# Combining network-guided GWAS to discover susceptibility mechanisms for breast cancer

Héctor Climente-González[1,2,3],Christine Lonjou[1,2,3],Fabienne Lesueur[1,2,3],GENESIS investigators,Dominique Stoppa-Lyonnet[4,5,6],Nadine Andrieu[1,2,3],Chloé-Agathe Azencott[3,1,2]

Systems biology provides a comprehensive approach to biomarker discovery and biological hypothesis building. It does so by jointly considering the statistical association between a gene and a phenotype, obtained experimentally, and the biological context of each gene, represented as a network. In this work we study the utility of six network methods to discover new biomarkers in GWAS data by searching subnetworks highly associated to a phenotype. We interrogate a familial breast cancer genome-wide association study (GWAS) focused on *BRCA1/2* negative French women. By trading statistical astringency for biological meaningfulness, most network methods get more compelling results than standard SNP- and gene-level analyses, recovering causal subnetworks tightly related to cancer susceptibility. We perform an in-depth benchmarking of the methods with regards to size of the solution subnetwork, their utility as biomarkers, and the stability and the runtime of the methods. Interestingly, a combination of solution subnetworks provided a concise subnetwork of 93 genes, enriched in known familial breast cancer susceptibility genes (*BABAM1*, *BLM*, *CASP8*, *FGFR2*, and *TOX3*, Fisher's exact test P-value $= 7.8 \times 10^{-5}$ ) and more central than average. Additionally, it includes subnetworks of mechanisms related to cancer, like protein folding (*HSPA1A*, *HSPA1B*, and *HSPA1L*) or mitochondrial ribosomes (*MRPS30*, *MRPS31*, *MRPS18B*). We also observed a general dysregulation in the neighborhood of *COPS5*, a gene related to multiple hallmarks of cancer.

## 1 Introduction

In human health, genome-wide association studies (GWAS) aim at quantifying how single-nucleotide polymorphisms (SNPs) predispose to complex diseases, like diabetes or some forms of cancer [1]. To that end, in a typical GWAS thousands of unrelated samples are genotyped: the cases, suffering the disease of interest, and the controls from the general population. Then, a per-SNP statistical association test is conducted (e.g. $\chi^2$). Those SNPs with a P-value lower than a conservative Bonferroni threshold are candidates to further studies in an independent cohort. Once the risk SNPs have been discovered, they can be used for risk assessment, and to deepen our understanding of the disease.

GWAS have successfully identified thousands of variants underlying many common diseases [2]. However, the experimental setting also presents intrinsic challenges. Some of them stem from the high-dimensionality of the problem, as every GWAS to date studies more variants than samples are genotyped. This limits the statistical power of the experiment, as only variants with large and moderate effects can be detected. And it is particularly problematic since the prevailing view is that most genetic architectures involve many variants with small effects [3]. Additionally, to avoid false positives, a conservative multiple test correction is applied, typically Bonferroni. However, Bonferroni is known to be overly conservative when the statistical tests are correlated, as is the case in GWAS [4]. Another open issue is the interpretation of the results, as the functional consequences of most common variants are not well understood. The fact that SNPs in non-coding regions often top the results of GWAS illustrates well this challenge. On top of that, recent large-sampled studies suggest that most of the genome contributes to a degree to any complex trait, in accordance with the infinitesimal model [5]. The omnigenic model [6] explains this by the dense functional

inter-relatedness between genes, influencing each other's behavior, which allows alterations in most genes to impact the "core" genes involved in a disease. A comprehensive statistical framework which includes the structure of biological data might address the aforementioned issues.

In this regard many authors turn to network biology to handle the complex interplay of biomolecules that lead to disease [7]. As its name suggests, network biology models biology as a network, where the biomolecules under study, often genes, are nodes, and selected functional relationships between them are the edges that link them. These functional relationships come from evidence that the genes jointly contribute to a biological function; for instance, their expression is correlated, or their gene products establish a protein-protein interaction. Under this view, complex diseases are not the consequence of a single altered gene, but of the interaction of multiple interdependent biomolecules [8]. In fact, an examination of biological networks shows that disease genes have differential properties [8] [9]. This is particularly true for cancer driver genes, which tend to be key players in connecting different, densely-connected communities of genes. Additionally, as genes that contribute to a disease tend to participate in similar biological functions, guilt-by-association strategies have proved effective at identifying disease genes [10].

Network-based, biomarker discovery methods exploit the guilt-by-association strategy to identify disease genes on GWAS data [11]. In essence, each SNP has a measure of association with the disease, given by the experiment, and functionally biological relationships, given by a network built on prior knowledge. Then, the problem becomes finding a functionally-related set of genes that is highly associated with the disease. Different solutions have been proposed to this problem, often stemming from divergent different mathematical frameworks and considerations of what the optimal solution looks like. Some methods strongly constrain the problem to certain kinds of subnetworks. Such is the extreme case of LEAN [12], which focuses on star subnetworks, i.e. instances were both a gene and its direct interactors are associated with the disease. Other algorithms, like dmGWAS [13] and heinz [14], focus on interconnected genes with high association with the disease. However, they differ in their tolerance to the inclusion of lowly associated nodes, and the possible number of disconnected subnetworks in the solution. Lastly, other methods also consider the topology of the network, favoring solutions that are densely interconnected; such is the case of HotNet2 [15], SConES [16], and SigMod [17].

In this work, we analyze the effectiveness of these six methods to discover new biomarkers on GWAS data. We focus on the GENESIS dataset [18], a study of familial breast cancer conducted in the French population. After following a classical GWAS approach, we use these network-based methods to recover additional familial breast cancer biomarkers. Some of them are known, while others are specific to this dataset. Lastly, we carry out a comparison of the solutions obtained by the different methods, and aggregate them to obtain a consensus network of predisposition to familial breast cancer.

## 2 Methods

### 2.1 GENESIS

The GENE Sisters (GENESIS) study was designed to investigate risk factors for familial breast cancer in the French population [18]. Index cases are patients with infiltrating mammary or ductal adenocarcinoma, who had a sister with breast cancer, and who have been tested negative for BRCA1 and BRCA2 pathogenic variants. Controls are unaffected colleagues and/or friends of the cases, born around the year of birth of the corresponding case ($\pm$ 3 years). We focused on the 2 577 samples of European ancestry, of which 1 279 are controls and 1 298 are cases. The genotyping was performed using the iCOGS array, a custom Illumina array designed to study genetic susceptibility of hormone-related cancers [19]. It contains 211 155 SNPs, including

SNPs putatively associated with breast, ovarian, and prostate cancers, SNPs associated with survival after diagnosis, and SNPs associated to other cancer-related traits, as well as functional candidate variants in selected genes and pathways.

## 2.2 Preprocessing and quality control

We discarded SNPs with a minor allele frequency lower than 0.1%, those not in Hardy - Weinberg equilibrium (P-value <0.001), and those with missing values on more than 10% of the samples. A subset of 20 duplicated SNPs in *FGFR2* were also removed. In addition, we removed the samples with more than 10% missing genotypes, and an additional 28 samples with TODO. The final dataset included 1 271 controls and 1 280 cases, genotyped over 197 083 SNPs.

We looked for population structure that could create confounding associations. A PCA revealed no differential population structure between cases and controls (Supplementary Figure 1). Independently, we did not find evidence of genomic inflation ($\lambda = 1.05$) either, further confirming the absence of confounding population structure.

## 2.3 High-score subnetwork search algorithms

### 2.3.1 SNP and gene association

To measure association between a genotype and the phenotype, we performed a per-SNP 1df $\chi^2$ allelic test using PLINK v1.90 [20]. Then, we used VEGAS2v2 to compute the gene-level association score [21] from the SNP P-values. In order to map SNPs to genes we used their overlap on the genome: all SNPs located within the boundaries of a gene, $\pm 50$ kb, were mapped to that gene. To compute the gene association we used the 10% of SNPs linked to the gene with lowest P-values. We used the 62 193 genes described in GENCODE 31 [22], although only 54 612 could be mapped to at least one SNP. Out of those, we focused exclusively on the 32 767 that could be mapped to an HGNC symbol. Out of the SNPs 197 083 remaining after quality control, 164 037 mapped to at least one of these genes.

We use such mapping to compare the outputs of methods who produce SNP- to those that produce gene-lists, and vice versa. In the former case, we consider any gene that can be mapped to any of the selected SNPs as selected as well. In the latter, we consider all the SNPs that can be mapped to that gene as selected by the method.

### 2.3.2 Mathematical notation

In this article, we use undirected, vertex-weighted networks, or graphs, $G = (V, E, w)$. $V = \{v_1, \ldots, v_n\}$ refers to the vertices, with weights $w : V \to \mathbb{R}$. Equivalently, $E \subseteq \{\{x, y\} | x, y \in V \wedge x \neq y\}$ refers to the edges. When referring to a subnetwork S, $V_S$ is the set of nodes in S and $E_S$ is the set of edges in S. A special case of subgraphs are *connected* subgraphs, which occur when every node in the subgraph can be reached from any other node.

On top of a weight, nodes have other properties provided by the topology of the graph. In this article we focus on two: degree centrality, and betweenness centrality. The degree centrality, or degree, is the number of edges that a node has. The betweenness centrality, or betweenness, is the number of times a node participates in the shortest paths between two other nodes.

In addition, we use several matrices that describe different properties of a graph. The described matrices are square, and have as many rows and columns as nodes are in the network. The element $(i, j)$ represents a selected relationship between $v_i$ and $v_j$. The *adjacency matrix* $W_G$ contains a 1 when the corresponding

nodes are connected through an edge, and 0 otherwise; the diagonal is zero. The *degree matrix* $D_G$ is a diagonal matrix which contains the degree of the different nodes. Lastly, the *Laplacian matrix* $L_G$ is defined as $L_G = D_G - W_G$.

### 2.3.3 Methods used

**Table 1:** Summary of the differences between the studied algorithms.

| Method Reference | Field | Nodes | Exhaustive | Solution | Components | Input | Scoring |
|---|---|---|---|---|---|---|---|
| dmGWAS [13] | GWAS | Gene | No | - | 1 | Summary | $-\log_{10}(P)$ |
| heinz [14] | Omics | Gene | Yes | - | 1 | Summary | BUM |
| HotNet2 [15] | Omics | Gene | Yes | Modular | $\geq 1$ | Summary | Local FDR |
| LEAN [12] | Omics | Gene | Yes | Star | $\geq 1$ | Summary | $-\log_{10}(P)$ |
| SConES [16] | GWAS | SNP | Yes | Modular | $\geq 1$ | Genotypes | $\chi^2$ |
| SigMod [17] | GWAS | Gene | Yes | Modular | 1 | Summary | $-\log_{10}(P)$ |

Beyond the assumption that genes that contribute to the same function will be nearby in the protein-protein interaction network (PPIN), they might be topologically related to each other in diverse ways (densely interconnected modules, nodes around a hub, a path, etc.). That is not the only choice to make: how to score the nodes, whether the affected mechanisms form a single connected component or several, how to frame the problem in a computationally efficient fashion, what is the best network to use, etc. In consequence, multiple solutions have been proposed. In this article, we examine six of them: five that explore the protein-protein interaction network, and one which explores SNP networks. We selected methods that were open source, had an implementation available, and an accessible documentation. Their main differences are summarized in Table 1.

**dmGWAS** dmGWAS searches the subgraph with the highest local density in low P-values [13]. To that end it searches candidate subnetwork solutions using a greedy, "seed and extend", heuristic:

1. Select a seed node.

2. Compute Stouffer's Z-score $Z_m$ for the current subgraph as

$$Z_m = \frac{\sum z_i}{\sqrt{k}}$$

where $k$ is the number of genes in the subgraph, $z_i = \phi^{-1}(1 - \text{P-value}_i)$, and $\phi^{-1}$ is the inverse normal distribution function.

3. Identify neighboring nodes i.e. nodes at distance $\leq d$. We set d = 2.

4. Add the neighboring nodes whose inclusion increases the $Z_{m+1}$ more than $Z_m \times (1 + r)$. In our experiments, we set r = 0.1.

5. Repeat 2-4 until no increment $Z_m \times (1 + r)$ is possible.

Lastly, the module's Z-score is normalized as

$$Z_N = \frac{Z_m - \text{mean}\,(Z_m(\pi))}{\text{SD}\,(Z_m(\pi))}$$

where $Z_m(\pi)$ represent a vector containing 100000 random subsets of the same number of genes.

We used the implementation of dmGWAS in the dmGWAS 3.0 R package [23]. We used the function *simpleChoose* to select the solution subnetwork, which aggregates the top 1% modules into the solution subnetwork.

**heinz** The goal of heinz is to identify the highest-scored connected subgraph on the network [14]. The authors propose a transformation of the genes' P-value into a score that is negative under no association with the phenotype, and positive value when there is. This transformation is achieved by modelling the distribution of P-values by a beta-uniform model (BUM) parameterized by the desired FDR. Thus formulated, the problem is NP-complete. To solve it efficiently it is re-casted as the Prize-Collecting Steiner Tree Problem (PCST), which seeks to select the connected subnetwork S that maximizes the *profit* p(S):

$$p(S) = \sum_{v \in V_S} p(v) - \sum_{e \in E_S} c(e).$$

were p(v) = w(v) - w' is the *profit* of adding a node, c(e) = w' is the *cost* of adding an edge, and $w' = min_{v \in V_G} w(v)$. All three are positive quantities. heinz implements the algorithm from [24], which in practice is often fast and optimal, neither is guaranteed. We used BioNet's implementation of heinz, available on Bioconductor [25, 26].

**HotNet2** HotNet2 was developed to find connected subgraphs of genes frequently mutated in cancer [15]. To that end, it considers both the local topology of the network and the scores of the nodes. The former is captured by an insulated heat diffusion process: at the beginning, the score of the node determines its initial heat; iteratively each node yields heat to its "colder" neighbors, and receives heat from its "hotter" neighbors, while retaining part of its own (hence, *insulated*). This process continues until equilibrium is reached, and results in a similarity matrix F. F is used to compute the similarity matrix E that accounts also for similarities in node scores as

$$E = F \, \text{diag}(w(V)),$$

where $\text{diag}(w(V))$ is a diagonal matrix with the node scores in its diagonal. We scored the nodes as in [27], assigning a score of 0 for the genes with low probability of being associated to the disease, and $-\log_{10}$(P-value) to those likely to be. In this dataset, the threshold separating both was a P-value of 0.125, which was obtained using a local FDR approach [28]. To obtain densely connected subnetworks,

HotNet2 prunes E, only preserving edges such that $w(E) > \delta$. Lastly, HotNet2 evaluates the statistical significance of the subnetworks by comparing their size to the size of networks obtained by permuting the node scores. HotNet2 has two parameters: the restart probability $\beta$, and the threshold heat $\delta$. Both parameters are set automatically by the algorithm, and are robust [15]. HotNet2 is implemented in Python [29].

**LEAN** LEAN searches disregulated "star" gene subnetworks, that is, subnetworks composed by one central node and all its interactors [12]. By imposing this restriction, LEAN is able to exhaustively test all such subnetworks (one per node). For a particular subnetwork of size $m$, the P-values corresponding to the involved nodes are ranked as $p_1 \leq \ldots \leq p_m$. Then, $k$ binomial tests are conducted, to compute the probability of having $k$ out of $m$ P-values lower or equal to $p_k$ under the null hypothesis. The minimum of these $k$ P-values is the score of the subnetwork. This score is transformed into a P-value through an empirical distribution obtained via a subsampling scheme, where sets of $m$ genes are selected randomly, and their score computed. Lastly, P-values are corrected for multiple testing through a Benjamini-Hochberg correction. We used the implementation of LEAN from the LEANR R package [30].

