver  
697 genes and mutations. *Cell*, 173:371–385.e18, 04 2018. doi: 10.1016/j.cell.2018.02.060.
- 698 Rebecca Basso, Dorit Hochbaum, and Fabio Vandin. Efficient algorithms to discover alterations with  
699 complementary functional association in cancer. *PLOS Comput Biol*, 15, 2019. doi: 10.1371/journal.  
700 pcbi.1006802.
- 701 G Bethune, D Bethune, N Ridgway, et al. Epidermal growth factor receptor (egfr) in lung cancer: an  
702 overview and update. *J of Thoracic Disease*, 2:48–51, 2010.
- 703 Sander Canisius, Lodewyk Wessels, and John W. M. Martens. A novel independence test for somatic  
704 alterations in cancer shows that biology drives mutual exclusivity but chance explains most co-occurrence.  
705 *Genome Biology*, 17(261):1–17, 2016. ISSN 1474-760X. doi: 10.1186/s13059-016-1114-x.
- 706 C. C Carson, S. J. Moschos, S. N. Edmiston, et al. Ii2 inducible t-cell kinase, a novel therapeutic target in  
707 melanoma. *Clin Cancer Res*, 21(9):55–65, 2015. doi: 10.1158/1078-0432.CCR-14-1826.
- 708 Ethan Cerami, Jianjiong Gao, Ugur Dogrusoz, et al. The cBio cancer genomics portal: An open platform  
709 for exploring multidimensional cancer genomics data. *Cancer Discovery*, 2(5):401–404, 2012. ISSN  
710 2159-8274. doi: 10.1158/2159-8290.CD-12-0095.
- 711 Shih-Wen Chang, Wan-Ru Chao, Alexandra Ruan, et al. A promising hypothesis of c-kit methylation/  
712 expression paradox in c-kit(+) squamous cell carcinoma of uterine cervix. *Diagnostic Pathology*, 10(1),  
713 2015. doi: 10.1186/s13000-015-0438-2.
- 714 Li Chen, Liu Yang, Ling Yao, et al. Characterization of pik3ca and pik3r1 somatic mutations in chinese  
715 breast cancer patients. *Nat Comm*, 9(1), 2018. doi: 10.1038/s41467-018-03867-9.
- 716 DT Cheng, TN Mitchell, A Zehir, et al. Msk-impact: A hybridization capture- based next-generation  
717 sequencing clinical assay for solid tumor molecular oncology. *J Mol Diagn*, 17(3):251–264, 2015.
- 718 G. Ciriello, E. Cerami, C. Sander, et al. Mutual exclusivity analysis identifies oncogenic network modules.  
719 *Genome Res.*, 22(2):398–406, Feb 2012.
- 720 A Clark and M Burleson. Spop and cancer: a systematic review. *American J of Cancer Res*, 10(3):704–726,  
721 2020.
- 722 Simona Constantinescu, Ewa Szczurek, Pejman Mohammadi, et al. TiMEx: a waiting time model for  
723 mutually exclusive cancer alterations. *Bioinformatics*, 32(7):968–975, 07 2015.
- 724 J. Das and H. Yu. HINT: High-quality protein interactomes and their applications in understanding human  
725 disease. *BMC Systems Biology*, 6:92, 2012.
- 726 Yulan Deng, Shangyi Luo, Chunyu Deng, et al. Identifying mutual exclusivity across cancer genomes:  
727 computational approaches to discover genetic interaction and reveal tumor vulnerability. *Brief in*  
728 *Bionform*, 2017. doi: 10.1093/bib/bbx109.
- 729 Christos M. Dimitrakopoulos and Niko Beerenwinkel. Computational approaches for the identification of  
730 cancer genes and pathways. *Wiley Interdiscip Rev Syst Biol Med*, 9(1):e1364, 2017. ISSN 1939-5094.  
731 doi: 10.1002/wsbm.1364.

