

Biological networks and GWAS: comparing and combining network methods to understand the genetics of familial breast cancer susceptibility in the GENESIS study

Héctor Climente-González^{1,2,3,4*}, Christine Lonjou^{1,2,3}, Fabienne Lesueur^{1,2,3}, Dominique Stoppa-Lyonnet^{5,6,7}✉, Nadine Andrieu^{1,2,3}, Chloé-Agathe Azencott^{3,1,2}
GENESIS study group¹

1 Institut Curie, PSL Research University, F-75005 Paris, France;

2 INSERM, U900, F-75005 Paris, France;

3 MINES ParisTech, PSL Research University, CBIO-Centre for Computational Biology, F-75006 Paris, France;

4 RIKEN Center for Advanced Intelligence Project (AIP), Tokyo, Japan;

5 Service de Génétique, Institut Curie, F-75005 Paris, France;

6 INSERM, U830, F-75005 Paris, France;

7 Université Paris Descartes.

✉For the GENESIS study group

* Membership list can be found in the Acknowledgments section.

* hector.climente(at)riken.jp

Abstract

Systems biology provides a comprehensive approach to biomarker discovery and biological hypothesis building. Indeed, it allows to jointly consider the statistical association between gene variation and a phenotype. Network approaches to disease use biological networks, which model functional relationships between the molecules in a cell, to generate hypotheses about the genetics of complex diseases. Several among them jointly consider gene scores, representing the association between each gene and the disease, and the biological context of each gene, represented as modeled by a network. In this work, we study six network methods which identify subnetworks with high association scores to a phenotype. Specifically, we examine their utility to discover new biomarkers for breast cancer susceptibility by interrogating such network methods using gene scores from GENESIS, a genome-wide association study (GWAS) focused on French women with a family history of breast cancer and tested negative for pathogenic variants in *BRCA1* and *BRCA2*. 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. By trading statistical stringency for biological meaningfulness, most network methods give non-BRCA familial breast cancer. We provide a critical comparison of these six methods, discussing the impact of their mathematical formulation and parameters. Using a biological network yields more compelling results than standard SNP- and gene-level analyses, recovering causal subnetworks tightly related to GWAS analyses. Indeed, we find significant overlaps between our solutions and the genes identified in the largest GWAS on breast cancer susceptibility. For instance, we show a general alteration of the neighborhood of *COPS5*. We further propose to combine these solutions into a consensus network, which brings further insights. The consensus network contains *COPS5*, a gene related to multiple hallmarks of cancer. Importantly, we find a significantly large overlap between the

genes in the solution networks and the genes significantly associated in the largest GWAS on susceptibility to breast cancer. Yet, network methods are notably unstable, producing different results when the input data changes slightly. To account for that, we produce, and 14 of its neighbors. The main drawback of network methods is that they are not robust to small perturbations in their inputs. Therefore, we propose a stable consensus subnetwork solution, formed by the most consistently selected genes. The stable consensus in multiple subsamples of the data. In GENESIS, it is composed of 68 genes, enriched in known breast cancer susceptibility genes (*BLM*, *CASP8*, *CASP10*, *DNAJC1*, *FGFR2*, *MRPS30*, and *SLC4A7*, Fisher's exact test P-value = 3×10^{-4}) and occupying more central positions in the network than average most genes. The network seems is organized around *CUL3*, encoding an ubiquitin ligase related protein that regulates the protein levels which is involved in the regulation of several genes involved in linked to cancer progression. In conclusion, this article shows the pertinence of network-based analyses to tackle known issues with GWAS, namely we showed how network methods help overcome the lack of statistical power and of interpretable solutions of GWAS and improve their interpretation. Project-agnostic implementations of each of the network all methods are available at <https://github.com/hclimente/gwas-tools> to facilitate their application to other GWAS datasets.

Author summary

In genome-wide Genome-wide association studies (GWAS), scan thousands of genomes are scanned to identify variants associated with a complex trait. Over the last 15 years, GWAS have advanced our understanding of the genetics of complex diseases, and in particular of hereditary cancers. Yet However, they have led to an apparent paradox: the more we perform such studies, the more it seems that the entire genome is involved in every disease. An elegant explanation has been proposed with the omnigenic model. The omnigenic model offers an appealing explanation: only a limited number of core core genes are directly involved in the disease, but gene functions are deeply interrelated, so that and so many other genes are able to can alter the function of the core genes. These interrelations are often modeled as networks, and multiple algorithms have been proposed to use these networks to identify the subset of core genes involved in a specific trait. In this study, we characterize This study applies and compares six such network methods on GENESIS, a GWAS dataset for familial breast cancer in the French population. Combining these approaches allows us to identify potentially novel breast cancer susceptibility genes, and provides a mechanistic explanation for their role in the development of the disease. Our pipeline can easily be applied to other diseases. We provide ready-to-use implementations of all the examined methods.

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 from the disease of interest, and the controls, taken from the general population. Then, a per SNP statistical association test is conducted statistical test of association (e.g. based on logistic regression) is conducted between each SNP and the phenotype. 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~~ and deepening our understanding of the disease.

GWAS have successfully identified thousands of variants underlying many common diseases [2]. However, this experimental setting also presents intrinsic inherent 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 it can only detect~~ variants with larger effects ~~can be detected~~ [3]. This is particularly problematic since the prevailing view is that most genetic architectures involve many variants with small effects [3]. Additionally, to avoid false positives, most GWAS apply a conservative multiple test correction ~~is applied~~, typically the previously mentioned Bonferroni correction. However, Bonferroni correction is overly conservative when the statistical tests ~~are correlated, as is the case correlate, as happens~~ in GWAS [4]. Another open issue is the interpretation of the results, as the functional consequences of most common variants are ~~not well understood~~ unknown. On top of that, recent large-sampled studies suggest that numerous loci spread all along the genome contribute to a degree to any complex trait, in accordance with the infinitesimal model [5]. The recently proposed omnigenic model [6] offers an explanation: genes are ~~very functionally inter-related strongly interrelated~~ and influence each other's ~~behavior~~ function, which allows alterations in most genes to impact the subset of "core" genes directly involved in the disease's mechanism. Hence, a comprehensive statistical framework ~~which that~~ includes the structure of biological data might ~~address the aforementioned issues help alleviate the issues above~~.

For this reason, many authors turn to network biology to handle the complex interplay of biomolecules that lead to disease ~~[7]~~ [7, 8]. As its name suggests, network biology models biology as a network, where the biomolecules under study, often genes, are nodes, and selected functional relationships are edges that link them. These relationships come from evidence that the genes jointly contribute to a biological function; for instance, their ~~expression is expressions are~~ correlated, or their 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 molecules [9]. In fact, an examination of biological networks shows that disease genes have differential properties ~~[9, 10]~~. ~~This is particularly true for cancer driver genes, which tend to be key players in connecting different, densely-connected communities of genes. Therefore, studying the neighborhood of disease-associated genes is effective at identifying new ones that are involved in the same biological functions~~ ~~[11]~~ [9, 12]: they tend to occupy central positions in the network (although not the most central ones); disease genes for the same pathology tend to cluster in modules; and often they are bottlenecks that interconnect modules.