**SConES** SConES searches the minimal, modular, and maximally associated subnetwork in a SNP graph [16]. Specifically, it solves the problem

$$\arg\max_{S \subseteq G} \underbrace{\sum_{v \in V_S} w(v)}_{\text{association}} + \lambda \underbrace{\sum_{v \in V_S} \sum_{u \notin V_S} L_{vu}}_{\text{connectivity}} - \underbrace{\eta |V_S|}_{\text{sparsity}} \tag{1}$$

where $\lambda$ and $\eta$ are hyperparameters that control the sparsity and the connectivity of the model. Given two hyperparameters, the aforementioned problem has a unique solution, that SConES finds using a graph min-cut procedure. We used the version on SConES implemented in the R package martini [31]. As in [16], we selected $\lambda$ and $\eta$ by cross-validation, choosing the values that produce the most stable solution across folds. Note that the solution to the above problem can consist of several connected subnetworks which are disconnected from each other. In this case, the selected hyperparameters were $\eta = 3.51$, $\lambda = 210.29$ for SConES GS; $\eta = 3.51$, $\lambda = 97.61$ for SConES GM; and $\eta = 3.51$, $\lambda = 45.31$ for SConES GI.

**SigMod** SigMod aims at identifying the most densely connected gene subnetwork that is most strongly associated to the phenotype [17]. It addresses an optimization problem similar to that of SConES (Equation 1), but using the Laplacian matrix rather than the adjacency matrix (Section 2.3.2), to quantify solutions containing many edges.

$$\arg\max_{S \in G} \underbrace{\sum_{v \in V_S} w(v)}_{\text{association}} + \lambda \underbrace{\sum_{v \in V_S} \sum_{u \in V_S} W_{vu}}_{\text{connectivity}} - \underbrace{\eta |V_S|}_{\text{sparsity}} \; .$$

As SConES, this optimization problem can also be solved by a graph min-cut approach.

SigMod presents three important differences with SConES. First it is designed for gene-gene networks. Second, by replacing the Laplacian by the adjacency matrix, it favors subnetworks containing many edges. SConES, instead, penalizes connections between the selected selected and unselected nodes. Third, it returns a single connected subnetwork, which it achieves by exploring a grid of hyperparameters

and processing their respective solutions. Specifically, for the range of $\lambda = \lambda_{\min}, \ldots, \lambda_{\max}$ for the same $\eta$, it prioritizes the solution with the largest change in size from $\lambda_n$ to $\lambda_{n+1}$. Such a large change implies that the network is strongly interconnected. This results in one candidate solution for each $\eta$, which are processed by removing any node not connected to any other. A score is assigned to each candidate solution by summing their node scores and normalizing by size. The candidate solution with the highest standardized score is the chosen solution. SigMod is implemented in an R package [32].

### 2.3.4 Gene-gene network

Out of the six methods tested, five use a gene-gene interaction network (Section 2.3.3). Although their respective statistical frameworks are compatible with any type of network (protein interactions, gene co-expression, regulatory, etc.), for practical reasons we focused on a PPIN, as they are interpretable, well characterized, and most of the methods were designed to scale appropriately to it. We built our PPIN from both binary and co-complex interactions stored in the HINT database (release April 2019) [33]. Unless specified otherwise, we used only interactions coming from high-throughput experiments to avoid biasing the topology of the network by well-studied genes with more known interactions on average. Out of the 146 722 interactions from high-throughput experiments that HINT stores, we were able to map 142 541 to a pair of HGNC symbols. The scoring function for the nodes changed from method to method (Section 2.3.3).

Additionally, we compared the results of the aforementioned PPIN with those obtained on another PPIN built using interactions coming from both high-throughput and targeted studies. In that case, out of the 179 332 interactions in HINT, we mapped 173 797 to a pair of HGNC symbols.

### 2.3.5 SNP networks

SConES [16] is the only of the studied methods designed to handle SNP networks. As in gene networks, two SNPs are linked in a SNP network when there is evidence of shared functionality between two SNPs. The authors suggested three ways of building these networks: connecting the SNPs consecutive in the genomic sequence ("GS network"); interconnecting all the SNPs mapped to the same gene, on top of GS ("GM network"); and interconnecting all SNPs mapped to two genes for which a protein-protein interaction exists ("GI network"). We focused on the GI network, as it is the network that fits better the scope of this article. However, at different stages of the manuscript we also used GS and GM. For the GM network, we used the mapping described in Section 2.3.1. For the GI network, we used the PPI as described in Section 2.3.4. For all three networks the node score used is the association of the individual SNPs with the phenotype; specifically, we used the 1 d.f. $\chi^2$.

### 2.3.6 Consensus network

The different high-weight subnetwork discovery algorithms make different assumptions on the properties of the solutions, and employ different strategies to find them. Hence, combining the outcome of the different approaches might provide a more complete outlook on the specific alterations on the GENESIS dataset. We built such consensus network by retaining the nodes that were selected by at least two of the methods. We combined the results of 6 methods: dmGWAS, heinz, HotNet2, LEAN, SConES GI, and SigMod.

## 2.4 Evaluation of methods

### 2.4.1 Classification accuracy of selected biomarkers

A desirable solution is one that is sparse, while offering a good predictive power on unseen samples. We evaluated the predicting power of the SNPs selected by the different methods through the performance of an L1-penalized logistic regression trained exclusively on those SNPs to predict the outcome (case/control). The L1 penalty helps to account for LD to reduce the number of SNPs included in the model (size of the active set), while improving the generalization of the classifier. The value of the regularization parameter, which controls both the magnitude and the sparsity of the coefficients, was set by cross-validation. To that end, we used the different network-methods on a random subset of 80% of the samples. On this same subset we trained our classifier exclusively on the SNPs selected by a particular method. When the method retrieved a list of genes (all of them except SConES), all the SNPs mapped to any of those genes were used. Then we evaluated performance of the classifier on the remaining 20% of the dataset. We repeated this procedure 5 times to estimate the average and the deviation of the different performance measures. The different performance measures we used were: size of the solution, size of the active set, specificity, and sensitivity. The size of the active set provides an estimate of a plausible, more sparse solution with a comparable predictive power to the original solution.

Additionally, for each of the methods, we evaluated their stability and their runtime. The stability of an algorithm is its sensitivity to small changes of the input, and is measured using the Pearson's correlation between different runs as suggested in [34]. To obtain a baseline, we also performed the procedure using all the SNPs. Lastly, another desirable property is that the method retrieves a good candidate causal subnetwork. In consequence, we compared the outcome of each of the methods to the consensus subnetwork of all the solutions (Section 2.3.6).

### 2.4.2 Biological relevance of the genes

An alternative way to validate the results is comparing our results to an external dataset. For that purpose, we recovered a list of 153 genes associated to familial breast cancer from DisGeNET [35]. Across this article we refer to these genes as *familial breast cancer genes*.

Additionally, we used the summary statistics from the Breast Cancer Association Consortium (BCAC) [36]. BCAC is one of the largest efforts in GWAS, genotyping over $120\,000$ women of European ancestry. As opposed to GENESIS, samples were not selected based on family history, and hence is enriched in sporadic breast cancers. Another difference is that BCAC is a relatively heterogeneous study on a pan-European sample, while GENESIS is a homogeneous dataset focused on the French population. Despite these differences, there should be shared genetic architecture. On top of that, that overlap should become larger when the results are aggregated at the gene level. For that purpose, we computed the gene association as in Section 2.3.1. iCOGS array was used for genotyping in BCAC [19], the same array as for GENESIS [18]. Although imputed data is available, we used exclusively the SNPs available on GENESIS after quality control to make the results comparable.

## 2.5 Code availability

This work required developing computational pipelines for several GWAS analyses, such physically mapping SNPs to genes, computing gene scores, and performing six different network analyses. For each of those processes, a streamlined, project-agnostic pipeline with a clear interface was created. They are compiled in the following GitHub repository: `https://github.com/hclimente/gwas-tools`. The code that applies

these pipelines to the GENESIS project, as well as the code that reproduces all the analyses in this article are available at `https://github.com/hclimente/genewa`. Although the GENESIS dataset is not publicly available, the published code should work on any other GWAS dataset. All the produced gene subnetworks were deposited on NDEx (`http://www.ndexbio.org`), under the UUID e9b0e22a-e9b0-11e9-bb65-0ac135e8bacf.

# 3  Results

## 3.1  FGFR2 is strongly associated with familial breast cancer

We conducted association analyses in the GENESIS dataset at both the SNP and the gene levels (Section 2.3.1). Two genomic regions have a P-value lower than the Bonferroni threshold in chromosomes 10 and 16 (Figure 1A). The former overlaps with gene FGFR2; the latter with CASC16, and it is located near the protein-coding gene TOX3. Variants in both FGFR2 and TOX3 were related to breast cancer susceptibility in other cohorts negative for BRCA1/2 [37]. Only the peak in chromosome 10 replicated in the gene-level analysis, with FGFR2 just above the threshold of significance (Figure 1B).

These results show the overlap in the genetic architecture of the disease between the studied French population sample and other populations, especially at the gene level. In addition, there are other SNPs whose P-values, although higher than the conventional threshold of significance, show a strong association with familial breast cancer. The most prominent of such regions are 3p24 and 8q24, both of which have been associated to breast cancer susceptibility in the past [38, 39]. This motivates exploring network methods, which trade statistical significance for biological relevance.
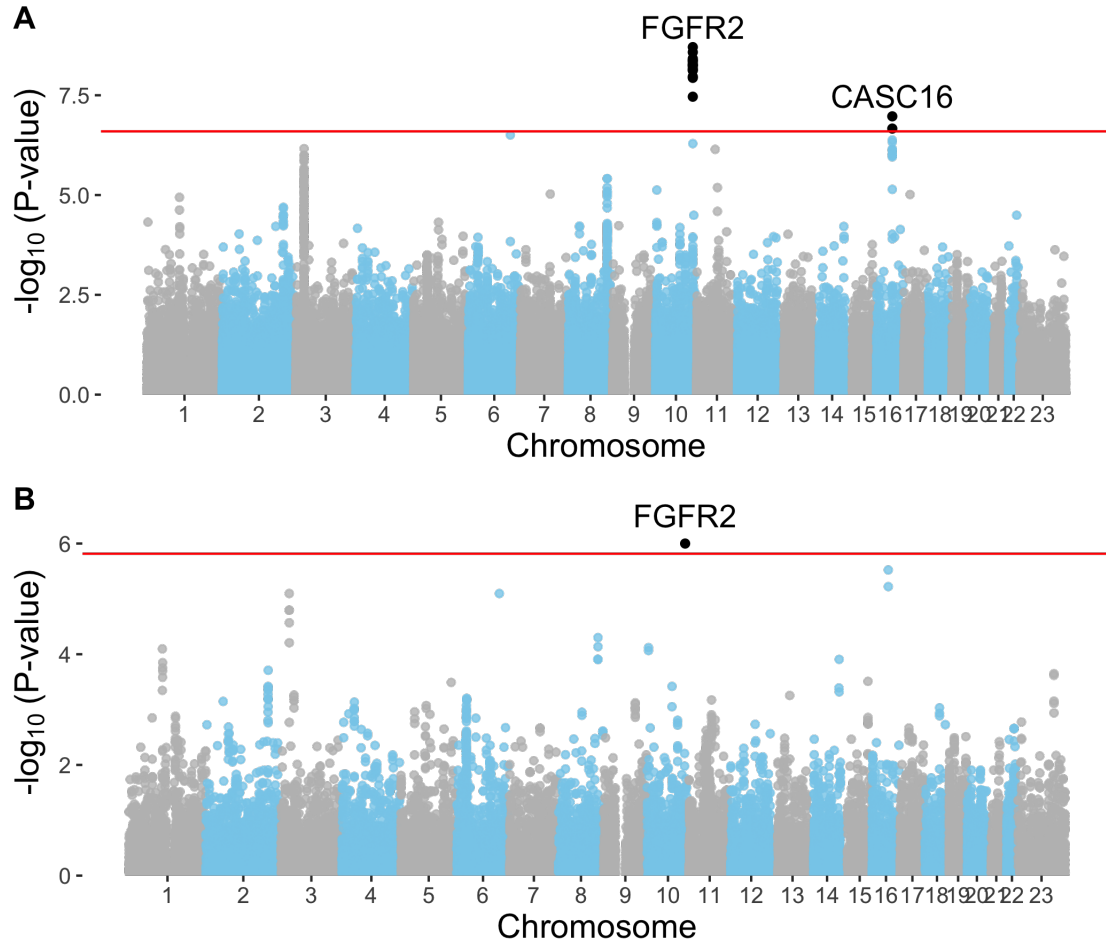
## 3.2  Network methods successfully identify genes associated with breast cancer

**Table 2:**  Summary statistics on the results of multiple network methods on the gene-gene interaction network. The first row contains the summary statistics on the whole network.

| Network | # genes | # edges | $\overline{\text{Betweenness}}$ | $\hat{P}_{gene}$ | $\rho_{consensus}$ |
|---|---|---|---|---|---|
| HINT HT | 13 619 | 142 541 | 16 706 | 0.46 | 0.066 |
| Consensus | 55 | 117 | 74 062 | 0.0051 | 1 |
| dmGWAS | 194 | 450 | 49 115 | 0.19 | 0.41 |
| heinz | 4 | 3 | 113 633 | 0.0012 | 0.21 |
| HotNet2 | 440 | 374 | 7 739 | 0.048 | 0.31 |
| LEAN | 0 | 0 | - | - | - |
| SConES GI | 0 (1) | 0 | - | - | - |
| SigMod | 142 | 249 | 92 603 | 0.0083 | 0.73 |

*# genes*: number of genes selected out of those that are part of the PPIN; in parentheses, the total number of genes. $\overline{\text{Betweenness}}$: mean betweenness of the selected genes in the PPIN. $\hat{P}_{gene}$: median P-value of the selected genes. $\rho_{consensus}$: Pearson's correlation between the subnetwork and the consensus network.