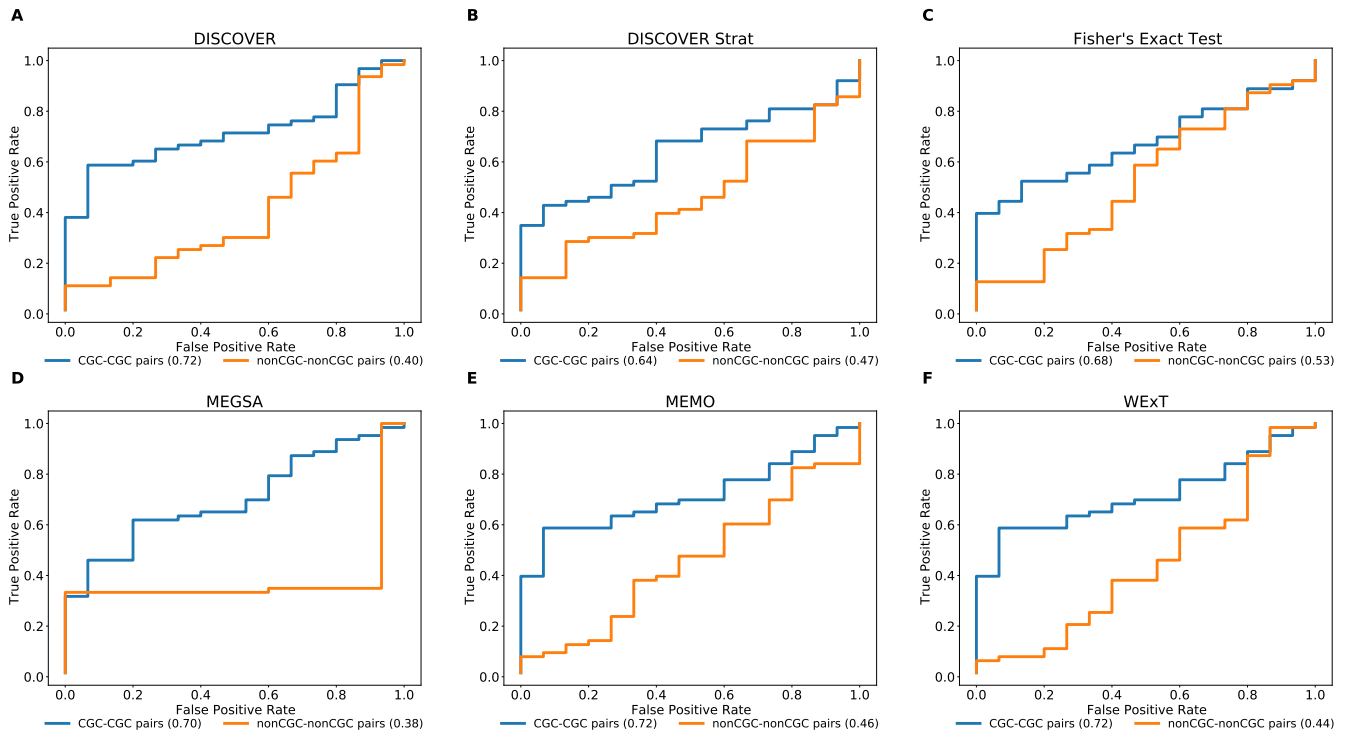
- 732 S.A. Forbes, D. Beare, H. Boutselakis, et al. Cosmic: somatic cancer genetics at high-resolution. *NAR*, 45:  
733 D777–D783, 2017.
- 734 Jianjiong Gao, Bülent Arman Aksoy, Ugur Dogrusoz, et al. Integrative analysis of complex cancer  
735 genomics and clinical profiles using the cbiportal. *Science Signaling*, 6(269):p11–p11, 2013. ISSN  
736 1945-0877. doi: 10.1126/scisignal.2004088.
- 737 SA Gayther, SJ Batley, L Linger, et al. Mutations truncating the ep300 acetylase in human cancers. *Nat*  
738 *Genet*, 24:300–303, 2000.
- 739 GTEX Consortium. The gtex consortium atlas of genetic regulatory effects across human tissues. *Science*,  
740 369(6509):1318–1330, 2020. ISSN 0036-8075. doi: 10.1126/science.aaz1776.
- 741 Justin Guinney, Rodrigo Dienstmann, Xin Wang, et al. The consensus molecular subtypes of colorectal  
742 cancer. *Nat Med*, 21, 10 2015. doi: 10.1038/nm.3967.
- 743 Jack P Hou and Jian Ma. Dawnrank: discovering personalized driver genes in cancer. *Genome Medicine*, 6  
744 (56):1–16, 2014.
- 745 Xing Hua, Paula Hyland, Jing Huang, et al. Megsa: A powerful and flexible framework for analyzing  
746 mutual exclusivity of tumor mutations. *The American J of Human Genetics*, 98, 02 2016. doi:  
747 10.1016/j.ajhg.2015.12.021.
- 748 Ryo Iwamoto, Naoki Mine, Hiroto Mizushima, et al. Erbb1 and erbb4 generate opposing signals regulating  
749 mesenchymal cell proliferation during valvulogenesis. *Development*, 144(8), 2017. doi: 10.1242/dev.  
750 152710.
- 751 Ze-Bin Jiang, Bing-Qiang Ma, Shao-Guang Liu, et al. mir-365 regulates liver cancer stem cells via rac1  
752 pathway. *Molecular Carcinogenesis*, 58(1):55–65, 2018. doi: 10.1002/mc.22906.
- 753 Yoo-Ah Kim, Dong-Yeon Cho, and T. M. Przytycka. MEMCover: integrated analysis of mutual exclusivity  