Network-based ~~biomarker~~ discovery methods exploit ~~this relatedness the differential properties described above~~ to identify disease genes on GWAS data ~~+13~~ [11, 13]. In essence, each SNP ~~has gene receives~~ a score of association with the disease, ~~given by the experiment, and functionally computed from the GWAS data, and a set of~~ biological relationships, given by a network built on prior knowledge. Then, the problem becomes finding a functionally-related set of highly-scoring genes. ~~Different Multiple~~ 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~~ For example, some methods restrict the problem to ~~certain kinds specific types~~ of subnetworks. Such is the ~~extreme~~ case of LEAN [14], which focuses on "star" subnetworks, i.e. ~~instances were~~, instances where both a gene and its direct interactors are associated with the disease. Other algorithms, like dmGWAS [15] and heinz [16], ~~focus on larger do not impose such strong constraints and search for~~ subnetworks interconnecting genes with high association scores. However,

they differ in their tolerance to the inclusion of low-scoring nodes, and the topology of the solution. Lastly, other methods also consider the topology of the network, favoring groups of nodes that are not only high-scoring, but also densely interconnected; such is the case of HotNet2 [17], SConES [18], and SigMod [19].

In this work, we analyze the effectiveness of studied the relevance of network-based approaches to genetics by applying these six network methods for biomarker discovery on to GWAS data. While all of them capture susceptibility mechanisms resembling that postulated by They use different interpretations of the omnigenic model, they do so in different ways, and provide a representative view of the field. We worked on the GENESIS dataset [20], a study on familial breast cancer conducted in the French population. After a classical GWAS approach, we used these network methods to recover identify additional breast cancer biomarkers susceptibility genes. Lastly, we carry out a comparison of compared the solutions obtained by the different methods, and aggregate and aggregated them to obtain a consensus network consensus solutions of predisposition to familial breast cancer that addressed their shortcomings.

2 Results

2.1 Conventional SNP- and gene-based analyses confirm that retrieve the *FGFR2* locus is associated with familial breast cancer in the GENESIS dataset

We conducted association analyses in the GENESIS dataset (Section 4.1) at both SNP and gene levels (Section ??). Two genomic regions have 4.2.1. At the SNP level, two genomic regions had a P-value lower than the Bonferroni threshold on chromosomes 10 and 16 (S2 FigA). The former overlaps with the gene *FGFR2*, the latter with *CASC16*, and it is located near and the protein-coding gene *TOX3*. Variants in both *FGFR2* and *TOX3* have been repeatedly associated with breast cancer susceptibility in other case-control studies [21], *BRCA1* and *BRCA2* carrier studies [22], and in hereditary breast and ovarian cancer families negative for mutations in *BRCA1* and *BRCA2* [23]. In our studied population At the gene level, only *FGFR2* was significantly associated with breast cancer at the gene level (S2 FigB).

Closer examination reveals two revealed two other regions (3p24 and 8q24) having low, albeit not genome-wide significant, P-values. Both of them have been associated to with breast cancer susceptibility in the past [24, 25]. We applied an L1-penalized logistic regression on the whole dataset (using all GENESIS genotypes as input and the phenotype (cancer/healthy) as the outcome (Section 4.5.2). It The algorithm selected 100 SNPs, both from all aforementioned regions mentioned above and new ones (S2 FigC). Yet, it is However, it was unclear why those SNPs were selected, as emphasized by the high P-value of some of them, which further complicates the biological interpretation. Moreover, and in opposition to what would be expected under the omnigenic model, the genes to which these SNPs map to (Section ??) are 4.3.5 were not interconnected in the PPIN (protein-protein interaction network (PPIN, Section 4.3.2). In addition Moreover, the classification performance of the method is very low, and L1-penalized methods select only one of several correlated variables and are prone to instability, which further complicates interpretation. This motivates model was low (sensitivity = 55%, specificity = 55%, Section 4.5). Together, these issues motivated exploring network methods, which trade statistical significance for biological relevance consider not only statistical association but also the location of each gene in a PPIN to find susceptibility subnetworks. In fact, such methods provided comparably (poor) classification performance to L1-penalized logistic regression (Fig 3B), while providing more interpretable solutions. genes.

2.2 Network methods successfully identify genes associated with breast cancer

We applied six network methods to the GENESIS dataset (Section 4.3.3). As none of the networks examined by LEAN was significant (BH Benjamini-Hochberg [BH] correction adjusted P-value ≤ 0.05), we obtained ~~six~~ five solutions (Fig 1): one for each of the remaining four gene-based methods, ~~and~~ one for SConES GI (which works at the SNP level), ~~and the consensus. These solutions differ.~~

~~These solutions differed~~ in many aspects, making it hard to draw joint conclusions. For starters, the overlap between the genes featured in each solution ~~is quite~~ was relatively small (Fig 1A). However, the methods tended to agree on the genes with the strongest signal: genes selected by more methods tended to have lower P-value of association (Fig 1B).

~~Another prominent difference is their~~

~~Another major difference was the solution size:~~ the largest solution, produced by HotNet2, ~~contains~~ contained 440 genes, while heinz's contained only 4 genes. While SConES GI failed to ~~did not~~ recover any protein coding gene, ~~by dealing~~ working with SNP networks ~~it retrieved~~ rather than gene networks allowed it to retrieve four subnetworks in intergenic regions, ~~and another one~~ and another subnetwork overlapping an RNA gene (*RNU6-420P*).

~~Their topologies also differ~~

~~The topologies of the five solutions differed as well (Fig 1C)~~, as measured by the median centrality (Table 1) and the number of connected components (Fig 1D). Only two methods have Table 1. Three methods yielded more than one connected component: SConES, as described above, SigMod, and HotNet2. HotNet2 produced 135 subnetworks, 115 of which have less ~~fewer~~ than five genes. The second largest subnetwork (13 nodes) ~~contains~~ contained the two breast cancer susceptibility genes *CASP8* and *BLM*. Lastly, a pathway enrichment analysis (Section 4.4) also showed similarities and differences in the underlying mechanisms. It linked different parts of SigMod's solution network to four processes: protein translation (including mitochondrial), mRNA splicing, protein misfolding, and keratinization (BH adjusted P-values < 0.03). Interestingly, dmGWAS solution is also related to protein misfolding (*attenuation phase*, BH adjusted P-value = 0.01). But, additionally, it includes submodules of proteins related to mitosis, DNA damage, and regulation of *TP53* (BH adjusted P-values < 0.05), which match previously known mechanisms of breast cancer susceptibility [26]. As with SigMod, the genes in HotNet2's solution are involved in mitochondrial translation (BH adjusted P-value = 1.87×10^{-4}), but also in glycogen metabolism and transcription of nuclear receptors (BH adjusted P-value < 0.04). 4.6).