We applied six network methods to the GENESIS dataset (Section 2.3.3), obtaining six solutions (Supplementary Figure 2, Supplementary Files 1 and 2): one for each of the five gene-based methods (Section 2.3.4),

**Fig. 1:** Association in GENESIS. The red line represents the Bonferroni threshold. **(A)** SNP association, measured from the outcome of a 1df $\chi^2$ allelic test. Significant SNPs that are within a coding gene, or within 50 kilobases of its boundaries, are annotated. The Bonferroni threshold is $2.54 \times 10^{-7}$. **(B)** Gene association, measured by P-value of VEGAS2v2 [21] using the 10% of SNPs with the lowest P-values. The Bonferroni threshold is $1.53 \times 10^{-6}$.
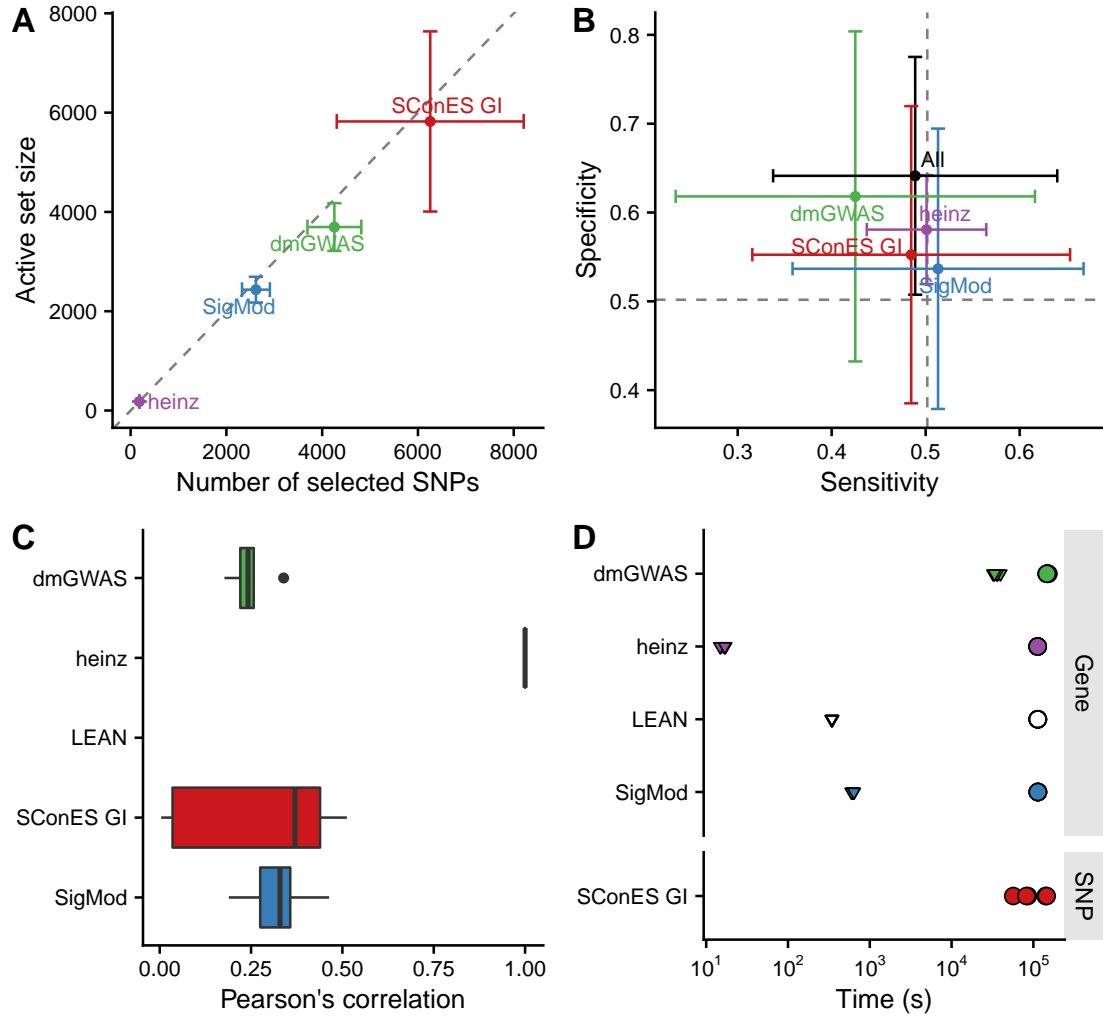
and one for SConES GI (Section 2.3.5). The solutions are very heterogeneous (Table 2 and Supplementary table 1): none of the subnetworks examined by LEAN is significant (adjusted P-value < 0.05), while HotNet2 produced the largest solution subnetwork with 440 genes. SConES GI failed to recover genes in the PPIN, but it recovered one genomic region mapped to RNA gene RNU6-420P. All solution subnetworks except LEAN's are, on average, more strongly associated to familial breast cancer than the whole PPIN (median P-values $\ll 0.46$), despite containing genes with higher P-values (Supplementary Figure 3). This exemplifies the trade-off between statistical significance and biological relevance. However, there are nuances between solutions: heinz strongly favored highly associated genes, while dmGWAS is less conservative (median gene P-values 0.0012 and 0.19, respectively); SConES tended to select whole LD-blocks; and HotNet2 and SigMod were less likely to select lowly associated genes.

The solution subnetworks present other desirable properties. First, four of the methods succeeded at recovering genes involved in the disease (Supplementary Figure 5), as their subnetworks were enriched in familial breast cancer genes (dmGWAS, heinz, HotNet2, and SigMod, Fisher's exact test one-sided P-value < 0.03). We also compared the outcome of the network methods to the association tests conducted on the population of European ancestry from the Breast Cancer Association Consortium (BCAC) [36] (Supplementary Figure 4). Encouragingly, every solution subnetwork is enriched in genes or SNPs that are Bonferroni-significant in BCAC. This confirms the capability of network methods to find the same signal as in more powered studies by leveraging on prior knowledge. Second, the genes in four solution subnetworks display on average a higher betweenness centrality than the rest of the genes, a difference that is significant in three solutions (dmGWAS, and SigMod, Wilcoxon rank-sum test P-value < $1.4 \times 10^{-21}$). This agrees with the notion that disease genes are more central than other, non-essential genes [9]. We observe that this conclusion holds in this disease, as familial breast cancer genes have higher betweenness centrality than others (one-tailed Wilcoxon rank-sum test P-value = $2.64 \times 10^{-5}$, Supplementary Figure 8C). Interestingly, SConES' selected SNPs are also more central than the average SNP (Supplementary table 1), suggesting that causal SNPs are also more central than unrelated SNPs. However, very central nodes are also more likely to be connecting a random pair of nodes, making then more likely to be selected by the examined methods. Hence, further work is needed draw conclusions.

As the solutions were quite different from each other it is hard to draw joint conclusions. The 4-gene solution selected by heinz includes the familial breast cancer gene TOX3, in region 16q12. By dealing with SNP networks, SConES studies the association of non-coding regions, as well as SNPs in any gene, coding or else. In fact, SConES GI, which adds to GM the interactions between genes, retrieves 4 subnetworks in intergenic regions, and 1 overlapping an RNA gene (RNU6-420P). SigMod, despite being related to SConES, produces a vastly different, large solution. On top of recovering three familial breast cancer genes, a keratin-based region of its subnetwork affects the cytoskeleton (*structural constituent of cytoskeleton*, GO enrichment's adjusted P-value = $9.10 \times 10^{-4}$), a potentially novel susceptibility mechanisms to cancer. Interestingly, dmGWAS solution is also related to cytoskeleton (*tubulin binding*, GO enrichment's adjusted P-value = 0.031). But, additionally, it includes a submodule of proteins related to *unfolded protein binding* (GO enrichment's adjusted P-value = 0.045), which was related to cancer susceptibility [40]. Lastly, HotNet2 produced 135 subnetworks, 115 of which have less than five genes. The second largest subnetwork (13 nodes), contains two familial breast cancer genes: CASP8 and BLM.

## 3.3  heinz retrieves a small, highly informative set of biomarkers in a fast and stable fashion

As the methods produced such different results, we compared their solutions in a 5-fold subsampling setting (Section 2.4.1). Specifically, we measured the following properties (Figure 2): (i) size of the solution

**Fig. 2:** Comparison of network-based GWAS methods on GENESIS. Each method was run 5 times of a random subset of the samples, and tested on the remaining samples (Section 2.4.1). **(A)** Number of SNPs selected by each method and number of SNPs on the active set used by the Lasso classifier. Points are the average over the 5 runs; lines represent the standard error of the mean. A grey diagonal line with slope 1 is added for comparison. For reference, the active set of Lasso using all the SNPs included, on average, 154 117.4 SNPs. **(B)** Sensitivity and specificity on test set of the L1-penalized logistic regression trained on the features selected by each of the methods. In addition, the performance of the classifier trained on all SNPs is displayed. Points are the average over the 5 runs; lines represent the standard error of the mean. **(C)** Pairwise Pearson's correlations of the solutions used by different methods. A Pearson's correlation of 1 means the two solutions are the same. A Pearson's correlation of 0 means that there is no SNP in common between the two solutions. **(D)** Runtime of the evaluated methods, by type of network used (gene or SNP). For gene network-based methods, inverted triangles represent the runtime of the algorithm itself, and circles the total time, which includes the algorithm themselves and the additional 119 980 seconds (1 day and 9.33 hours) which took VEGAS2v2 on average to compute the gene scores from SNP summary statistics.

12

subnetwork; (ii) sensitivity and specificity of an L1-penalized logistic regression on the selected SNPs; (iii) stability; and (iv) computational runtime.

Both solution size and active set of SNPs selected by Lasso varies greatly between the different methods (Figure 2A). heinz has the smallest solutions, with an average of 182 selected SNPs, out of which 5.6% (10.2) are selected by Lasso. The largest solutions come from SConES GI (6256.6 SNPs), and dmGWAS (4255.0 SNPs). Interestingly, heinz has the highest proportion of the selected SNPs that go into the active set (99.9%), although it is high for all the methods ($> 86\%$). This suggests methods are selecting informative SNPs on average.

The sensitivity and specificity of the classifier on the testing data informs us about the usefulness of the selected SNPs as patient classification (Figure 2B). All classifiers' sensitivities were in the $0.42 - 0.51$ range; the specificities, between 0.54 and 0.62. On average, SigMod had the highest sensitivity (0.51); dmGWAS, the highest specificity (0.52). Both heinz and SigMod had on average better sensitivity than the classifier trained on all the SNPs, but none had superior specificity. However, the differences are negligible, well within the 95% confidence interval.

Another desirable quality of an algorithm is stability (Section 2.4.1). Both heinz and LEAN displayed a high stability in our benchmark, consistently selecting the same genes and no genes over the 5 subsamples, respectively (Figure 2C). The other methods displayed similarly low stabilities.
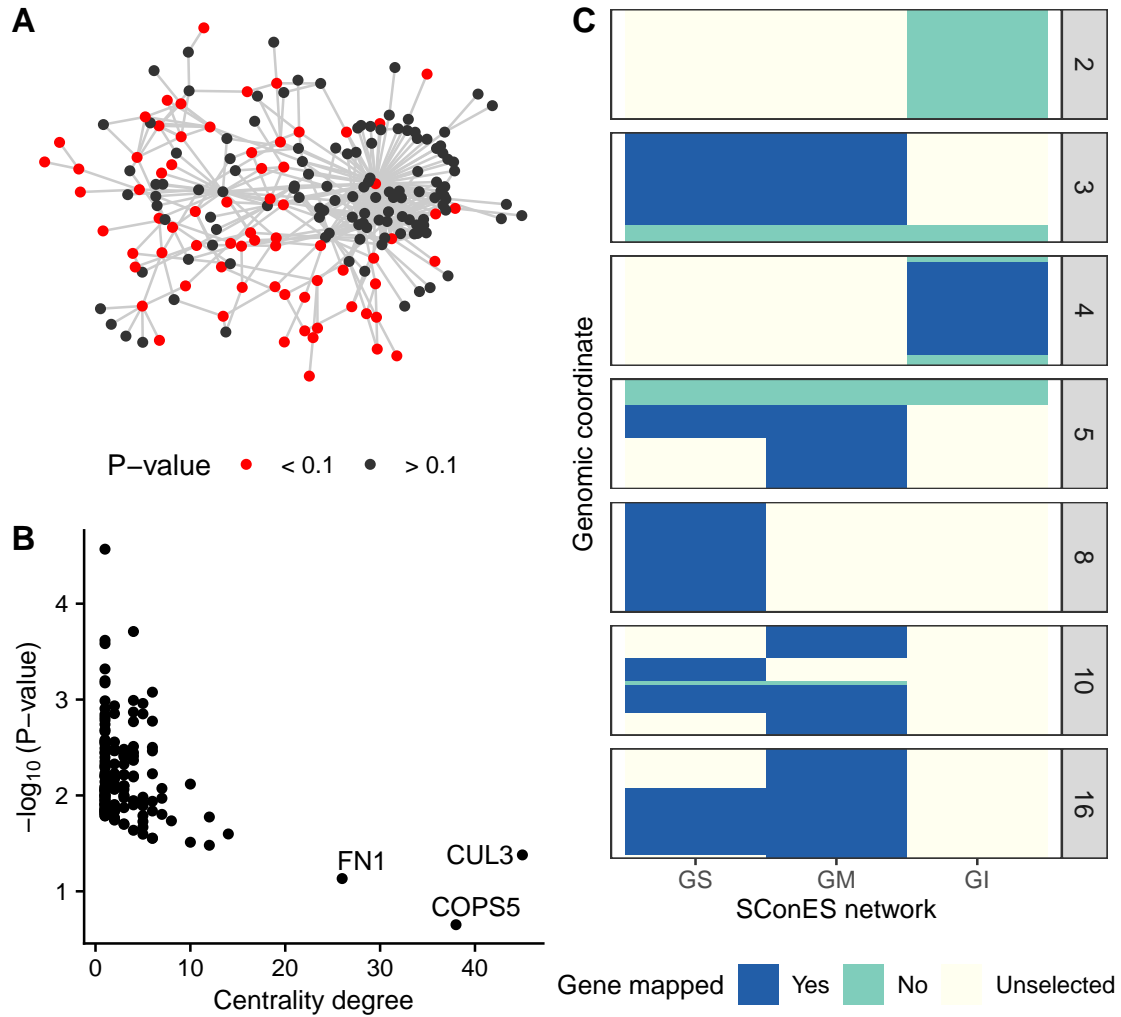
In terms of computational runtime, the fastest method was heinz (Figure 2D), which leverages on its ability to find efficiently the solution in a few seconds. The slowest method was dmGWAS (1 day and 17 hours on average) followed by SConES GI (1 day and 4.32 hours on average). Including the time required to compute the gene scores, however, slows down considerably gene-based methods; on this benchmark, that step took on average 1 day and 9.33 hours. Considering that, it took 3 days and 2.4 hours on average for dmGWAS' to produce results.

## 3.4 No solution is perfect

In practice, and despite their similarities and their involvement in cancer mechanisms, the solutions are remarkably different (Supplementary Figure 6A). That is due to the particulars of the methods, and directly or indirectly, they provide information about the dataset. For instance, the fact that LEAN did not provide any biomarkers implies that there is no gene such that both itself and its environment are on average strongly associated with the disease.

In this dataset, heinz's solution is very conservative, providing a small solution with the lowest median P-value for the subnetwork (Table 2). Due to this parsimonious and highly associated solution, it was the best method to select a set of good biomarkers for classification. (Figure 2B). Its conservativeness stems from its preprocessing step, which models the gene p-values as a mixture model of a beta and a uniform distribution, controlled by an FDR parameter. Due to the limited signal at the gene level in this dataset (Figure 1B), only 36 of them are retain a positive score after applying the BUM model (Section 2.3.3). Hence, heinz's solution subnetwork consists only of 4 genes, which does not provide much insight of the biology of cancer. Importantly, it ignores genes that are strongly associated to cancer in this dataset like FGFR2.