754 and functional network reveals dysregulated pathways across multiple cancer types. *Bioinformatics*, 31  
755 (12):i284–i292, 2015.
- 756 Yoo-Ah Kim, Sanna Madan, and Teresa M Przytycka. WeSME: uncovering mutual exclusivity of cancer  
757 drivers and beyond. *Bioinformatics*, 33(6):814–821, 05 2017. ISSN 1367-4803.
- 758 M. D. Leiserson, D. Blokh, R. Sharan, et al. Simultaneous identification of multiple driver pathways in  
759 cancer. *PLoS Comput. Biol.*, 9(5):e1003054, 2013a.
- 760 M. D. M. Leiserson, H. T. Wu, F. Vandin, et al. Comet: A statistical approach to identify combinations  
761 of mutually exclusive alterations in cancer. *Genome Biol*, 16, 2015a. ISSN 1474-7596. doi: 10.1186/  
762 s13059-015-0700-7.
- 763 Mark Leiserson, Matthew Reyna, and Ben Raphael. A weighted exact test for mutually exclusive mutations  
764 in cancer. *Bioinformatics*, 32:i736–i745, 07 2016. doi: 10.1093/bioinformatics/btw462.
- 765 Mark D. M. Leiserson, Dima Blokh, Roded Sharan, et al. Simultaneous identification of multiple driver  
766 pathways in cancer. *PLOS Comput Biol*, 9(5):1–15, 05 2013b. doi: 10.1371/journal.pcbi.1003054.
- 767 Mark D. M. Leiserson, Fabio Vandin, Hsin-Ta Wu, et al. Pan-cancer network analysis identifies  
768 combinations of rare somatic mutations across pathways and protein complexes. *Nat Genet*, 47(2):  
769 106–114, 2015b. ISSN 1546-1718. doi: 10.1038/ng.3168.
- 770 Sisheng Liu, Jinpeng Liu, Yanqi Xie, et al. MEScan: a powerful statistical framework for genome-scale  
771 mutual exclusivity analysis of cancer mutations. *Bioinformatics*, 11 2020.
- 772 Katja Luck, Dae-Kyum Kim, Luke Lambourne, et al. A reference map of the human binary protein  
773 interactome. *Nature*, 580:402–408, 04 2020. doi: 10.1038/s41586-020-2188-x.
- 774 F Martínez-Jiménez, F Muinos, I Senis, et al. A compendium of mutational cancer driver genes. *Nat Rev*  
775 *Cancer*, 20:555–572, 2020.