~~Overview of the solutions produced by the different network methods (Section 4.3.3) on the GENESIS dataset. As LEAN did not produce any significant gene (BH adjusted P-value > 0.05), it was excluded. Unless indicated otherwise, results refer to genes, except for SConES GI which are at the SNP level. (A) Manhattan plots of SNPs/genes; in black, the method's solution. The Bonferroni threshold is indicated by a red line (2.54×10^{-7} for SNPs, 1.53×10^{-6} for genes). (B) Overlap between the genes selected by each of the methods, measured by Pearson correlation between indicator vectors. (C) Precision and recall of the evaluated methods with respect to Bonferroni-significant SNPs/genes in BCAC. For reference, we added a gray line with a slope of 1. (D) Solution networks produced by the different methods. Lastly, a pathway enrichment analysis (Section 4.4) also showed similarities and differences between the methods' solutions. It linked different parts of SigMod's solution to four processes (S2 Table): protein translation (including mitochondrial), mRNA splicing, protein misfolding, and keratinization (BH adjusted P-values < 0.03).~~

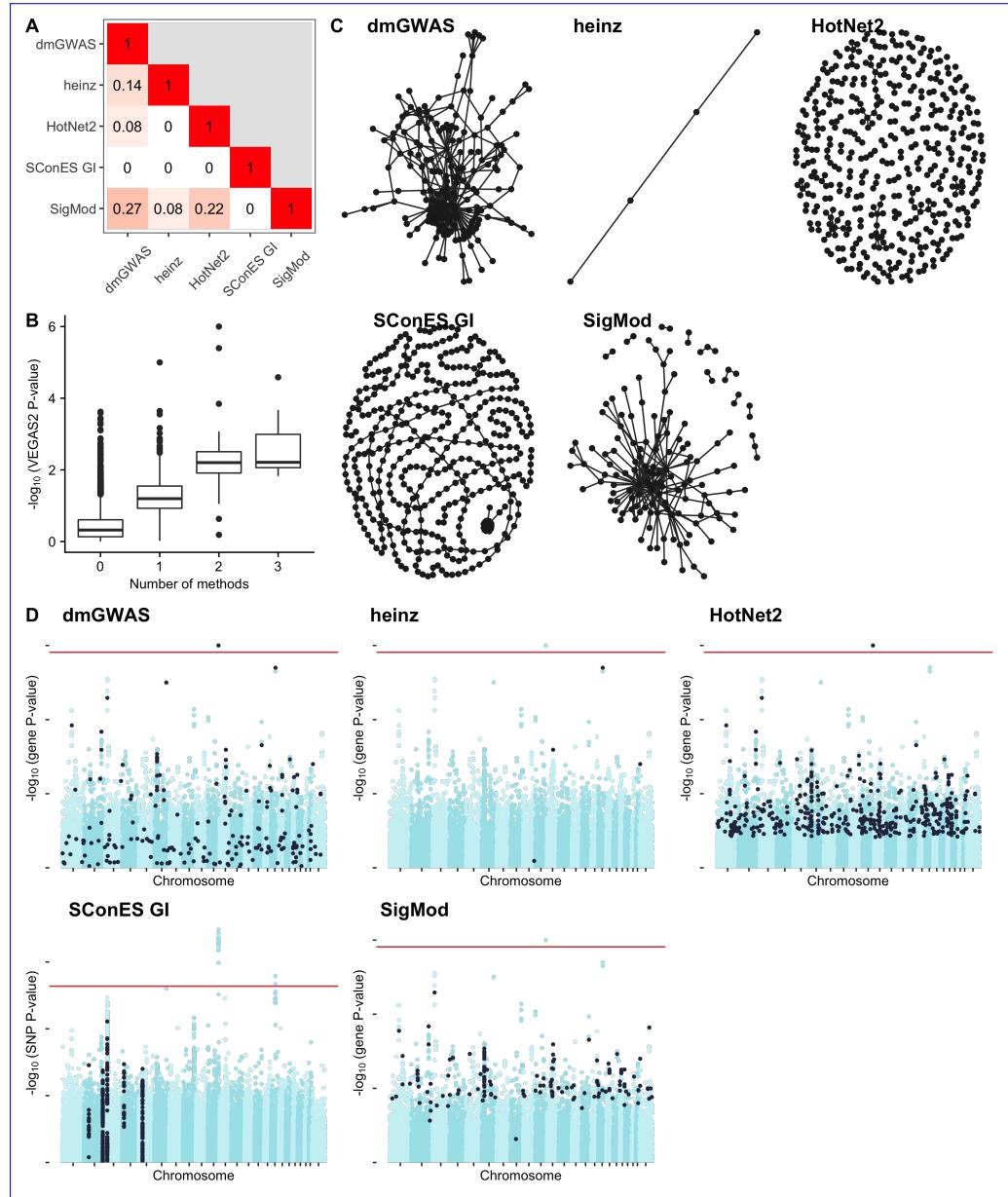


Fig 1. Overview of the solutions produced by the different network methods (Section 4.3.3) on the GENESIS dataset. As LEAN did not produce any significant solution (BH adjusted P-value < 0.05), it is not shown. Unless indicated otherwise, results refer to SNPs for SConES GI, and to genes for the other methods. (A) Overlap between the genes selected by each method, measured by Pearson correlation between indicator vectors (Sections 4.5.1 and 4.3.5). (B) Distribution of VEGAS2 P-values of the genes in the PPIN not selected by any network method (12 213), and of those selected by 1 (575), 2 (73), or 3 (20) methods. (C) Solution networks produced by the different methods. (D) Manhattan plots of SNPs/genes; in black, the method's solution. The red line indicates the Bonferroni threshold (2.54×10^{-7} for SNPs, 1.53×10^{-6} for genes).

Table 1. Summary statistics on the solutions of multiple network methods on the gene-gene interaction network PPIN. The first row contains the summary statistics on the whole network PPIN.

Network	# genes	# edges	# components	Betweenness	\hat{P}_{gene}	$p_{\text{consensus}}$	# genes in consensus
HINT HT	13 619	142 541	15	16 706	0.46		0.06693/93
dmGWAS	194	450	1	49 115	0.19		0.4155/93
heinz	4	3	1	113 633	0.001		0.214/93
HotNet2	440	374	130	7 739	0.048		0.3163/93
LEAN	0	0	-0	-	-		0/93
SConES GI	0 (1)	0	-0	-	-		0/93
SigMod	142	249	11	92 603	0.008		0.7384/93
Consensus	93	186	21	50 737	0.006		1/93/93
Stable consensus	68	49	32	94 854	0.005		0.5443/93

genes: number of genes selected out of those that are part of the PPIN; for SConES GI, the total number of genes, including RNA genes, was added in parentheses. # components: number of connected components.

Betweenness: mean – median betweenness of the selected genes in the PPIN. \hat{P}_{gene} : median – median VEGAS2 P-value of the selected genes. $p_{\text{consensus}}$: Pearson correlation between the subnetwork and the # genes in consensus network: number of genes in common between the method's solution and the 93 genes in the consensus solution.

Interestingly, the dmGWAS solution (S3 Table) was also related to protein misfolding (*attenuation phase*, BH adjusted P-value = 0.01). However, it additionally included submodules of proteins related to mitosis, DNA damage, and regulation of TP53 (BH adjusted P-values < 0.05), which match previously known mechanisms of breast cancer susceptibility [26]. As with SigMod, the genes in HotNet2's solution (S4 Table) were involved in mitochondrial translation (BH adjusted P-value = 1.87×10^{-4}), but also in glycogen metabolism and transcription of nuclear receptors (BH adjusted P-value < 0.04).