On the other end of the spectrum, we have large solutions provided by dmGWAS, HotNet2, and SigMod. dmGWAS' subnetwork is the least associated subnetwork on average. This is due to the greedy framework it uses, which considers all nodes at distance 2 of the examined, and accepts weakly associated genes if they are linked to another, strongly associated one. This is exacerbated when the results of successive greedy searches are aggregated, leading to a large, tightly connected cluster of unassociated genes (Figure 3A). SigMod displays the same tendency, as the most central genes are the least associated to the disease (Figure

**Fig. 3:** Drawbacks confronted when using network guided methods. **(A)** dmGWAS solution subnetwork. Genes with a P-value $< 0.1$ are highlighted in red. **(B)** Centrality degree and -$\log_{10}$ of the VEGAS P-value for the nodes in SigMod solution subnetwork. **(C)** Genomic regions where either SConES GS, GM or GI select SNPs.
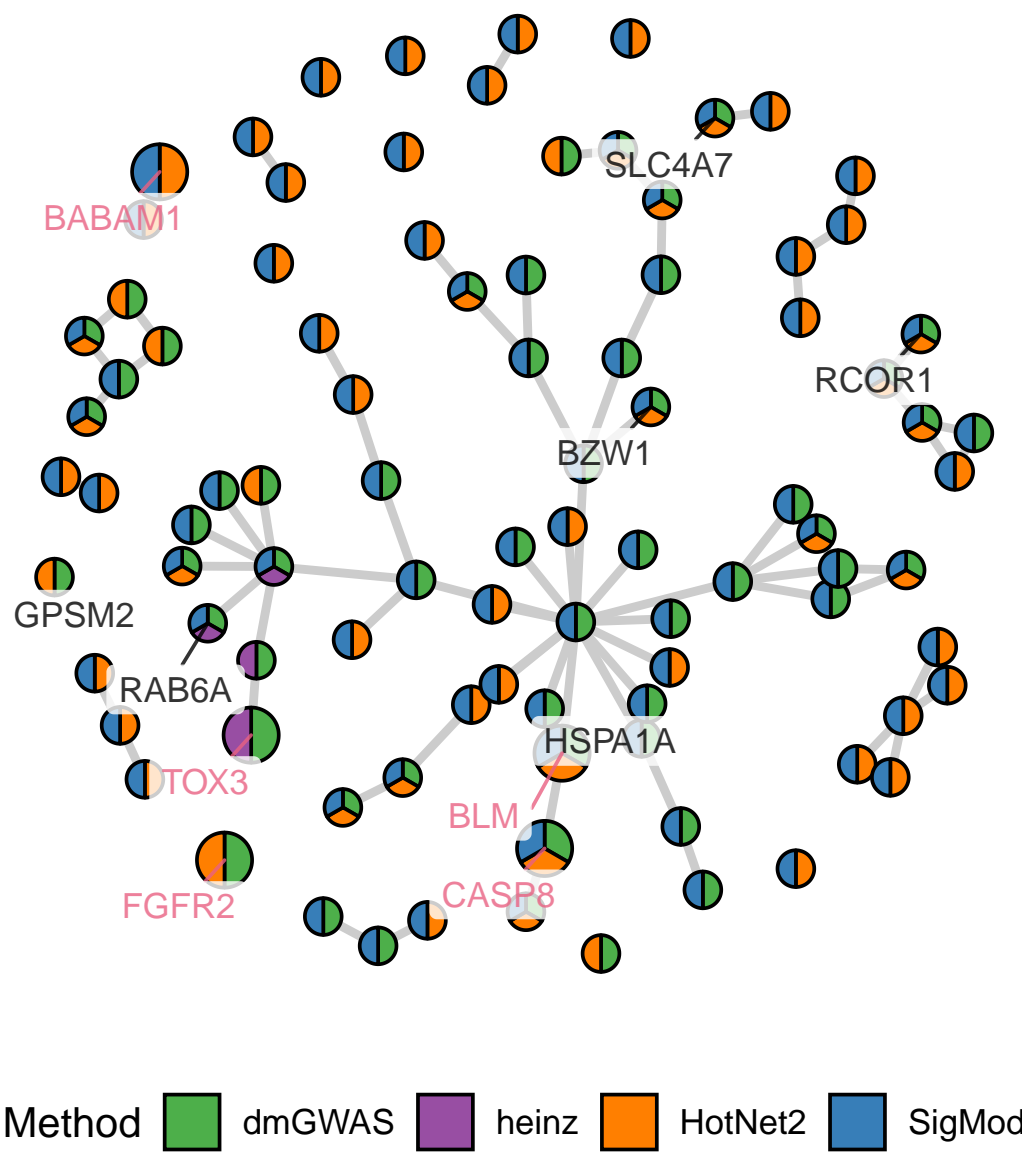
3B). This relatively low signal-to-noise ratio combined with the large solution requires additional analyses to draw conclusions, such as enrichment analyses. In the same line, HotNet2's subnetwork is even harder to interpret, being composed of 440 genes divided into 135 subnetworks. Lastly, SigMod misses some of the most strongly associated, familial breast cancer genes in the dataset, like FGFR2 and TOX3.

By virtue of using a SNP subnetwork, SConES analyzes each SNP in their context. It therefore selects SNPs in genes none of whose interactors are associated to the disease, as well as SNPs in non-coding regions or in non-interacting genes. In fact, due to linkage disequilibrium, such genes are favored by SConES, as selecting an SNPs in an LD block which overlaps with a gene favors selecting the rest of the gene. This might explain why the GS and GM networks, heavily affected by linkage disequilibrium, produce similar results (Supplementary Figure 6B). On the other hand, SConES penalizes selecting SNPs and not their neighbors. This makes it conservative regarding SNPs with many interactions, for instance those mapped to hubs in the PPIN. For this reason, SConES GI did not select any protein coding gene, despite selecting similar regions as SConES GS (Figure 3C). In fact SConES GS and SConES GM select regions related to breast cancer, like 16q12 (TOX3, Section 3.1), 3p24 (SLC4A7/NEK10 [41]), 5p12 (FGF10, MRPS30 [42]), and 10q26 (FGFR2, Section 3.1). On top of that only SConES GS selects region 8q24 (POU5F1B [43]). We hypothesize that the lack of results on the PPIN network of SConES GI and LEAN due to the same cause: the absence of joint association of a module. Although in the case of SConES other hyperparameters could lead to a more informative solution (lower $\lambda$), finding the appropriate configuration is hard. In addition, the iCOGS platform is not a real GWAS experiment: the genome is not unbiasedly surveyed, some regions are fine-mapped — which might distort gene structure in GM and GI networks — while others are under studied - hurting the accuracy with which the GS network captures the genome structure.

## 3.5  Aggregating solutions provides insights into the biology of cancer

To leverage on the strengths of each of the methods and compensate their respective weaknesses, we built a consensus subnetwork that captures the mechanisms most shared among the solution subnetworks (Section 2.3.6). The consensus subnetwork (Figure 4) contains 93 genes and is enriched in familial breast cancer genes (Fisher's exact test P-value = $7.8 \times 10^{-5}$). Due to the limited overlap between methods, only 20 genes were common to more than two of them (Supplementary Figure 8A). Encouragingly, the more methods selected a gene, the higher its association was (Supplementary Figure 8B). Globally, a GO enrichment shows the involvement of two cellular processes: unfolded protein binding, and structural constituent of cytoskeleton (adjusted P-values of 0.001, 0.001, respectively), which were already observed in different solutions (Section 3.2). Remarkably, many of the selected genes are related to mitochondrial translation. For instance, MRPS30 (VEGAS P-value = 0.001), encodes a mitochondrial ribosomal protein and was also linked to breast cancer susceptibility [42]. Albeit disconnected from MRPS30, the consensus network includes a 2-node subnetwork composed of two mitochondrial ribosomal protein (MRPS31 - VEGAS P-value = $7.67 \times 10^{-3}$ - and MRPS18B - VEGAS P-value = $7.92 \times 10^{-3}$), which suggests an involvement of mitochondrial ribosomes in carcinogenesis [?].

We also examined the topological properties of the nodes. The genes in the consensus network have higher betweenness centrality than the rest of the genes (Wilcoxon rank-sum test P-value = $4.29 \times 10^{-18}$). Interestingly, within genes in the consensus network, cancer genes are as central as non-cancer genes (Wilcoxon rank-sum test P-value = 0.57). Centrality, however, is weakly anti-correlated with association to the disease (Pearson correlation coefficient = -0.26, Supplementary Figure 8D), which suggests that some highly central genes were selected because they were on the shortest path between two highly associated genes. In view of this, we hypothesize that highly central genes might contribute to the heritability through consistent alterations of their neighborhood, consistent with the omnigenic model of disease [6]. For instance, the

**Fig. 4:** Consensus subnetwork on GENESIS (Section 2.3.6). Each node is represented by a pie chart, which accounts the methods that selected it. The labeled genes have a VEGAS2v2 P-value < 0.001 and/or are known familial breast cancer genes (colored in pink).

most central node in the consensus network is COPS5 (Supplementary Figure 7), a gene related to multiple hallmarks of cancer and which is overexpressed in multiple tumors, including breast and ovarian cancer [44]. Despite its lack of association in GENESIS (VEGAS P-value = 0.22), its neighbors in the consensus subnetwork have consistently low P-values (median VEGAS P-value = 0.006).

The consensus subnetwork is not completely connected: out of the 93 genes, the largest connected subnetwork includes only 49. A GO enrichment analysis showed that this component is related to three major cellular processes: unfolded protein binding, structural constituent of cytoskeleton, and poly(U) RNA binding (adjusted P-values of 0.01, 0.04, and 0.04, respectively). We found support in the literature of the involvement of each of these functions in the development of cancer, as discussed next. The consensus network also contains a protein directly involved in caspase-mediated apoptosis, CASP8 (VEGAS P-value = $1.95 \times 10^{-4}$). This is related to the enriched activity, *unfolded protein binding*, which inhibits caspase-dependent apoptosis, raising the chances of developing cancer [40]. It involves three Hsp70 chaperones of the consensus subnetwork: HSPA1A, HSPA1B, and HSPA1L. They all are closely encoded in 6p21. In fact out of the 22 SNPs that map to any of these three genes, 9 map to all of them, and 4 to two, making hard to disentangle their association. HSPA1A was the most strongly associated one (VEGAS P-value = $8.37 \times 10^{-4}$). Remarkably, 14 of the 93 genes are in subnetworks of size 1 (isolated) or two, as they do not have a consistently altered neighborhood. One of them is the familial breast cancer FGFR2 (Section 3.1). Another one is SLC4A7 (VEGAS P-value = $2.70 \times 10^{5}$), a gene encoding a sodium bicarbonate cotransporter. In the past the genomic region containing both SLC4A7 and nearby gene NEK10 (VEGAS P-value = $1.56 \times 10^{-5}$) were associated with familial breast cancer [41]. NEK10 is a gene that might be involved in cell-cycle control, but it is absent from the PPIN and hence it could not be studied by gene methods. Despite that, the fact that both dmGWAS, HotNet2 and SigMod link SLC4A7 in their different subnetwork supports the notion that this gene is the responsible for cancer susceptibility.

## 3.6 Limitations of network analyses

The strength of network-based analyses comes from leveraging prior knowledge to boost discovery. In consequence, they falter in front of understudied genes, especially those not in the network. Out of the 32767 genes that we can map the genotyped SNPs, 60.7% (19887) are not in the protein-protein interaction network. The majority of those (14660) are non-coding genes, mainly lncRNA, miRNA, and snRNA (Supplementary Figure 9). The importance of these genes, like CASC16, is highlighted in Section 3.1. Among the excluded protein-coding genes we find genes like NEK10 (P-value $1.6 \times 10^{-5}$) or POU5F1B, both linked to breast cancer susceptibility [41]. However, on average protein coding genes absent from the PPIN are less associated with this phenotype (Wilcoxon rank-sum P-value = $2.79 \times 10^{-8}$, median P-values of 0.43 and 0.47). As we are using interactions from high-throughput experiments, such difference cannot be due to well-known genes having more known interactions. As disease genes tend to be more central [9], we hypothesize that it is due to interactions between central genes being more likely. It is worth noting that network approaches that do not use PPIs, like SConES GS and GM, did recover SNPs in NEK10 and CASC16. Lastly, all the methods rely heavily on how SNPs are mapped to genes. In Section 3.1 we highlight ambiguities that appear when genes overlap or are in linkage disequilibrium.

As not all databases compile the same interactions, the choice of the PPIN determines the final output. In this work we used exclusively interactions from HINT from high-throughput experiments. This responds to concerns of some authors about biases introduced by adding interactions coming from targeted studies in the literature [45, 33]. It is a "rich getting richer" phenomenon, where popular genes have a higher proportion of their interactions described. On the other hand, one study found that the best predictor of the performance of a network for disease gene discovery is the size of the network [10]. This would support using

the largest amount of interactions. To clarify their impact on this study, we compared the impact of using only physical interactions from high-throughput experiment versus interactions from both high-throughput and the literature (Section 2.3.4. For most of the methods using a larger network did not greatly impact the size or the stability of the solution, the classification accuracy, or the runtime (Supplementary Figure 10).

# 4    Discussion

In this article we evaluate the viability of systems biology approaches to GWAS, and examine a GWAS dataset on familial breast cancer focused on BRCA1/2 negative French women. Systems biology addresses two of the largest GWAS issues: interpretability and an overly conservative statistical framework that hinders discovery. This is achieved by considering the biological context of each of the genes and SNPs. Based on divergent considerations of what the desired set of biomarkers is, several methods for network-guided biomarker discovery have been proposed. We reviewed the performance of six of them on GWAS. Despite their differences, most of them produced a relevant subset of biomarkers, recovering known familial breast cancer genes. We also discuss the limitations of such analyses, related to the lack of known interactions around some genes. A crucial step for the gene based methods is the computation of the gene score. In the manuscript we used VEGAS2v2 [21] due to the flexibility it offers to use user-specified gene annotations. However, it presents known problems (selection of an appropriate percentage of top SNPs, long runtimes and P-value precision limited to the number of permutations [27]), other algorithms might have more statistical power.

The network methods we studied differ in what the optimal solution subnetwork looks like. On the one hand, SConES and heinz prefer small highly associated solutions. On the other hand, SigMod and dmGWAS gravitate towards larger, less associated solutions which provide a wide overview of the biological context. While the former provide a reduced set of biomarkers, the latter deepen our understanding of the disease and provide biological hypotheses. They are not exempt of limitations. dmGWAS and SigMod's solution's size require further analyses, which risk oversimplifying their richness. Also, incautious practitioners might be misled by some genes, which are very central in the solution subnetworks, while being weakly associated. Nonetheless, they are pushed into the solution by their privileged topological properties. On the other end, conservative solutions, like SConES GI and heinz might not shed much light on the etiology of the disease.

To overcome the problems posed by the individual methods while exploiting their strengths, we propose combining them into a consensus subnetwork. We use a straightforward aggregation to generate it, including any node that was recovered by at least two methods. The resulting network is a synthesis of the altered mechanism: it is smaller than the largest solutions (SigMod and dmGWAS), which makes it more manageable, and includes the majority of the strongly associated smaller solutions (SConES and heinz). The consensus subnetwork captures mechanisms and genes known to be related to cancer, recovering familial breast cancer genes as well as genome regions associated to cancer susceptibility. However, thanks to its small size and its network structure, it provides compelling hypotheses of non-canonical mechanisms involved in carcinogenesis, like mitochondrial translation and chaperone activity.