- 776 Trine O. Nielsen, Steen S. Poulsen, Fabrice Journe, et al. Her4 and its cytoplasmic isoforms are associated  
777 with progression-free survival of malignant melanoma. *Melanoma Research*, 24(1):88–91, 2014. doi:  
778 10.1097/cmr.0000000000000040.
- 779 Sandra Orchard, Mais Ammari, Bruno Aranda, et al. The mintact project—intact as a common curation  
780 platform for 11 molecular interaction databases. *NAR*, 42(D):D358—63, January 2014. ISSN 0305-1048.  
781 doi: 10.1093/nar/gkt1115.
- 782 P Porras, E Barrera, A Bridge, et al. Towards a unified open access dataset of molecular interactions. *Nat*  
783 *Comm*, 11(6144), 2020. doi: 10.1038/s41467-020-19942-z.
- 784 Todd D Prickett, Neena S Agrawal, Xiaomu Wei, et al. Analysis of the tyrosine kinome in melanoma  
785 reveals recurrent mutations in erbb4. *Nat Genet*, 41(10):1127–1132, 2009. doi: 10.1038/ng.438.
- 786 Satyendra Kumar Singh, Sunita Sinha, et al. Myh9 suppresses melanoma tumorigenesis, metastasis and  
787 regulates tumor microenvironment. *Medical Oncology*, 37(10), 2020. doi: 10.1007/s12032-020-01413-6.
- 788 Zbyslaw Sondka, Sally Bamford, Charlotte G. Cole, et al. The cosmic cancer gene census: describing  
789 genetic dysfunction across all human cancers. *Nature Reviews Cancer*, 18(11):696–705, 2018. ISSN  
790 1474-1768. doi: 10.1038/s41568-018-0060-1.
- 791 Junrong Song, Wei Peng, and Feng Wang. An entropy-based method for identifying mutual exclusive  
792 driver genes in cancer. *IEEE ACM Trans. Comput. Biol. Bioinform.*, 17(3):758–768, 2020.
- 793 Ewa Szczurek and Niko Beerenwinkel. Modeling mutual exclusivity of cancer mutations. *PLOS Comp Biol*,  
794 10(3):1–12, 03 2014. doi: 10.1371/journal.pcbi.1003503. URL [https://doi.org/10.1371/  
795 journal.pcbi.1003503](https://doi.org/10.1371/journal.pcbi.1003503).
- 796 Lyon D Szklarczyk D, Gable AL et al. String v11: protein–protein association networks with increased  
797 coverage, supporting functional discovery in genome-wide experimental datasets. *NAR*, 47(D1), 2018.  
798 doi: 10.1093/nar/gky1131.
- 799 E Sügis, J Dauvillier, A Leontjeva, et al. Hena, heterogeneous network-based data set for alzheimer’s  
800 disease. *Scientific Data*, 6(151), 2019. doi: 10.1038/s41597-019-0152-0.
- 801 Roman Thomas, Alissa Baker, Ralph Debiasi, et al. High-throughput oncogene mutation profiling in  
802 human cancer. *Nat Genet*, 39:347–51, 04 2007. doi: 10.1038/ng1975.
- 803 Collin Tokheim, Nickolas Papadopoulos, Kenneth Kinzler, et al. Evaluating the evaluation of cancer  
804 driver genes. *Proceedings of the National Academy of Sciences*, 113:201616440, 11 2016. doi:  
805 10.1073/pnas.1616440113.
- 806 Joris van de Haar, Sander Canisius, Michael Yu, et al. Identifying epistasis in cancer genomes: A delicate  
807 affair. *Cell*, 177:1375–1383, 05 2019. doi: 10.1016/j.cell.2019.05.005.
- 808 F. Vandin, E. Upfal, and B. J. Raphael. De novo discovery of mutated driver pathways in cancer. *Genome*  
809 *Res.*, 22(2):375–385, Feb 2012.
- 810 Bert Vogelstein, Nickolas Papadopoulos, Velculescu, et al. Cancer genome landscapes. *Science*, 339(6127):  
811 1546–1558, March 2013.
- 812 Tao Wang, Shasha Ruan, Xiaolu Zhao, et al. Oncovar: an integrated database and analysis platform for  
813 oncogenic driver variants in cancers. *NAR*, 49(D1), 2020. doi: 10.1093/nar/gkaa1033.
- 814 Chen-Hsiang Yeang and Arnold Levine. Combinatorial patterns of somatic gene mutations in cancer.  
815 *FASEB J.*, 22:2605–22, 05 2008. doi: 10.1096/fj.08-108985.
- 816 H. Yu, L. Tardivo, S. Tam, et al. Next-generation sequencing to generate interactome datasets. *Nature*  
817 *Methods*, 8:478–480, 2011.
- 818 J. Zhang and S. Zhang. The discovery of mutated driver pathways in cancer: Models and algorithms.  
819 *IEEE/ACM Transactions on Computational Biology and Bioinformatics*, 15(3):988–998, May 2018.  
820 ISSN 1545-5963. doi: 10.1109/TCBB.2016.2640963.