Despite their differences, there are additional common themes. All obtained solution subnetworks have solutions had lower association P-values than the whole PPIN (median VEGAS2 P-value $\ll 0.46$, Table 1), despite containing genes with higher P-values as well (Fig 1AD). This exemplifies illustrates the trade-off between statistical significance controlling for type I error and biological relevance. However, there are nuances between solutions in this regard: heinz strongly favored genes with lower P-values, while dmGWAS was less conservative (median VEGAS2 P-values 0.0012 and 0.19, respectively); SConES tended to select whole LD-blocks; and HotNet2 and SigMod were less likely to select low scoring genes.

Additionally, the solution subnetworks solutions presented other desirable properties. First, five four of them were enriched in known breast cancer susceptibility genes (consensus, dmGWAS, heinz, HotNet2, and SigMod, Fisher's exact test one-sided P-value < 0.03). Second, the genes in four solution subnetworks displayed on average a three solutions displayed, on average, a significantly higher betweenness centrality than the rest of the genes, a difference that is significant in four solutions (consensus, (dmGWAS, HotNet2, 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 [10], an observation that holds in breast cancer (one-tailed Wilcoxon rank-sum test P-value = 2.64×10^{-5} when comparing the betweenness of known susceptibility genes versus the rest). Interestingly, SConES selected SNPs that are the SNPs in SConES' solution were also more central than the average SNP (Supplementary table S1 Table), suggesting that causal SNPs are also more central than non-associated non-associated SNPs.

2.3 A case study: the consensus network solution

Despite the heterogeneity of the solutions, their shared properties suggest that each method captures the differences between the solutions suggested that each of them captured different aspects of cancer susceptibility. Indeed, only 20 genes are common to more than two solutions (A), but encouragingly, out of the 668 genes that were selected by at least one method, only 93 were selected by at least two, 20 by three, and none by four or more. Encouragingly, the more methods selected a gene, the higher its association score to the phenotype (B). To leverage on Fig 1B, a relationship that plateaued at 2. Hence, to leverage their strengths and compensate for their respective weaknesses, we built a consensus subnetwork that captures the mechanisms most shared among the solution subnetworks solution using the genes shared among at least two solutions (Section 4.3.3). This subnetwork solution (Fig 2) contains contained 93 genes and exhibits exhibited the aforementioned properties of the individual solutions: enrichment in breast cancer susceptibility genes and higher betweenness centrality than the rest of the genes.

A pathway enrichment analysis of the genes in the consensus network also shows solution also showed similar pathways as the individual solutions (S5 Table.). We found two involved mechanisms: *mitochondrial translation* and *attenuation phase*. The former is supported by genes like *MRPS30* (VEGAS2 P-value = 0.001), which encode a mitochondrial ribosomal protein and was also linked to breast cancer susceptibility [27]. Interestingly, increased mitochondrial translation has been found in cancer cells [28], and its inhibition was proposed as a therapeutic target. With regards to the attenuation phase of heat shock response, it involves involved three Hsp70 chaperones: *HSPA1A*, *HSPA1B*, and *HSPA1L*. The genes encoding these proteins are all near each other at 6p21, in the region known as HLA. In fact, out of the 22 SNPs that map mapped to any of these three genes, 9 map to all of them mapped to all three, and 4 to two, making which made it hard to disentangle their effects. *HSPA1A* was the most strongly associated gene (VEGAS2 P-value = 8.37×10^{-4}).

Topologically the consensus consists, the consensus consisted of a connected component composed of 49 genes, and multiple smaller subnetworks (Fig 2B and C). Among the latter, 19 genes are were in subnetworks containing a single gene or two connected nodes, implying that they do. This implied that they did not have a consistently altered neighborhood, but are but were strongly associated themselves and hence picked by at least two methods. The opposite would be the case of highly central genes large connected component contained genes that are highly central in the PPIN, a property which is weakly anti-correlated. This property weakly anticorrelated with the P-value of association to the disease (Pearson correlation coefficient = -0.26, S3 FigD). This suggests that they anticorrelation suggested that these genes were selected because they were on the shortest path between two highly associated genes. In view high scoring genes. Because of this, we hypothesize that highly central genes might contribute to the heritability through alterations of their neighborhood, eonsistently consistent with the omnigenic model of disease [6]. For instance, the most central node in the consensus network is COPS5 solution was *COPS5*, a component of the COP9 signalosome which regulates multiple signalling that regulates multiple signaling pathways. *COPS5* is related to multiple hallmarks of cancer and is overexpressed in multiple tumors, including breast and ovarian cancer [29]. Despite its lack of association in GENESIS or BCAC in studies conducted by the Breast Cancer Association Consortium (BCAC) [21] (VEGAS2 P-value of 0.22 and 0.14, respectively), its neighbors in the consensus subnetwork have solution had consistently low P-values (median VEGAS2 P-value = 0.006).