In order to produce the consensus network, we had to face the different interfaces, preprocessing steps, and unexpected behaviors of the various methods. To facilitate that other authors apply them to new datasets and aggregate their solutions, we built six nextflow pipelines [46] with a consistent interface and, whenever possible, parallelized computation. They are available on GitHub: `https://github.com/hclimente/gwas-tools`. Importantly, those methods that had a permissive license were compiled into a Docker image for easier use, which is available on Docker Hub hclimente/gwas-tools.

# 5    Funding and acknowledgments

# References

[1] W. S. Bush and J. H. Moore, "Chapter 11: Genome-Wide Association Studies," *PLoS Computational Biology*, vol. 8, p. e1002822, Dec. 2012. 00001.

[2] A. Buniello, J. A. MacArthur, M. Cerezo, L. W. Harris, J. Hayhurst, C. Malangone, A. McMahon, J. Morales, E. Mountjoy, E. Sollis, D. Suveges, O. Vrousgou, P. L. Whetzel, R. Amode, J. A. Guillen, H. S. Riat, S. J. Trevanion, P. Hall, H. Junkins, P. Flicek, T. Burdett, L. A. Hindorff, F. Cunningham, and H. Parkinson, "The NHGRI-EBI GWAS Catalog of published genome-wide association studies, targeted arrays and summary statistics 2019," *Nucleic Acids Research*, vol. 47, pp. D1005–D1012, Jan. 2019. 00092.

[3] P. M. Visscher, N. R. Wray, Q. Zhang, P. Sklar, M. I. McCarthy, M. A. Brown, and J. Yang, "10 Years of GWAS Discovery: Biology, Function, and Translation," *The American Journal of Human Genetics*, vol. 101, pp. 5–22, July 2017. 00634.

[4] M. H. Wang, H. J. Cordell, and K. Van Steen, "Statistical methods for genome-wide association studies," *Seminars in Cancer Biology*, May 2018. 00001.

[5] N. Barton, A. Etheridge, and A. Véber, "The infinitesimal model: Definition, derivation, and implications," *Theoretical Population Biology*, vol. 118, pp. 50–73, Dec. 2017. 00054.

[6] E. A. Boyle, Y. I. Li, and J. K. Pritchard, "An Expanded View of Complex Traits: From Polygenic to Omnigenic," *Cell*, vol. 169, pp. 1177–1186, June 2017. 00586.

[7] L. I. Furlong, "Human diseases through the lens of network biology," *Trends in Genetics*, vol. 29, pp. 150–159, Mar. 2013. 00128.

[8] A.-L. Barabási, N. Gulbahce, and J. Loscalzo, "Network medicine: a network-based approach to human disease," *Nature Reviews Genetics*, vol. 12, pp. 56–68, Jan. 2011. 02826.

[9] J. Piñero, A. Berenstein, A. Gonzalez-Perez, A. Chernomoretz, and L. I. Furlong, "Uncovering disease mechanisms through network biology in the era of Next Generation Sequencing," *Scientific Reports*, vol. 6, p. 24570, Apr. 2016. 00016.

[10] J. K. Huang, D. E. Carlin, M. K. Yu, W. Zhang, J. F. Kreisberg, P. Tamayo, and T. Ideker, "Systematic Evaluation of Molecular Networks for Discovery of Disease Genes," *Cell Systems*, vol. 6, pp. 484–495.e5, Apr. 2018. 00024.

[11] C.-A. Azencott, "Network-Guided Biomarker Discovery," in *Machine Learning for Health Informatics*, vol. 9605, pp. 319–336, Cham: Springer International Publishing, 2016. 00000.

[12] F. Gwinner, G. Boulday, C. Vandiedonck, M. Arnould, C. Cardoso, I. Nikolayeva, O. Guitart-Pla, C. V. Denis, O. D. Christophe, J. Beghain, E. Tournier-Lasserve, and B. Schwikowski, "Network-based analysis of omics data: The LEAN method," *Bioinformatics*, p. btw676, Oct. 2016. 00007.

[13] P. Jia, S. Zheng, J. Long, W. Zheng, and Z. Zhao, "dmGWAS: dense module searching for genome-wide association studies in protein–protein interaction networks," *Bioinformatics*, vol. 27, pp. 95–102, Jan. 2011. 00205.

[14] M. T. Dittrich, G. W. Klau, A. Rosenwald, T. Dandekar, and T. Muller, "Identifying functional modules in protein-protein interaction networks: an integrated exact approach," *Bioinformatics*, vol. 24, pp. i223–i231, July 2008. 00429.

[15] M. D. M. Leiserson, F. Vandin, H.-T. Wu, J. R. Dobson, J. V. Eldridge, J. L. Thomas, A. Papoutsaki, Y. Kim, B. Niu, M. McLellan, M. S. Lawrence, A. Gonzalez-Perez, D. Tamborero, Y. Cheng, G. A. Ryslik, N. Lopez-Bigas, G. Getz, L. Ding, and B. J. Raphael, "Pan-cancer network analysis identifies combinations of rare somatic mutations across pathways and protein complexes," *Nature Genetics*, vol. 47, pp. 106–114, Feb. 2015. 00411.

[16] C.-A. Azencott, D. Grimm, M. Sugiyama, Y. Kawahara, and K. M. Borgwardt, "Efficient network-guided multi-locus association mapping with graph cuts," *Bioinformatics*, vol. 29, pp. i171–i179, July 2013. 00047.

[17] Y. Liu, M. Brossard, D. Roqueiro, P. Margaritte-Jeannin, C. Sarnowski, E. Bouzigon, and F. Demenais, "SigMod: an exact and efficient method to identify a strongly interconnected disease-associated module in a gene network," *Bioinformatics*, p. btx004, Jan. 2017. 00007.

[18] O. M. Sinilnikova, M.-G. Dondon, S. Eon-Marchais, F. Damiola, L. Barjhoux, M. Marcou, C. Verny-Pierre, V. Sornin, L. Toulemonde, J. Beauvallet, D. Le Gal, N. Mebirouk, M. Belotti, O. Caron, M. Gauthier-Villars, I. Coupier, B. Buecher, A. Lortholary, C. Dugast, P. Gesta, J.-P. Fricker, C. Noguès, L. Faivre, E. Luporsi, P. Berthet, C. Delnatte, V. Bonadona, C. M. Maugard, P. Pujol, C. Lasset, M. Longy, Y.-J. Bignon, C. Adenis, L. Venat-Bouvet, L. Demange, H. Dreyfus, M. Frenay, L. Gladieff, I. Mortemousque, S. Audebert-Bellanger, F. Soubrier, S. Giraud, S. Lejeune-Dumoulin, A. Chevrier, J.-M. Limacher, J. Chiesa, A. Fajac, A. Floquet, F. Eisinger, J. Tinat, C. Colas, S. Fert-Ferrer, C. Penet, T. Frebourg, M.-A. Collonge-Rame, E. Barouk-Simonet, V. Layet, D. Leroux, O. Cohen-Haguenauer, F. Prieur, E. Mouret-Fourme, F. Cornélis, P. Jonveaux, O. Bera, E. Cavaciuti, A. Tardivon, F. Lesueur, S. Mazoyer, D. Stoppa-Lyonnet, and N. Andrieu, "GENESIS: a French national resource to study the missing heritability of breast cancer," *BMC Cancer*, vol. 16, p. 13, Dec. 2016. 00005.

[19] L. C. Sakoda, E. Jorgenson, and J. S. Witte, "Turning of COGS moves forward findings for hormonally mediated cancers," *Nature Genetics*, vol. 45, pp. 345–348, Apr. 2013. 00060.

[20] C. C. Chang, C. C. Chow, L. C. Tellier, S. Vattikuti, S. M. Purcell, and J. J. Lee, "Second-generation PLINK: rising to the challenge of larger and richer datasets," *GigaScience*, vol. 4, p. 7, Dec. 2015. 01610.

[21] A. Mishra and S. Macgregor, "VEGAS2: Software for More Flexible Gene-Based Testing," *Twin Research and Human Genetics*, vol. 18, pp. 86–91, Feb. 2015. 00125.

[22] A. Frankish, M. Diekhans, A.-M. Ferreira, R. Johnson, I. Jungreis, J. Loveland, J. M. Mudge, C. Sisu, J. Wright, J. Armstrong, I. Barnes, A. Berry, A. Bignell, S. Carbonell Sala, J. Chrast, F. Cunningham,

T. Di Domenico, S. Donaldson, I. T. Fiddes, C. García Girón, J. M. Gonzalez, T. Grego, M. Hardy, T. Hourlier, T. Hunt, O. G. Izuogu, J. Lagarde, F. J. Martin, L. Martínez, S. Mohanan, P. Muir, F. C. Navarro, A. Parker, B. Pei, F. Pozo, M. Ruffier, B. M. Schmitt, E. Stapleton, M.-M. Suner, I. Sycheva, B. Uszczynska-Ratajczak, J. Xu, A. Yates, D. Zerbino, Y. Zhang, B. Aken, J. S. Choudhary, M. Gerstein, R. Guigó, T. J. Hubbard, M. Kellis, B. Paten, A. Reymond, M. L. Tress, and P. Flicek, "GENCODE reference annotation for the human and mouse genomes," *Nucleic Acids Research*, vol. 47, pp. D766–D773, Jan. 2019. 00063.

[23] Q. Wang and P. Jia, "dmgwas 3.0." `https://bioinfo.uth.edu/dmGWAS/`, 2014. Accessed: 2019-07-16.

[24] I. Ljubić, R. Weiskircher, U. Pferschy, G. W. Klau, P. Mutzel, and M. Fischetti, "An Algorithmic Framework for the Exact Solution of the Prize-Collecting Steiner Tree Problem," *Mathematical Programming*, vol. 105, pp. 427–449, Feb. 2006. 00223.

[25] D. Beisser, G. W. Klau, T. Dandekar, T. Muller, and M. T. Dittrich, "BioNet: an R-Package for the functional analysis of biological networks," *Bioinformatics*, vol. 26, pp. 1129–1130, Apr. 2010. 00188.

[26] M. Dittrich and D. Beisser, "Bionet." `https://bioconductor.org/packages/BioNet/`, 2008. Accessed: 2019-07-16.

[27] P. Nakka, B. J. Raphael, and S. Ramachandran, "Gene and Network Analysis of Common Variants Reveals Novel Associations in Multiple Complex Diseases," *Genetics*, vol. 204, pp. 783–798, Oct. 2016. 00015.

[28] S. Scheid and R. Spang, "twilight; a Bioconductor package for estimating the local false discovery rate," *Bioinformatics*, vol. 21, pp. 2921–2922, June 2005. 00054.

[29] M. D. M. Leiserson, F. Vandin, H.-T. Wu, J. R. Dobson, J. V. Eldridge, J. L. Thomas, A. Papoutsaki, Y. Kim, B. Niu, M. McLellan, M. S. Lawrence, A. Gonzalez-Perez, D. Tamborero, Y. Cheng, G. A. Ryslik, N. Lopez-Bigas, G. Getz, L. Ding, and B. J. Raphael, "Hotnet2." `https://github.com/raphael-group/hotnet2`, 2018. Accessed: 2019-07-16.

[30] F. Gwinner, "Leanr." `https://cran.r-project.org/web/packages/LEANR/`, 2016. Accessed: 2019-07-16.

[31] H. Climente-González and C.-A. Azencott, "martini." `https://www.bioconductor.org/packages/martini/`, 2019. Accessed: 2019-07-16.

[32] Y. Liu, "Sigmod v2." `https://github.com/YuanlongLiu/SigMod`, 2018. Accessed: 2019-07-16.

[33] J. Das and H. Yu, "HINT: High-quality protein interactomes and their applications in understanding human disease," *BMC Systems Biology*, vol. 6, no. 1, p. 92, 2012. 00204.

[34] S. Nogueira and G. Brown, "Measuring the Stability of Feature Selection," in *Machine Learning and Knowledge Discovery in Databases*, vol. 9852, pp. 442–457, Cham: Springer International Publishing, 2016. 00000.

[35] J. Piñero, Bravo, N. Queralt-Rosinach, A. Gutiérrez-Sacristán, J. Deu-Pons, E. Centeno, J. García-García, F. Sanz, and L. I. Furlong, "DisGeNET: a comprehensive platform integrating information on human disease-associated genes and variants," *Nucleic Acids Research*, vol. 45, pp. D833–D839, Jan. 2017. 00369.

[36] K. Michailidou, J. Beesley, S. Lindstrom, S. Canisius, J. Dennis, M. J. Lush, M. J. Maranian, M. K. Bolla, Q. Wang, M. Shah, B. J. Perkins, K. Czene, M. Eriksson, H. Darabi, J. S. Brand, S. E. Bojesen, B. G. Nordestgaard, H. Flyger, S. F. Nielsen, N. Rahman, C. Turnbull, O. Fletcher, J. Peto, L. Gibson, I. dos Santos-Silva, J. Chang-Claude, D. Flesch-Janys, A. Rudolph, U. Eilber, S. Behrens, H. Nevanlinna, T. A. Muranen, K. Aittomäki, C. Blomqvist, S. Khan, K. Aaltonen, H. Ahsan, M. G. Kibriya, A. S. Whittemore, E. M. John, K. E. Malone, M. D. Gammon, R. M. Santella, G. Ursin, E. Makalic, D. F. Schmidt, G. Casey, D. J. Hunter, S. M. Gapstur, M. M. Gaudet, W. R. Diver, C. A. Haiman, F. Schumacher, B. E. Henderson, L. Le Marchand, C. D. Berg, S. J. Chanock, J. Figueroa, R. N. Hoover, D. Lambrechts, P. Neven, H. Wildiers, E. van Limbergen, M. K. Schmidt, A. Broeks, S. Verhoef, S. Cornelissen, F. J. Couch, J. E. Olson, E. Hallberg, C. Vachon, Q. Waisfisz, H. Meijers-Heijboer, M. A. Adank, R. B. van der Luijt, J. Li, J. Liu, K. Humphreys, D. Kang, J.-Y. Choi, S. K. Park, K.-Y. Yoo, K. Matsuo, H. Ito, H. Iwata, K. Tajima, P. Guénel, T. Truong, C. Mulot, M. Sanchez, B. Burwinkel, F. Marme, H. Surowy, C. Sohn, A. H. Wu, C.-c. Tseng, D. Van Den Berg, D. O. Stram, A. González-Neira, J. Benitez, M. P. Zamora, J. I. A. Perez, X.-O. Shu, W. Lu, Y.-T. Gao, H. Cai, A. Cox, S. S. Cross, M. W. R. Reed, I. L. Andrulis, J. A. Knight, G. Glendon, A. M. Mulligan, E. J. Sawyer, I. Tomlinson, M. J. Kerin, N. Miller, A. Lindblom, S. Margolin, S. H. Teo, C. H. Yip, N. A. M. Taib, G.-H. Tan, M. J. Hooning, A. Hollestelle, J. W. M. Martens, J. M. Collée, W. Blot, L. B. Signorello, Q. Cai, J. L. Hopper, M. C. Southey, H. Tsimiklis, C. Apicella, C.-Y. Shen, C.-N. Hsiung, P.-E. Wu, M.-F. Hou, V. N. Kristensen, S. Nord, G. I. G. Alnaes, G. G. Giles, R. L. Milne, C. McLean, F. Canzian, D. Trichopoulos, P. Peeters, E. Lund, M. Sund, K.-T. Khaw, M. J. Gunter, D. Palli, L. M. Mortensen, L. Dossus, J.-M. Huerta, A. Meindl, R. K. Schmutzler, C. Sutter, R. Yang, K. Muir, A. Lophatananon, S. Stewart-Brown, P. Siriwanarangsan, M. Hartman, H. Miao, K. S. Chia, C. W. Chan, P. A. Fasching, A. Hein, M. W. Beckmann, L. Haeberle, H. Brenner, A. K. Dieffenbach, V. Arndt, C. Stegmaier, A. Ashworth, N. Orr, M. J. Schoemaker, A. J. Swerdlow, L. Brinton, M. Garcia-Closas, W. Zheng, S. L. Halverson, M. Shrubsole, J. Long, M. S. Goldberg, F. Labrèche, M. Dumont, R. Winqvist, K. Pylkäs, A. Jukkola-Vuorinen, M. Grip, H. Brauch, U. Hamann, T. Brüning, P. Radice, P. Peterlongo, S. Manoukian, L. Bernard, N. V. Bogdanova, T. Dörk, A. Mannermaa, V. Kataja, V.-M. Kosma, J. M. Hartikainen, P. Devilee, R. A. E. M. Tollenaar, C. Seynaeve, C. J. Van Asperen, A. Jakubowska, J. Lubinski, K. Jaworska, T. Huzarski, S. Sangrajrang, V. Gaborieau, P. Brennan, J. McKay, S. Slager, A. E. Toland, C. B. Ambrosone, D. Yannoukakos, M. Kabisch, D. Torres, S. L. Neuhausen, H. Anton-Culver, C. Luccarini, C. Baynes, S. Ahmed, C. S. Healey, D. C. Tessier, D. Vincent, F. Bacot, G. Pita, M. R. Alonso, N. Álvarez, D. Herrero, J. Simard, P. P. D. P. Pharoah, P. Kraft, A. M. Dunning, G. Chenevix-Trench, P. Hall, and D. F. Easton, "Genome-wide association analysis of more than 120,000 individuals identifies 15 new susceptibility loci for breast cancer," *Nature Genetics*, vol. 47, pp. 373–380, Apr. 2015. 00000.