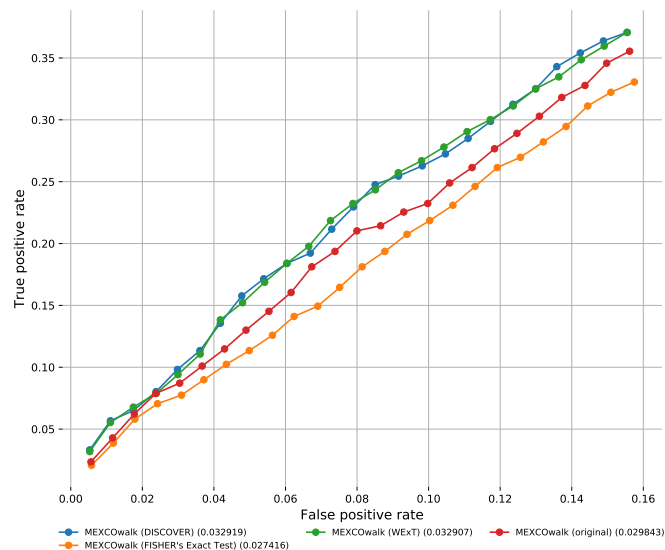
821 Zeyu Zhang, Yaning Yang, Yinsheng Zhou, et al. A forward selection algorithm to identify mutually  
822 exclusive alterations in cancer studies. *J of Human Genetics*, November 2020. ISSN 1434-5161. doi:  
823 10.1038/s10038-020-00870-1.