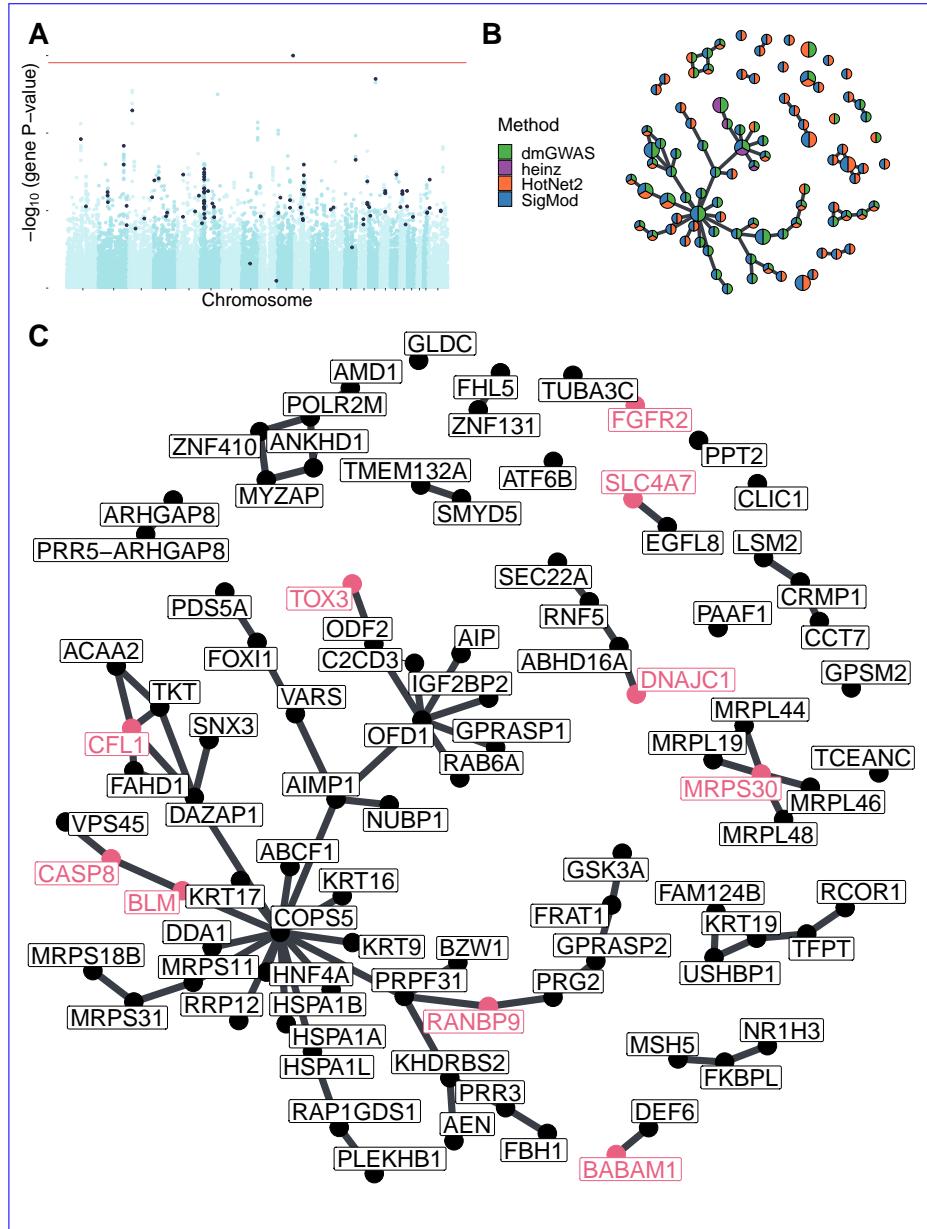


Fig 2. Consensus subnetwork on GENESIS Consensus solution on GENESIS (Section 4.3.3). (A) Manhattan plot of genes; in black, the ones in the consensus solution. The red line indicates the Bonferroni threshold (Section 1.53 4.3.3 $\times 10^{-6}$ for genes). (B) Consensus network. Each node-gene is represented by a pie chart, which shows the methods that selected it. We labeled (and enlarged) the two most central genes (*COPS5* and *OFD1*) and those genes that are known breast cancer susceptibility genes, and for significantly associated with breast cancer susceptibility in the BCAC dataset BCAC-significant genes (Section 4.6). (C) The latter ones nodes are also in the same disposition as in panel B, but we indicated every gene name. We colored in pink All gene the names are indicated in known breast cancer susceptibility genes and BCAC-significant genes.

2.4 Network methods boost biomarker discovery

We compared the results obtained with different network methods to the European sample of the Breast Cancer Association Consortium (BCAC) [21]BCAC, the largest GWAS to date on breast cancer (Section 4.6). Although BCAC case-control studies do not necessarily target cases with a family history of breast cancer like GENESIS does, this comparison is pertinent since we expect a shared genetic architecture at the gene level, at which most network methods operate. This shared genetic architecture, together methods operate. Together with BCAC's scale (90 times more samples than GENESIS) provides, this shared genetic architecture provided a reasonable counterfactual of what we would expect if GENESIS had a larger sample size. We computed a gene association score on BCAC, in an equivalent way to the one described in Section ??(Section 4.6). The solutions provided by the different network approaches overlap methods overlapped significantly with BCAC findings hits (Fisher's exact test P-value < 0.019). The gene-based network methods achieve methods achieved comparable precision (2%-25%) and recall (1.3-12.1%) at recovering BCAC-significant genes (Fig 1CS4 Fig.A). Interestingly, while SConES GI at the SNP-level achieves achieved a similar recall at the SNP-level (8.6%), it shows showed a much higher precision (47.3%).

2.5 Network methods share limitations

We compared the six network methods in a 5-fold subsampling setting (Section ??). Specifically, we measured five 4.5. In this comparison we measured properties (Fig 3) and S4 Fig.); the size of the solution subnetwork; sensitivity and the specificity of an L1-penalized logistic regression classifier on the selected SNPs; stability; and the stability of the methods; and their computational runtime. The solution size varies greatly between the different methods (Fig 3A). Heinz produced the smallest solutions, with an average of 182 selected SNPs. The largest solutions (Section 4.3.5) while the largest ones came from SConES GI (6 256.6 SNPs) and dmGWAS (4 255.0 SNPs). LEAN did not produce any solution in any of the subsamples.

To determine whether the selected SNPs could be used for patient classification predict cancer susceptibility, we computed the performance of the classifier on the classifiers' performances on test dataset sets (Fig 3S4 Fig.B). The different classifiers displayed similarly poor low sensitivities and specificities, all in the 0.52 – 0.56 range. Interestingly, the classifier trained on all the SNPs had a similar performance, despite being the only method aiming only at minimizing to minimize prediction error. It should be considered that Of course, although these performances are were low, we do did not expect to separate cases from controls well using exclusively genetic data [30].

Another desirable quality of an a selection algorithm is the stability of the solution with regards respect to small changes in the input (Section 4.5.1). Both heinz and LEAN displayed a high stability Heinz was highly stable in our benchmark, consistently selecting the same genes and no genes over the 5 subsamples, respectively (Fig 3C). Conversely, the while the other methods displayed similarly low stabilities (Fig 3B).

In terms of computational runtime, the fastest method was heinz (Fig 3DC), which returned a solution in a few seconds. HotNet2 was the slowest (3 days and 14 hours on average). Including the time required to compute the gene scores, however, slows slowed down considerably gene-based methods; on this benchmark, that step took on average 1 day and 9.33 hours. Considering that Including this first step, it took 5 days on average for HotNet2 to produce a result.

Using different combinations of parameters (Section 4.3.4), we computed how good each of the methods was at recovering the results of a conventional GWAS on BCAC (Section 4.6, Fig 3D). SConES exhibits the largest area under the curve since, when $\lambda = 0$ (i.e., network topology is disregarded), it is equivalent to a Bonferroni correction. The remaining network methods have similar areas under the curve, with heinz having the largest one.

2.6 Network topology matters, and association scores matter and might lead to ambiguous results

As shown above, and despite their similarities, the network methods produced different ways of modeling the problem led to remarkably different solutions. This is due to the particularities of each methods, and directly or indirectly provide information about the dataset. Importantly, understanding which assumptions the methods made allowed us to understand the results more in depth. For instance, the fact that LEAN did not return any biomarker implies that there is gene implied that there was no gene such that both itself and its environment are on average were, on average, strongly associated with the disease.

In this-the GENESIS dataset, heinz's solution is-was very conservative, providing a small solution with the lowest median P-value for the subnetwork (Table 1). Due to this parsimonious and highly associated solution, it was the best method to stably select a set of biomarkers. By repeatedly selecting this compact solution, heinz was the most stable method (Fig 3GB). Its conservativeness stems from its preprocessing step, which models modeled the gene P-values as a mixture model of a beta distribution and a uniform distribution, controlled by an FDR parameter. Due to the limited signal at the gene level in this dataset (S2 FigB), only 36 of all the genes retain genes retained a positive score after that transformation. Yet However, this small solution does-did not provide much insight into the susceptibility mechanisms to cancer. Importantly, it ignores genes that are associated to ignored genes that were associated with cancer in this dataset, like FGFR2.

On the other end of the spectrum, dmGWAS, HotNet2, and SigMod produced large solutions. dmGWAS' subnetwork is the least associated subnetwork on average. This is due to the greedy framework it uses, which has a bias for DmGWAS' solution was the lowest scoring solution on average because of its greedy framework, which is biased towards larger solutions [31]. It considered all nodes at distance 2 of the examined subnetwork, and accepted a weakly associated genes gene if it was linked to another, strongly associated one. This is exacerbated when high scoring one. Aggregating the results of successive greedy searches are aggregated exacerbates this bias, leading to a large, tightly connected cluster of unassociated genes (Fig 4A). 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 solution was even harder to interpret, being composed of 440 genes divided into 135 subnetworks. Lastly, SigMod misses missed some of the most strongly associated, highest scoring breast cancer susceptibility genes in the dataset, like FGFR2 and TOX3.