[37] E. S. Rinella, Y. Shao, L. Yackowski, S. Pramanik, R. Oratz, F. Schnabel, S. Guha, C. LeDuc, C. L. Campbell, S. D. Klugman, M. B. Terry, R. T. Senie, I. L. Andrulis, M. Daly, E. M. John, D. Roses, W. K. Chung, and H. Ostrer, "Genetic variants associated with breast cancer risk for Ashkenazi Jewish women with strong family histories but no identifiable BRCA1/2 mutation," *Human Genetics*, vol. 132, pp. 523–536, May 2013. 00019.

[38] A. G. Brisbin, Y. W. Asmann, H. Song, Y.-Y. Tsai, J. A. Aakre, P. Yang, R. B. Jenkins, P. Pharoah, F. Schumacher, D. V. Conti, D. J. Duggan, M. Jenkins, J. Hopper, S. Gallinger, P. Newcomb, G. Casey, T. A. Sellers, and B. L. Fridley, "Meta-analysis of 8q24 for seven cancers reveals a locus between NOV and ENPP2 associated with cancer development," *BMC Medical Genetics*, vol. 12, p. 156, Dec. 2011. 00033.

[39] SEARCH, The GENICA Consortium, kConFab, Australian Ovarian Cancer Study Group, S. Ahmed, G. Thomas, M. Ghoussaini, C. S. Healey, M. K. Humphreys, R. Platte, J. Morrison, M. Maranian, K. A. Pooley, R. Luben, D. Eccles, D. G. Evans, O. Fletcher, N. Johnson, I. dos Santos Silva, J. Peto, M. R. Stratton, N. Rahman, K. Jacobs, R. Prentice, G. L. Anderson, A. Rajkovic, J. D. Curb, R. G. Ziegler, C. D. Berg, S. S. Buys, C. A. McCarty, H. S. Feigelson, E. E. Calle, M. J. Thun, W. R. Diver, S. Bojesen, B. G. Nordestgaard, H. Flyger, T. Dörk, P. Schürmann, P. Hillemanns, J. H. Karstens, N. V. Bogdanova, N. N. Antonenkova, I. V. Zalutsky, M. Bermisheva, S. Fedorova, E. Khusnutdinova, D. Kang, K.-Y. Yoo, D. Y. Noh, S.-H. Ahn, P. Devilee, C. J. van Asperen, R. A. E. M. Tollenaar, C. Seynaeve, M. Garcia-Closas, J. Lissowska, L. Brinton, B. Peplonska, H. Nevanlinna, T. Heikkinen, K. Aittomäki, C. Blomqvist, J. L. Hopper, M. C. Southey, L. Smith, A. B. Spurdle, M. K. Schmidt, A. Broeks, R. R. van Hien, S. Cornelissen, R. L. Milne, G. Ribas, A. González-Neira, J. Benitez, R. K. Schmutzler, B. Burwinkel, C. R. Bartram, A. Meindl, H. Brauch, C. Justenhoven, U. Hamann, J. Chang-Claude, R. Hein, S. Wang-Gohrke, A. Lindblom, S. Margolin, A. Mannermaa, V.-M. Kosma, V. Kataja, J. E. Olson, X. Wang, Z. Fredericksen, G. G. Giles, G. Severi, L. Baglietto, D. R. English, S. E. Hankinson, D. G. Cox, P. Kraft, L. J. Vatten, K. Hveem, M. Kumle, A. Sigurdson, M. Doody, P. Bhatti, B. H. Alexander, M. J. Hooning, A. M. W. van den Ouweland, R. A. Oldenburg, M. Schutte, P. Hall, K. Czene, J. Liu, Y. Li, A. Cox, G. Elliott, I. Brock, M. W. R. Reed, C.-Y. Shen, J.-C. Yu, G.-C. Hsu, S.-T. Chen, H. Anton-Culver, A. Ziogas, I. L. Andrulis, J. A. Knight, J. Beesley, E. L. Goode, F. Couch, G. Chenevix-Trench, R. N. Hoover, B. A. J. Ponder, D. J. Hunter, P. D. P. Pharoah, A. M. Dunning, S. J. Chanock, and D. F. Easton, "Newly discovered breast cancer susceptibility loci on 3p24 and 17q23.2," *Nature Genetics*, vol. 41, pp. 585–590, May 2009. 00000.

[40] S. K. Calderwood and J. Gong, "Heat Shock Proteins Promote Cancer: It's a Protection Racket," *Trends in Biochemical Sciences*, vol. 41, pp. 311–323, Apr. 2016. 00105.

[41] S. Ahmed, G. Thomas, M. Ghoussaini, C. S. Healey, M. K. Humphreys, R. Platte, J. Morrison, M. Maranian, K. A. Pooley, R. Luben, D. Eccles, D. G. Evans, O. Fletcher, N. Johnson, I. dos Santos Silva, J. Peto, M. R. Stratton, N. Rahman, K. Jacobs, R. Prentice, G. L. Anderson, A. Rajkovic, J. D. Curb, R. G. Ziegler, C. D. Berg, S. S. Buys, C. A. McCarty, H. S. Feigelson, E. E. Calle, M. J. Thun, W. R. Diver, S. Bojesen, B. G. Nordestgaard, H. Flyger, T. Dörk, P. Schürmann, P. Hillemanns, J. H. Karstens, N. V. Bogdanova, N. N. Antonenkova, I. V. Zalutsky, M. Bermisheva, S. Fedorova, E. Khusnutdinova, D. Kang, K.-Y. Yoo, D. Y. Noh, S.-H. Ahn, P. Devilee, C. J. van Asperen, R. A. E. M. Tollenaar, C. Seynaeve, M. Garcia-Closas, J. Lissowska, L. Brinton, B. Peplonska, H. Nevanlinna, T. Heikkinen, K. Aittomäki, C. Blomqvist, J. L. Hopper, M. C. Southey, L. Smith, A. B. Spurdle, M. K. Schmidt, A. Broeks, R. R. van Hien, S. Cornelissen, R. L. Milne, G. Ribas, A. González-Neira, J. Benitez, R. K. Schmutzler, B. Burwinkel, C. R. Bartram, A. Meindl, H. Brauch, C. Justenhoven, U. Hamann, J. Chang-Claude, R. Hein, S. Wang-Gohrke, A. Lindblom, S. Margolin, A. Mannermaa, V.-M. Kosma, V. Kataja, J. E. Olson, X. Wang, Z. Fredericksen, G. G. Giles, G. Severi, L. Baglietto, D. R. English, S. E. Hankinson, D. G. Cox, P. Kraft, L. J. Vatten, K. Hveem, M. Kumle, A. Sigurdson, M. Doody, P. Bhatti, B. H. Alexander, M. J. Hooning, A. M. W. van den Ouweland, R. A. Oldenburg, M. Schutte, P. Hall, K. Czene, J. Liu, Y. Li, A. Cox, G. Elliott, I. Brock, M. W. R. Reed, C.-Y. Shen, J.-C. Yu, G.-C. Hsu, S.-T. Chen, H. Anton-Culver, A. Ziogas, I. L. Andrulis, J. A. Knight, J. Beesley, E. L. Goode, F. Couch, G. Chenevix-Trench, R. N. Hoover, B. A. J. Ponder, D. J. Hunter, P. D. P. Pharoah, A. M. Dunning, S. J. Chanock, and D. F. Easton, "Newly discovered breast cancer susceptibility loci on 3p24 and 17q23.2," *Nature Genetics*, vol. 41, pp. 585–590, May 2009. 00000.

[42] D. A. Quigley, E. Fiorito, S. Nord, P. Van Loo, G. G. Alnaes, T. Fleischer, J. Tost, H. K. Moen Vollan, T. Tramm, J. Overgaard, I. R. Bukholm, A. Hurtado, A. Balmain, A.-L. Børresen-Dale, and V. Kristensen, "The 5p12 breast cancer susceptibility locus affects MRPS30 expression in estrogen-receptor positive tumors," *Molecular Oncology*, vol. 8, pp. 273–284, Mar. 2014. 00000.

[43] J. Breyer, D. Dorset, T. Clark, K. Bradley, T. Wahlfors, K. McReynolds, W. Maynard, S. Chang, M. Cookson, J. Smith, J. Schleutker, W. Dupont, and J. Smith, "An Expressed Retrogene of the Master Embryonic Stem Cell Gene POU5f1 Is Associated with Prostate Cancer Susceptibility," *The American Journal of Human Genetics*, vol. 94, pp. 395–404, Mar. 2014. 00018.

[44] G. Liu, F. X. Claret, F. Zhou, and Y. Pan, "Jab1/COPS5 as a Novel Biomarker for Diagnosis, Prognosis, Therapy Prediction and Therapeutic Tools for Human Cancer," *Frontiers in Pharmacology*, vol. 9, p. 135, Feb. 2018. 00005.

[45] J. J. Cai, E. Borenstein, and D. A. Petrov, "Broker Genes in Human Disease," *Genome Biology and Evolution*, vol. 2, pp. 815–825, Jan. 2010. 00060.

[46] P. Di Tommaso, M. Chatzou, E. W. Floden, P. P. Barja, E. Palumbo, and C. Notredame, "Nextflow enables reproducible computational workflows," *Nature Biotechnology*, vol. 35, pp. 316–319, Apr. 2017. 00176.

# Supplementary materials

**Table 1:** Summary statistics on the results of SConES on the three SNP-SNP interaction networks. The first row within each block contains the summary statistics on the whole network.
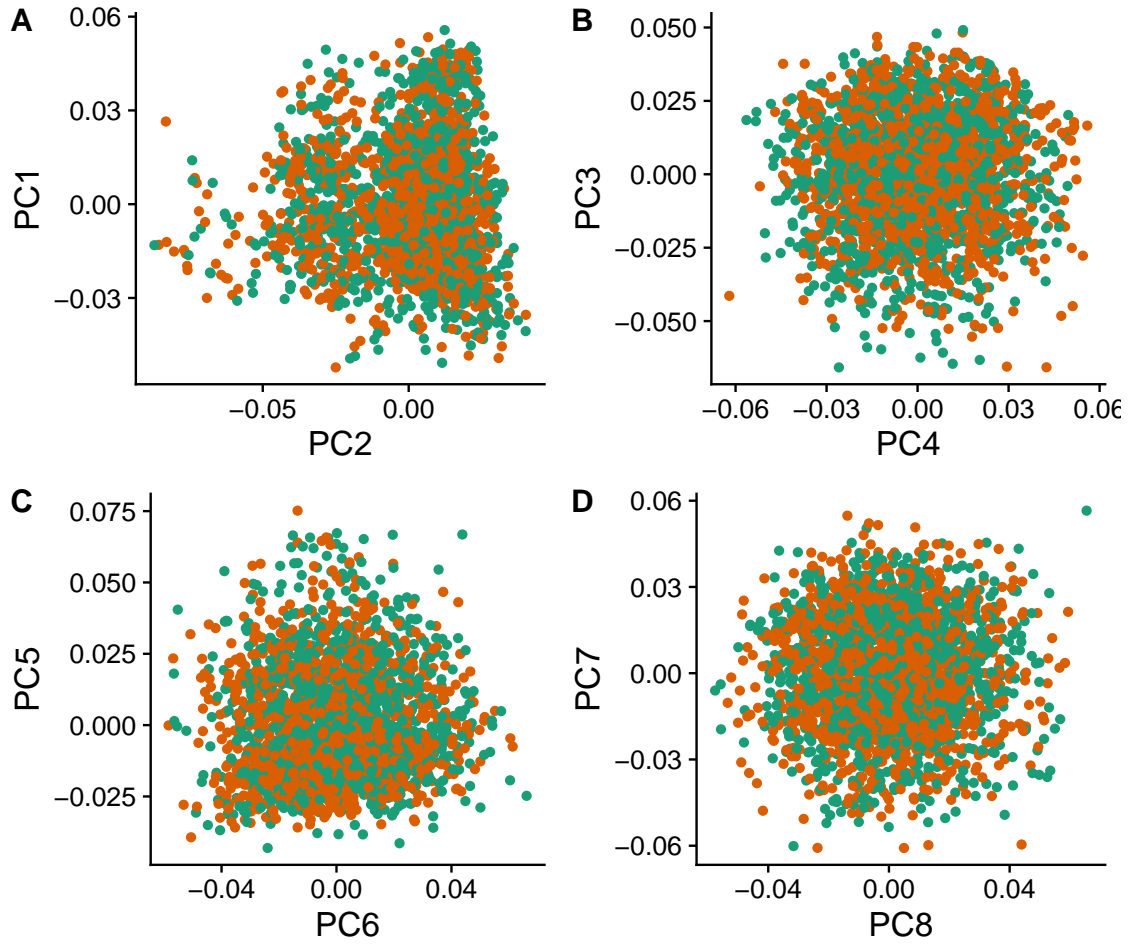
| Network | SNPs | Edges | Subnetworks | $\overline{\text{Betweenness}}$ | $\hat{\text{P}}_{\text{SNP}}$ |
|---|---|---|---|---|---|
| GS | 197 083 | 197 060 | - | $2.03 \times 10^7$ | 0.49 |
| SConES GS | 1 590 | 1 585 | 5 | $2.52 \times 10^7$ | 0.023 |
| GM | 197 083 | 6 442 446 | - | $3.99 \times 10^6$ | 0.49 |
| SConES GM | 1 692 | 177 611 | 5 | $4.40 \times 10^6$ | 0.055 |
| GI | 197 083 | 28 733 720 | - | $1.46 \times 10^6$ | 0.49 |
| SConES GI | 408 | 539 | 5 | $9.33 \times 10^6$ | 0.076 |