**Figure 1.** Comparison of mutual exclusivity results of DISCOVER and DISCOVER Strat on COADREAD cohort (498 samples) (A) The scatterplot of percentage significance of mutual exclusivity runs ( $p$ -value; $>0.05$ ) of DISCOVER on COADREAD data where tests are performed between a CGC gene and all other CGC genes. (B) The scatter plot of percentage significance of mutual exclusivity runs of DISCOVER where tests are performed between a CGC gene and its PPI neighbors that are in CGC (red) compared with (A) in gray. (C) The scatterplot of percentage significance of mutual exclusivity runs of DISCOVER Strat where tests are performed between a CGC gene and all other CGC genes (blue) compared with the results from (A) in gray. (D) The scatterplot of percentage significance of mutual exclusivity runs of DISCOVER Strat where tests are performed between a CGC gene and its PPI neighbors that are in CGC (red) compared with (C) in blue.



**Figure 2.** Performance of selected ME tests in terms of discriminating TSN and non-TSN gene pairs based on estimated ME p-values on COADREAD data. Blue curve is plotted with CGC gene pairs and red curve is plotted with non-CGC gene pairs. Mutual exclusivities are estimated with a) DISCOVER, b) DISCOVER Strat, c) Fisher's Exact Test, d) MEGSA, e) MEMO and g) WeXT respectively.



**Figure 3.** The number of recovered CGC genes for the original MEXCOwalk as well as for its modified versions where mutual exclusivity values are estimated with DISCOVER, Fisher's Exact Test and WeXT. COADREAD dataset is used with  $t = 5$  setting. The numbers in parentheses indicate the area under the ROC curve for the corresponding curve.

**Table 1.** Results of network-centric ME evaluation framework with control group  $\mathcal{X}_1$  COADREAD t20 (498 samples, 196 CGC-CGC pairs)

Method	Precision	Sensitivity	F1 Score	Precision <sub>strict</sub>	Sensitivity <sub>strict</sub>	F1 Score <sub>strict</sub>
DISCOVER	0.661	0.220	0.331	0.708	0.183	0.291
DISCOVER Strat	0.727	0.041	0.078	0.727	0.041	0.078
Fisher's Exact Test	0.500	0.031	0.058	0.500	0.031	0.058
MEGSA	0.611	0.056	0.103	0.588	0.051	0.094
MEMO	0.658	0.329	0.439	0.647	0.237	0.347
WExT	0.676	0.403	0.505	0.725	0.329	0.453

**Table 2.** Results of network-centric ME evaluation framework with control group  $\mathcal{X}_2$  COADREAD t20 (498 samples, 107 CGC-CGC pairs)

Method	Precision	Sensitivity	F1 Score	Precision <sub>strict</sub>	Sensitivity <sub>strict</sub>	F1 Score <sub>strict</sub>
DISCOVER	0.537	0.276	0.365	0.579	0.210	0.308
DISCOVER Strat	0.455	0.048	0.086	0.400	0.038	0.069
Fisher's Exact Test	0.444	0.038	0.069	0.375	0.028	0.052
MEGSA	0.571	0.075	0.133	0.538	0.066	0.118
MEMO	0.566	0.388	0.460	0.495	0.215	0.300
WExT	0.575	0.438	0.497	0.596	0.295	0.395

**Table 3.** Results of network-centric ME evaluation framework with control group  $\mathcal{X}_1$  COADREAD t5 (498 samples, 1748 CGC-CGC pairs)

Method	Precision	Sensitivity	F1 Score	Precision <sub>strict</sub>	Sensitivity <sub>strict</sub>	F1 Score <sub>strict</sub>
DISCOVER	0.647	0.052	0.096	0.658	0.046	0.086
DISCOVER Strat	0.618	0.012	0.024	0.618	0.012	0.024
Fisher's Exact Test	0.583	0.008	0.016	0.565	0.007	0.014
WExT	0.645	0.121	0.203	0.668	0.102	0.177

**Table 4.** Results of network-centric ME evaluation framework with control group  $\mathcal{X}_2$  COADREAD t5 (498 samples, 1625 CGC-CGC pairs)

Method	Precision	Sensitivity	F1 Score	Precision <sub>strict</sub>	Sensitivity <sub>strict</sub>	F1 Score <sub>strict</sub>
DISCOVER	0.721	0.052	0.097	0.746	0.048	0.090
DISCOVER Strat	0.641	0.013	0.025	0.641	0.013	0.025
Fisher's Exact Test	0.619	0.008	0.016	0.619	0.008	0.016
WExT	0.670	0.118	0.200	0.712	0.103	0.180