Another peculiarity of network methods is-was their relationship to degree centrality. On the one hand we observed that We studied random rewirings of the PPIN that preserved node centrality (Section 4.7). In this setting, network methods favored central genes (Fig 4B) even though highly central genes often had no association to disease breast cancer susceptibility (Fig 4B). On the other, network methods favor central genes, as they often connect high scoring nodes. This was specially the case of SigMod C. We found this bias especially in SigMod (S4-S6 Fig), which selected three highly central, unassociated genes in both the PPIN and in many of the random

rewirings: *COPS5*, *CUL3*, and *FN1*. As However, as we showed in Section 2.3, 2.3 and will show in 2.8, there is evidence in the literature of the contribution of the first two to breast cancer susceptibility. With regards to *FN1*, it encodes a fibronectin, a protein of the extracellular matrix involved in cell adhesion and migration. Overexpression of *FN1* has been observed in breast cancer [32], and it is negatively correlated anticorrelates with poor prognosis in other cancer types [33, 34].

By virtue of using a SNP subnetwork, SConES analyzes analyzed each SNP in their functional context. It therefore can select SNPs in genes without any associated interactor, as well as SNPs in Therefore, it could select SNPs located in genes not included in the PPIN and in non-coding regions or in non-interacting genes. In fact, due to linkage disequilibrium, SConES favors such genes, as selecting SNPs in an LD block which overlaps with a gene favors selecting the rest of the gene. This might explain why SConES produces We compared the solution of SConES in the GI network (using PPIN information), to the one using only positional information (GS network) and to the one using positional and gene annotations (GM network). Importantly, SConES produced similar results on the GS and GM networks, heavily affected by linkage disequilibrium (S5 Fig.). On the other hand, SConES penalizes selecting SNPs and not their neighbors. This makes it conservative regarding SNPs with many interactions, like those mapped to hub genes in the PPIN. For this reason, SConES GI did not select any protein coding gene, despite selecting similar regions as SConES GS. While the solutions on those two considerably overlap with SConES GI's, they contained additional gene-coding segments (Fig 4C). In fact, both SConES GS and SConES GM select GM selected chromosome regions related to breast cancer, like 3p24 (*SLC4A7/NEK10* [35]), 5p12 (*FGF10, MRPS30* [27]), 10q26 (*FGFR2*), and 16q12 (*TOX3*). On top of those SConES GS selects In addition to those, SConES GS selected region 8q24, also linked to breast cancer (*POU5F1B* [36]).

We hypothesize that the lack of results on the PPIN network of SConES GI and LEAN are due to the same cause: the absence of joint association of a gene and a majority of its neighbors. Although in the case of SConES other hyperparameters could lead to a more informative solution (e. g. a lower λ in Equation 5), it is unclear what the best strategy to find them is. In addition, due to the design of the iCOGS array, the genome of GENESIS participants has not been unbiasedly surveyed: some regions are fine-mapped — which might distort gene structure in GM and GI networks — while others are under studied — hindering the accuracy with which the GS network captures the genome structure.

2.7 Different parameters produce similarly-sized solutions

We explored methods' parameter space by running them under different combinations of parameters (Section 4.3.4). In agreement with their formulations (Section 4.3.3), larger values of specific parameters produced less stringent solutions (S5-S7 Fig A): for HotNet2 and heinz, this is the threshold above which genes receive a positive score; for dmGWAS, it is the d parameter, which controls how far neighbors could be added; for SigMod, it is nmax, which specifies the maximum size of the solution; and for LEAN, it is the P-value threshold to consider a solution significant. Two parameters had the opposite effect (the larger, the more stringent): SigMod's maxjump, which sets the threshold to consider an increment in λ "large enough"; and SConES' η , where higher values produce smaller solutions. However, two of the parameters did not have the expected effect: dmGWAS' r, which controls the minimum increment in the score required to add a gene; and SigMod's maxjump, which sets the threshold to consider an increment in λ "large enough". In both cases, the size of the solution was very similar across the different values. Despite the differences in size, the solutions' size was relatively robust to the choice of parameters (S5-S7 Fig B).

We computed the Pearson correlation between the different solutions as in Section 4.5.1 to study how the parameters affected which genes and SNPs were selected (S8 Fig.). This analysis showed that dmGWAS and SigMod were robust to two parameters: the parameter d determined dmGWAS' output more than r ; for SigMod, it was $nmax$ rather than $maxjump$.

SConES presented an interesting case in terms of feature selection: most of the explored combinations of parameters led to trivial solutions (they included either all the SNPs or none of them) (S8 Fig.). To explore a more meaningful parameter space, we selected the parameters in two rounds in our experiments. First, we explored the whole sample space. Then, we focused on a range of η and λ 1.5 orders of magnitude above and below the best parameters, respectively. This second parameter space was more diverse, which allowed to find more interesting solutions.

2.8 Adjusting for instability Building a stable consensus network preserves global network properties

Most of the network methods, including the consensus, were highly unstable (Fig 3B), raising questions about the reliability of the results. We built a new, stable consensus network solution using the genes selected most often across the 30 solutions obtained by running the 6 methods on 5 different splits of the data (Section ?? 4.5). Such a network is expected to capture the subnetworks more often found altered, and hence should be more resistant to noise. We used only genes selected in at least 7 of the solutions, which corresponded to 1% of all genes selected at least once. The resulting consensus was composed of 68 genes (Fig 5). This network shares most of the properties of the consensus: breast cancer susceptibility genes are overrepresented (P -value = 3×10^{-4}), as well as genes involved in mitochondrial translation and the attenuation phase (adjusted P -values 0.001 and 3×10^{-5} respectively); the selected genes are more central than average (P -value = 1.1×10^{-14}); and a considerable number of nodes (19) are isolated (Fig 5B and C).

However, although this new network exhibits similar global properties as the previous one, the lack of stability results in different genes being selected. In this case Despite these similarities, the consensus and the stable consensus included different genes. In the stable consensus network, the most central gene is *CUL3*, which is absent from the previous consensus network and has a low association score in both GENESIS and BCAC (P -values of 0.04 and 0.26, respectively). This gene is a component of Cullin-RING ubiquitin ligases. Encouragingly, it impacts the protein levels of multiple genes relevant for cancer progression [37], and its overexpression was also linked to increased sensitivity to carcinogens [38].

3 Discussion

In recent years, the ability of GWAS to unravel the mechanisms leading to complex diseases has been called into question [6]. On the one hand, the omnigenic model proposes that gene functions are interwoven with each other in a dense co-function network. The practical consequence is that larger and larger GWAS will lead to the discovery of discovering an uninformative wide-spread pleiotropy. On the other hand, discovery in GWAS is hindered by a Second, its conservative statistical framework hinders GWAS discovery. Network methods tackle elegantly address these two issues by using both the association score and an interaction network to take into consideration the biological context of each of the genes

and SNPs. Based on what could be considered diverse interpretations of the omnigenic model, several methods for network-guided biomarker discovery have been proposed in recent years. In this article we evaluated the relevance of six of them by examining six of these methods (Section 4.3.3) by applying them to the GENESIS study, a GWAS dataset on familial breast cancer (Section 4.1).