$\overline{\text{Betweenness}}$: mean betweenness of the selected SNPs in the corresponding full network; $\hat{\text{P}}_{\text{SNP}}$: median P-value of the selected SNPs.
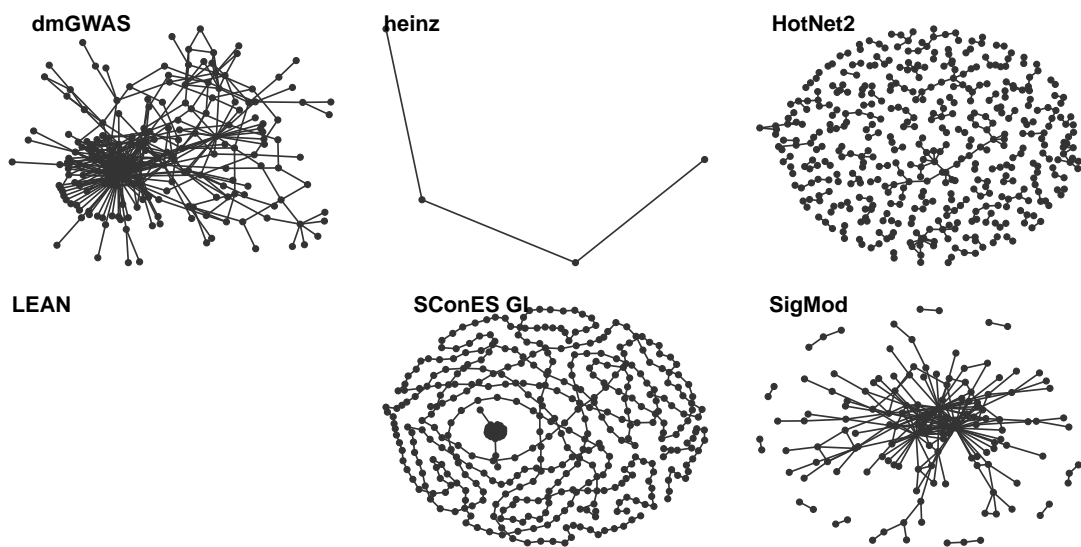
**Table 2:** Summary statistics on the results of multiple network methods on the gene-gene interaction network. The first row contains the summary statistics on the whole network.

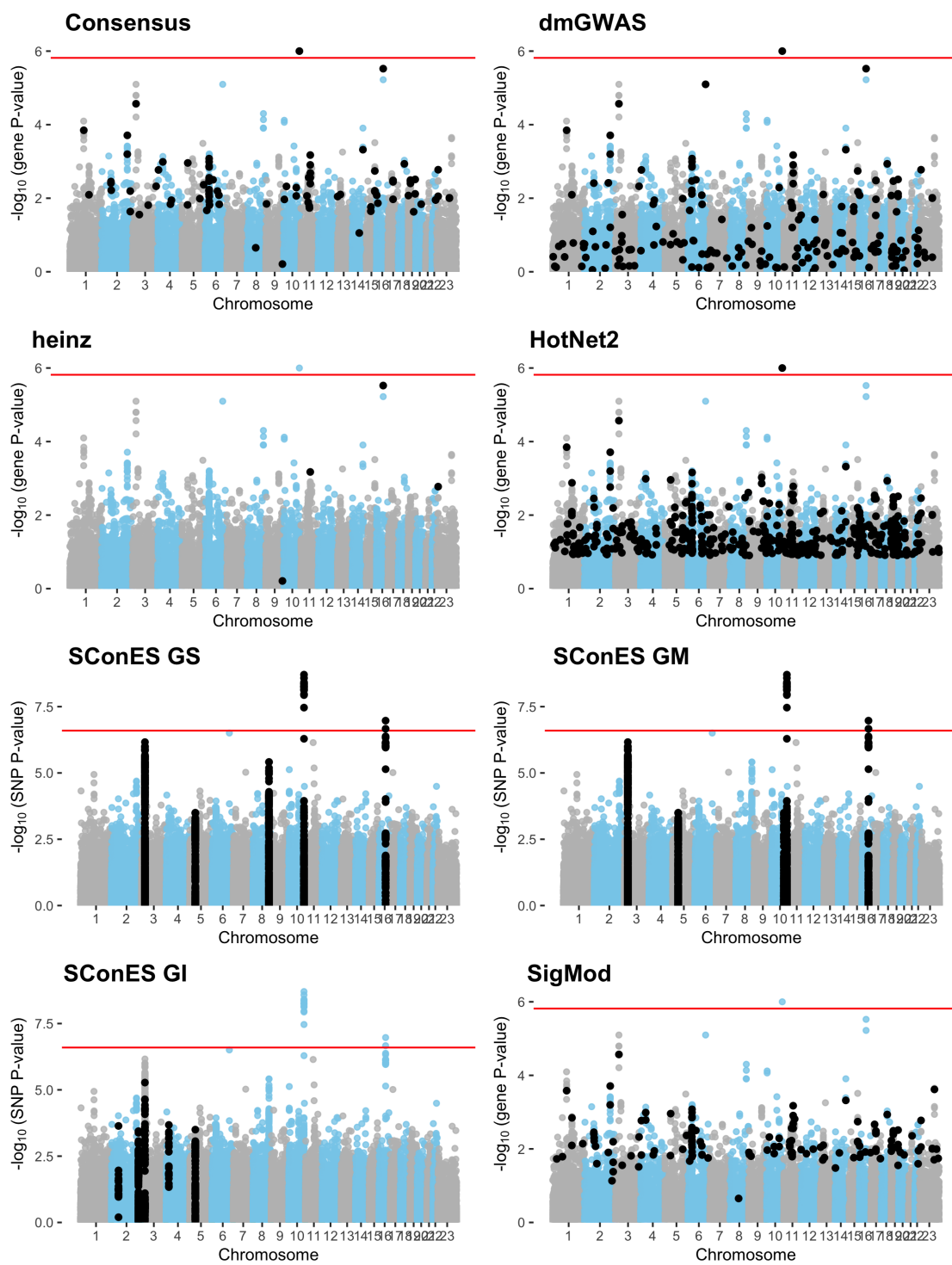| Network | Genes | Edges | $\overline{\text{Betweenness}}$ | $\hat{\text{P}}_{\text{gene}}$ | $\rho_{consensus}$ |
|---|---|---|---|---|---|
| SConES GS | 5 | 0 | 9805 | $2.7 \times 10^{-5}$ | 0.19 |
| SConES GM | 28 | 2 | 4267 | 0.067 | 0.12 |

$\overline{\text{Betweenness}}$: mean betweenness of the selected genes in the full network; $\hat{\text{P}}_{\text{gene}}$: median P-value of the selected genes; $\rho_{consensus}$: Pearson's correlation with the consensus network.
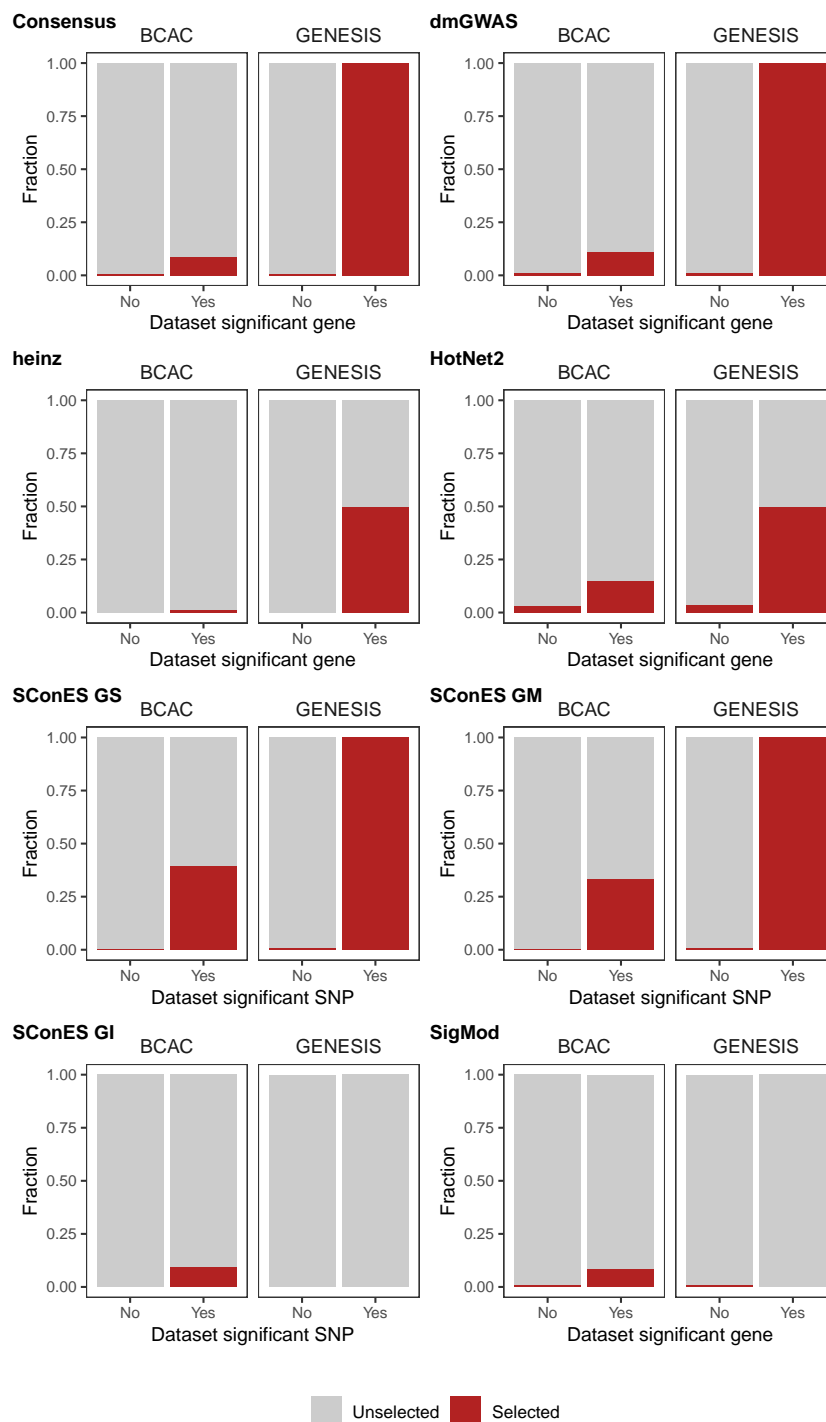
**Fig. 1:** GENESIS shows no differential population structure between cases and controls. **(A,B,C,D)** Eight main principal components computed on the genotypes of GENESIS. Cases are colored in green, controls in orange.
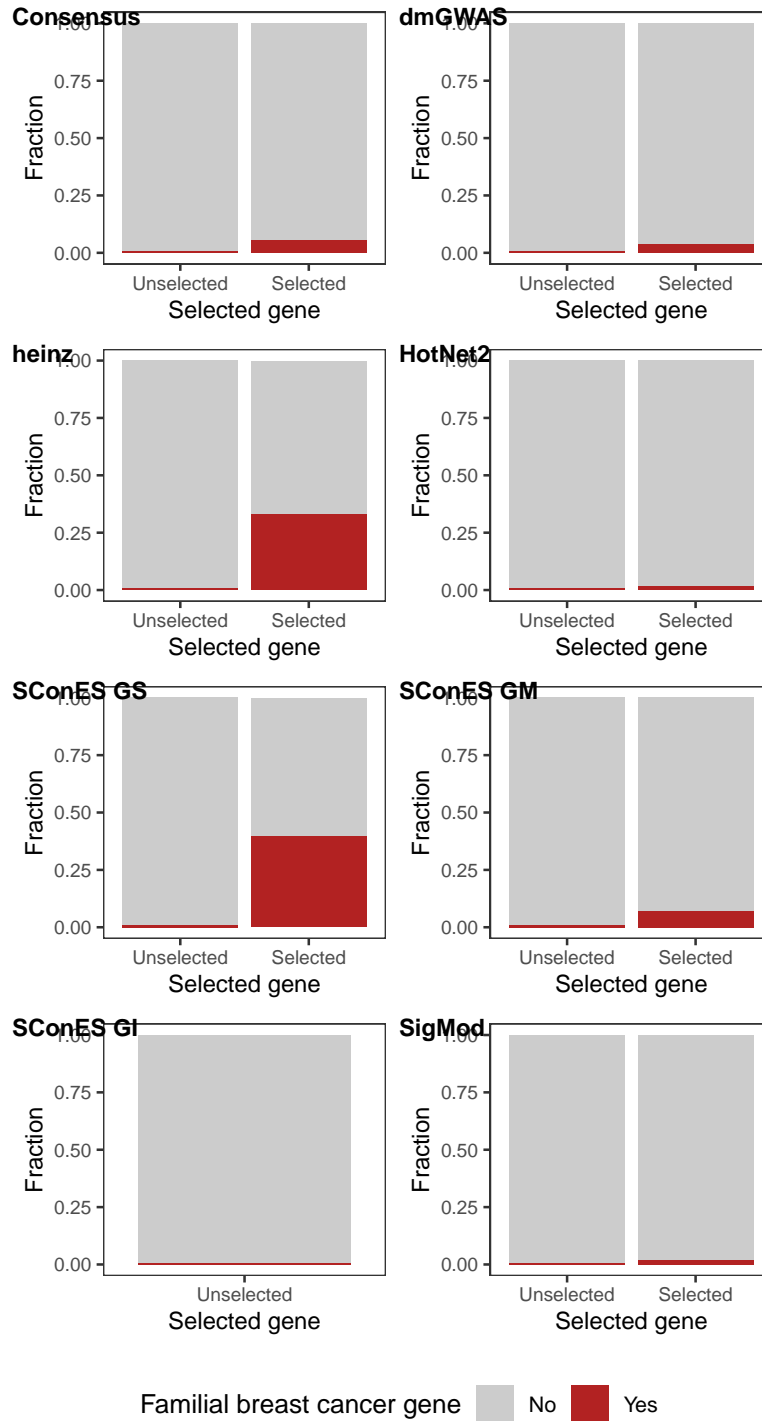
**Fig. 2:** Overview of the subnetworks produced by the different network methods. **(dmGWAS, heinz, HotNet2, LEAN, and SigMod)** contain gene subnetworks; **(SConES GI)**, SNP subnetworks.
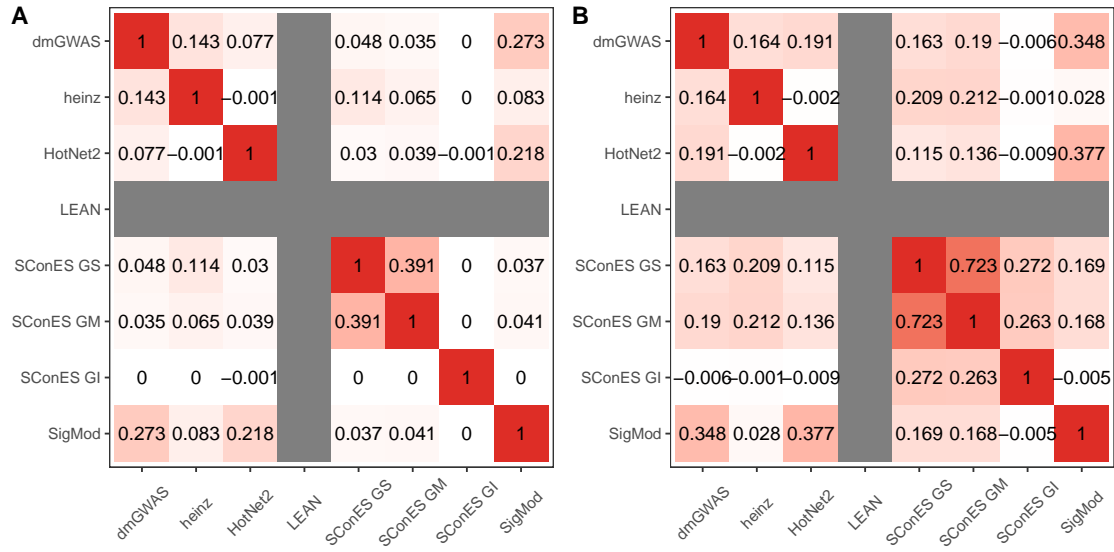
**Fig. 3:** Manhattan plots showing the biomolecules selected by each method. In **(Consensus, dmGWAS, heinz, HotNet2, and SigMod)** datapoints are genes; in **(SConES GS, GM, and GI)**, SNPs. LEAN was excluded, as it did not select any gene.

**Fig. 4:** Proportion of the Bonferroni significant biomolecules (in either the GENESIS or the BCAC datasets) selected by each of the methods on the GENESIS data. **(Consensus, dmGWAS, heinz, HotNet2, and SigMod)** involve significant genes, only among those present in the protein-protein interaction network. **(SConES GS, GM and GI)** involve significant SNPs. LEAN was excluded, as it did not select any gene. The presented network methods recover a higher proportion of significant genes than of non-significant genes in both datasets, despite their lack of significance in GENESIS.
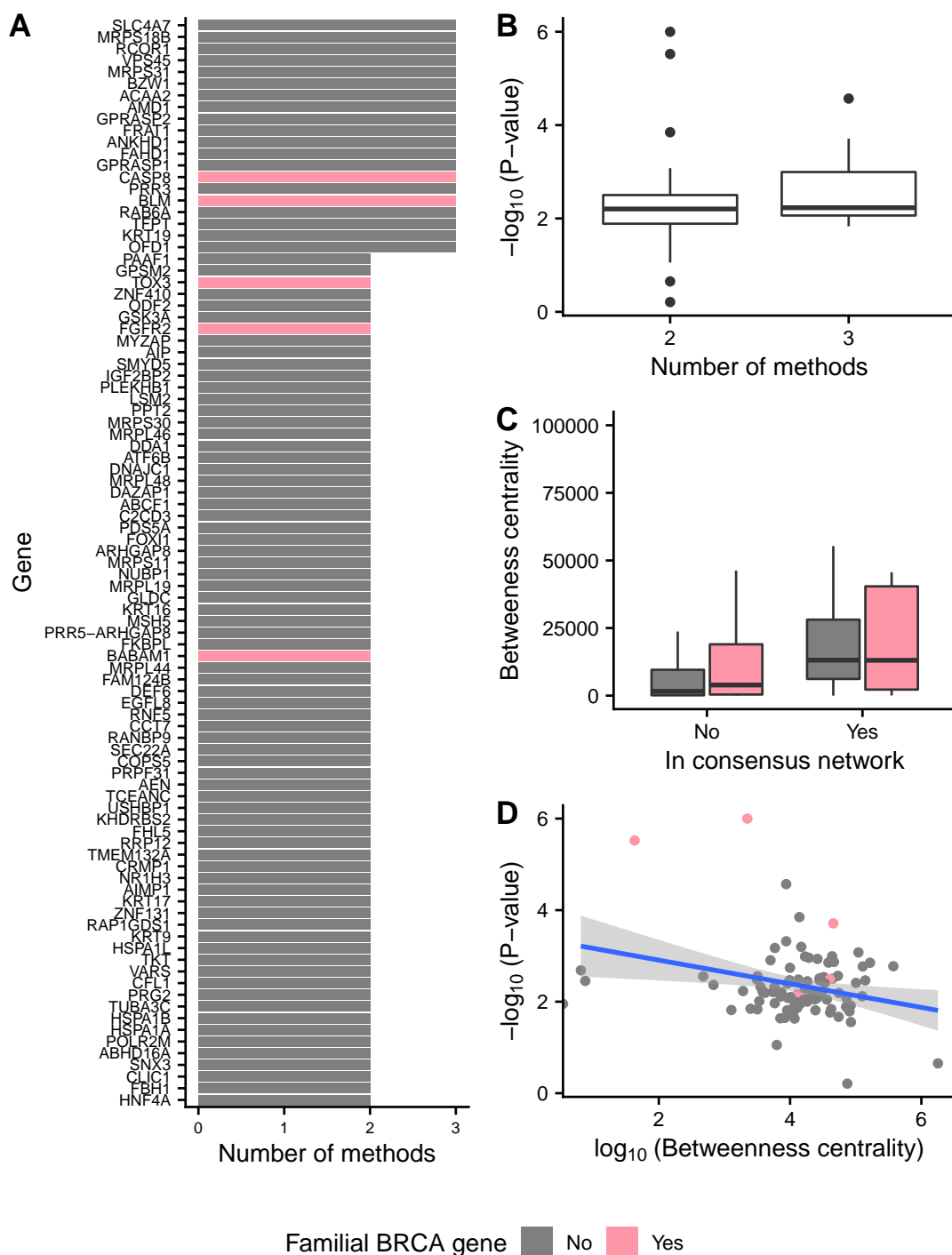
**Fig. 5:** :Proportion of the selected genes by each of the methods on the GENESIS data that is a known familial breast cancer gene (Section 2.4.2). Only genes present in the protein-protein interaction network were considered. LEAN was excluded, as it did not select any gene. The presented network methods recover a higher proportion of familial breast cancer genes than of other genes, despite their lack of significance in GENESIS.
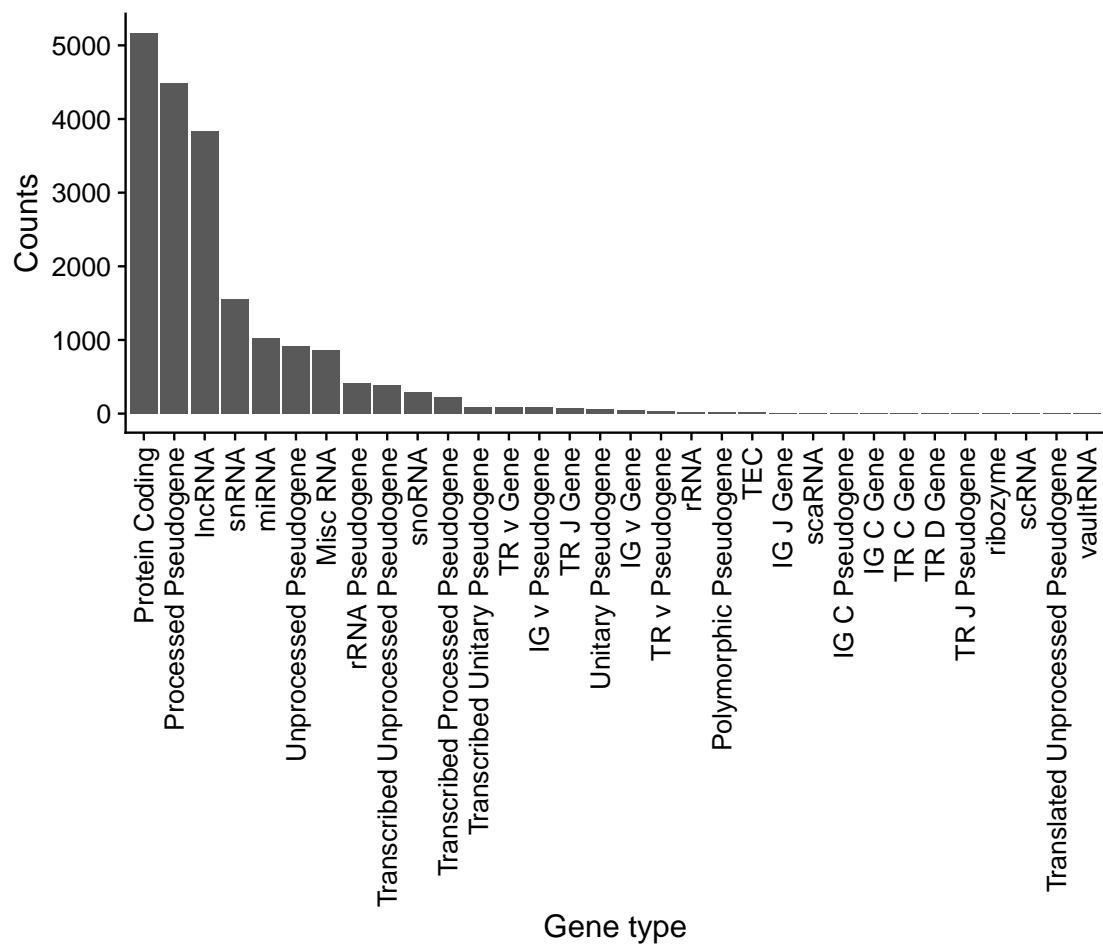
**Fig. 6:** Pearson's correlation between the different solution subnetworks. **(A)** Correlation between selected SNPs. **(B)** Correlation between selected genes. In general, the solutions display a very low overlap.

**Fig. 7:** Consensus subnetwork on GENESIS (Section 2.3.6). **(A)** Each node is represented by a pie chart, which accounts the methods that selected it. The labeled genes have a VEGAS2v2 P-value < 0.001 and/or are known familial breast cancer genes (colored in pink). This panel is equivalent to Figure 4. **(B)** The name of every gene is indicated.
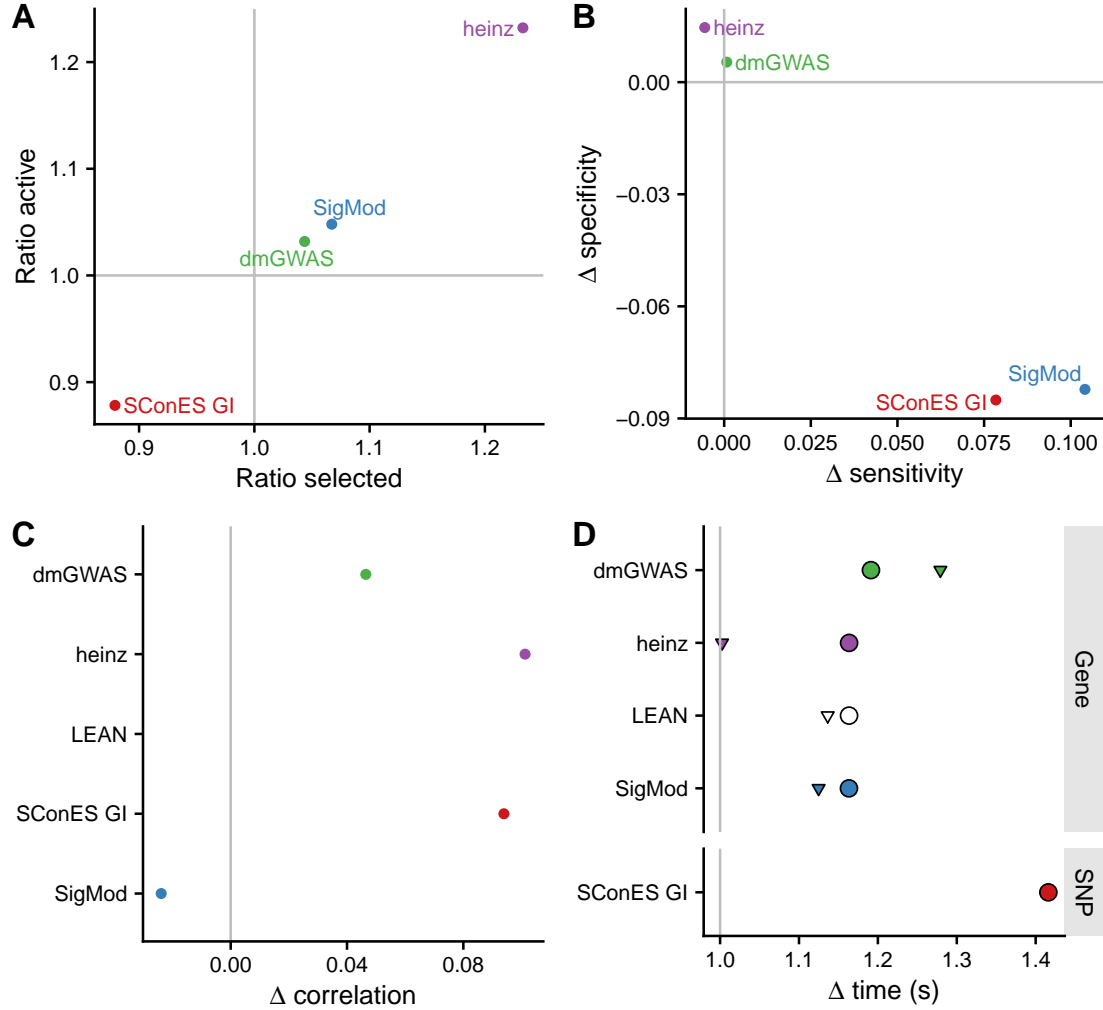
**Fig. 8:** Genes on the consensus network. Familial breast cancer genes are colored in pink; the rest are colored in grey. **(A)** Number of methods selecting every gene in the subnetwork. **(B)** VEGAS P-values of association of the genes, with regards to the number of methods that selected them. **(C)** Comparison of betweenness centrality of the genes in the consensus network and the other genes in the PPIN and not in the consensus network. To improve visualization, we removed outliers. **(D)** Relationship between the $\log_{10}$ of the betweenness centrality and the $-\log_{10}$ of the VEGAS P-value of the genes in the consensus network. The blue line represents a fitted generalized linear model.

**Fig. 9:** Biotypes of genes from the annotation that are not present in the HINT protein-protein interaction network.

**Fig. 10:** Comparison of benchmark on high-throughput interactions to benchmark on both high-throughput and literature curated interactions. Grey lines represent no change between the benchmarks (1 for ratios, 0 for differences). **(A)** Ratios of the selected features between both benchmarks and of the active set. **(B)** Shifts in sensitivity and specificity. **(C)** Shift in Pearson's correlation between benchmarks. **(D)** Ratio between the runtimes of the benchmarks. For gene network-based methods, inverted triangles represent the ratio of runtimes of the algorithms themselves, and circles the total time, which includes the algorithm themselves and the additional 119980 seconds (1 day and 9.33 hours) which took VEGAS2v2 on average to compute the gene scores from SNP summary statistics. In general, adding additional interactions slightly improves the stability of the solution, but increases the solution size, has mixed effects on the sensitivity and specificity, and impacts negatively the required runtime of the algorithms.