Most of the network methods produced a relevant subset of biomarkers, reCAPITULATING DmGWAS, Heinz, HotNet2, SConES, and SigMod all yielded compelling solutions, which include (but are not limited to) known breast cancer susceptibility genes (Section 2.2). In general, the selected genes and SNPs were more central than average, in accordance most other genes and SNPs, agreeing with the observation that disease genes are more relatively central [10]. However, very central nodes are also more likely to be connecting any given random pair of nodes, making them more likely to be selected by these network methods. Across this article we show that network methods (Section 2.6). However, we found support in the literature for the involvement of the selected highly central genes that were selected (*COPS5*, *CUL3* and *FN1*) could plausibly be involved, and *CUL3* in breast cancer susceptibility. Yet, further work is needed to characterize the impact of centrality on network methods' outputs (Sections 2.3, 2.6, and 2.8). Despite these similarities, the methods' solutions were notably different. At one end of the spectrum, SConES and heinz preferred small, highly associated solutions, providing a conveniently short list of biomarkers, at the expense of not shedding high scoring solutions, which were also small and hence did not shed much light on the etiology of the disease's etiology. On the other end, SigMod and dmGWAS gravitated towards larger, less associated solutions which provide dmGWAS, HotNet2 and SigMod gravitated towards lower scoring but larger solutions, which provided a wide overview of the biological context. While this deepens our understanding of the disease and provides breast cancer susceptibility and provided biological hypotheses, they require further analyses, which might deter unexperienced practitioners. HotNet2 balances both approaches at the expense of producing the largest solution: a constellation of many, highly associated, small subnetworks. Additionally required further analyses. For instance, we examined the centrality of the selected genes to understand how much that property was driving their selection (Section 2.6). However, all solutions share shared two drawbacks. First, they are were all equally bad at discriminating cases from controls. Yet However, the classification accuracy of network methods is was similar to that of a machine learning classifier trained on the entire genome (Section 2.5), which suggests that cases and controls are difficult to separate in this the GENESIS dataset. This might may be due to unaccounted for environmental factors, and limited statistical power, which reduces the ability to identify relevant SNPs. However, in any event, we do not expect to separate people who have or will develop cancer from others on the sole basis of their genomes, ignoring all environmental factors and chance events. Hence, network methods were preferable to the logistic regression classifier since they did "no worse" at classification while providing an interpretable solution. Second, all methods are were remarkably unstable, yielding different solutions for slightly different inputs. This might partly be have been caused by the instability of the P-values themselves in low statistical power settings [39]. Hence, heinz's conservative transformation of P-values, which favors favored only the most extreme ones, leads led to improved stability. Another source of instability might be have been the redundancy inherent to biological networks, a consequence of an evolutionary pressure to avoid single points of failure [40]. Hence, biological networks will often have multiple paths connecting two high-scoring nodes.

To overcome the limitations of the individual methods while exploiting their these limitations while exploiting the each method's strengths, we proposed combining them

into a consensus ~~subnetwork~~. We use a ~~solution~~. We used the straightforward strategy of including any node that was recovered by ~~multiple~~ at least two methods. We proposed two networks~~s~~ thus proposed two solutions (Sections 2.3 and 2.8): a consensus network that tackled ~~solution~~, which addressed the heterogeneity of the solutions~~in the full dataset~~, and a stable consensus network, that addressed ~~solution~~, which addressed the instability of the methods. They both synthesized the altered mechanisms: they both included the majority of the strongly associated smaller solutions and captured genes and broader mechanisms related to cancer, thus synthesizing the mechanisms altered in breast cancer cases. Thanks to their smaller size and ~~their~~ network structure, they provided compelling hypotheses on genes like *COPS5* and *CUL3*, which lack genome-wide association with the disease, ~~but who but~~ are related to cancer at the expression level and ~~whose neighborhood has consistent high association scores. Crucially, the consensus consistently interact with high scoring genes. Notably, while the consensus approach was as unstable as the tested methods, while individual network-guided methods, the stable consensus shared these properties while accounting for instability. This supports that instability might be caused by~~ network retained the ability to provide compelling hypotheses and had better stability. This supported that redundant but equivalent biological mechanisms, and hence validates mechanisms might cause instability and supported the conclusions obtained on the individual solutions~~and the consensus~~.

In this work, we have compared our results to significant genes and SNPs in the BCAC study [21]. Network methods showed modest precision but much higher recall at recovering BCAC hits (Section 2.4). While precision might be desirable when a subset of useful markers is required (for instance, for diagnosis), higher recall is desirable in exploratory settings. Nonetheless, BCAC was not an ideal ground truth. First, the studied populations are not entirely overlapping: BCAC focused on a pan-European cohort, while GENESIS targeted the French population. Second, the study designs differed: a high proportion of breast cancer cases investigated in BCAC were sporadic (not selected according to family history), while GENESIS was a homogeneous dataset not included in BCAC focused on the French high-risk population attending the family cancer clinics. Finally, and this is indeed the motivation for this study, GWAS are unlikely to identify all genes relevant for the disease: some might only show up in rare-variant studies; others might have too small effect sizes. Network methods account for this by including genes with low association scores but with relevant topological properties. Hence, network methods and GWAS, even when well-powered, are unlikely to capture exactly the same sets of genes. This might partly excuse the low precisions displayed in Section 2.4 and the low AUC displayed in Section 2.5.

As not all PPIN databases compiled the same interactions, the choice of the PPIN determines the final output. In this work, we used only interactions from HINT from high-throughput experiments (Section 4.3.2). This responded to concerns about adding interactions identified in targeted studies and falling into “rich getting richer” problems: since popular genes have a higher proportion of their interactions described [12, 41], they might bias discovery towards themselves by reducing the average shortest path length between two random nodes. On the other hand, Huang et al. [11] found that larger networks were more useful than smaller networks to identify disease genes. This would support using the largest networks in our experiments. However, when we compared the impact of using a larger PPIN containing interactions from both high-throughput experiments and the literature (Section 4.3.2), for most of the methods it did not change much the size or the stability of the solution, the classification accuracy, or the runtime (S8–S10 Fig). This supports using only interactions from high-throughput experiments, which produced similar solutions and

avoided falling into “circular reasonings”, where the best-known genes were artificially pushed into the solutions, as we observed in Section 2.6.

The strength of network-based analyses comes from leveraging prior knowledge to boost discovery. In consequence, they show their shortcomings ~~with respect to on~~ understudied genes, especially those not in the network. Out of the 32 767 genes ~~that we can map to which we mapped~~ the genotyped SNPs~~to~~, 60.7% (19 887) ~~are were~~ not in the protein-protein interaction network~~PPIN~~. The majority of those (14 660) are non-coding genes, mainly lncRNA, miRNA, and snRNA (~~S7 S9~~ Fig). ~~Yet Nevertheless~~, RNA genes like *CASC16* ~~are were~~ associated to breast cancer (Section 2.1), reminding us of the importance of using networks beyond coding genes. ~~In addition Besides~~, even protein-coding genes linked to breast cancer susceptibility [35], like *NEK10* (P-value 1.6×10^{-5} , ~~located near overlapping with~~ *SLC4A7*) or *POU5F1B*, were absent from the ~~network~~~~PPIN~~. However, on average protein-coding genes absent from the PPIN ~~are were~~ less associated with ~~this phenotype breast cancer susceptibility~~ (Wilcoxon rank-sum P-value = 2.79×10^{-8} , median P-values of 0.43 and 0.47). ~~As we are using interactions from high-throughput experiments, such difference cannot~~ This could not be due to well-known genes having more known interactions ~~because we only used interactions from high-throughput experiments~~. As disease genes tend to be more central [10], we hypothesize that it ~~is was~~ due to interactions between central genes being more likely. It is worth noting that network ~~approaches methods~~ that do not use PPIs, like SConES GS and GM, ~~did recover recovered~~ SNPs in *NEK10* and *CASC16*. Moreover, both SConES GM and GI recovered intergenic regions, which might contain key regulatory elements [42]~~and, yet, but~~ are excluded from gene-centric approaches. This shows the potential of SNP networks, in which SNPs are linked when there is evidence of co-function, to perform network-guided GWAS even in the absence of gene-level interactions. Lastly, all the methods are heavily affected by how SNPs are mapped to genes. ~~In Section 2.3 we highlight ambiguities that appear when genes overlap or are in linkage disequilibrium. In fact, the presented case is paradigmatic, since the genes are in the most gene-dense region of the genome [43]. Network methods are prone to selecting such genes when they are functionally related, and hence interconnected in the network, but might be more resilient to them when the overlapping genes are unrelated. Making use of more targeted mappings of SNPs to genes, and other strategies (e.g., eQTLs, SNPs associated to the expression of a gene), altogether with a stringent LD pruning, might address such problems. gene expression) might lead to different results.~~

~~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 about adding interactions identified in targeted studies and prone to a “rich getting richer” phenomenon: popular genes have a higher proportion of their interactions described [12, 41], and they might bias discovery by reducing the average shortest path length between two random nodes. On the other hand, Huang et al. [11] found that the best predictor of the performance of a network for disease gene discovery is the size of the network, which supports using the largest amount of interactions. When we compared the impact of using a larger network containing interactions from both high-throughput experiments and the literature (Section 4.3.2), we found that for most of the methods it did not greatly change the size or the stability of the solution, the classification accuracy, or the runtime (.). This supports using only interactions from high-throughput experiments, which produces apparently similar solutions and avoids falling into “circular reasonings”, where the best known genes are artificially pushed into the solutions.~~

A crucial step for the ~~gene-based~~gene-based methods is the computation of ~~the~~

~~gene score~~gene scores. In this work, we used VEGAS2 [44] 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 [45], ~~and other algorithms~~ [45–47]. Additionally, ~~other algorithms like PEGASUS~~ [45], SKAT [46] or COMBAT [47] might have more statistical power.

~~Another important decision is how to handle LD in a GWAS.~~

~~How to handle linkage disequilibrium (LD) is often a concern among GWAS practitioners. Often, the question is whether an LD-based pruning of the genotypes will improve the results.~~ VEGAS2 accounts for LD patterns, and hence an LD pruning step would not impact gene-based network methods, although it would speed up VEGAS2’s computation time. ~~In Section 2.3 we highlighted ambiguities that appear when genes overlap or are in LD. The presented case is paradigmatic since all three genes are in the HLA region, the most gene-dense region of the genome [43]. Network methods are prone to selecting such genes when they are functionally related, and hence interconnected in the PPIN. But the opposite case is also true: when genes are not functionally related (and hence disconnected in the PPIN), network methods might disregard them even if they have high association scores.~~ With regards to SConES, fewer SNPs would lead to simpler SNP networks and, possibly, shorter runtimes. However, ~~as mentioned in LD patterns also affect SConES’ in other ways, since its formulation penalizes selecting a SNP and not its neighbors, via a nonzero parameter η in Eq 5.~~ Due to LD, nearby SNPs’ P-values correlate; since positional information determines SNP networks, nearby SNPs are likely to be connected. Hence, SConES tends to select LD-blocks formed by low P-value SNPs. This might explain why SConES produced similar results on the GS and GM networks, heavily affected by LD (Section 2.6, ~~LD patterns seem paramount to SConES’ solutions, and an LD pruning step could potentially alter them~~). However, this same behavior raises the burden of proof required to select SNPs with many interactions, like those mapped to hub genes in the PPIN. For this reason, SConES GI did not select any protein coding gene. This could be caused by the absence of joint association of a gene and most of its neighbors, a hypothesis supported by LEAN’s lack of results. Yet, a different combination of parameters could lead to a more informative SConES’ solution (e.g., a lower λ in Eq 5), although it is unclear how to find it. In addition, due to the design of the iCOGS array (Section 4.1), the genome of GENESIS participants has not been unbiasedly surveyed: some regions are fine-mapped — which might distort gene structure in GM and GI networks — while others are understudied — hindering the accuracy with which the GS network captures the genome structure. A strong LD pruning might address such problems.

~~In order to produce the consensus networks~~ To produce the two consensus solutions, we faced ~~the different practical challenges due to the differences in~~ interfaces, preprocessing steps, and unexpected behaviors of the various methods. To ~~facilitate that other authors apply them~~ make it easier for others to apply these methods to new datasets and aggregate their solutions, we built six `nextflow` pipelines [48] with a consistent interface and, whenever possible, parallelized computation. They are available on GitHub: [hclimente/gwas-tools](#) (Section 4.9). Importantly, ~~those methods that had we compiled those methods with~~ a permissive license ~~were compiled~~ into a Docker image for easier use, ~~which is~~ available on Docker Hub [hclimente/gwas-tools](#).

4 Materials and methods

4.1 GENESIS dataset, preprocessing, and quality control

The GENE Sisters (GENESIS) study was designed to investigate investigated risk factors for familial breast cancer in the French population [20]. Index cases are were 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 / were unaffected colleagues or friends of the cases born around the year of birth of their corresponding case (± 3 years). We focused on the 2 577 samples of European ancestry, of which 1 279 are controls were controls, and 1 298 are were cases. The genotyping was performed using platform was the iCOGS array, a custom Illumina array designed to study genetic susceptibility of the genetic susceptibility to hormone-related cancers [49]. It contains contained 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 candidate functional variants in selected genes and pathways.

4.2 Preprocessing and quality control

We discarded SNPs with a minor allele frequency lower than 0.1%, those not in Hardy–Weinberg equilibrium in controls (P -value < 0.001), and those with genotyping data missing on more than 10% of the samples. A We also removed a subset of 20 duplicated SNPs in *FGFR2* were also removed. In addition, we removed. We excluded the samples with more than 10% missing genotypes. After controlling for relatedness, we excluded 17 additional samples were removed (6 for sample identity error, 6 controls related to other samples, 2 cases being related to an index case, and 3 additional controls having a high relatedness score). Lastly, based on study selection criteria, 11 other samples were removed (1 control having cancer, 4 index cases with no affected sister, 3 half-sisters, 1 sister with lobular carcinoma *in situ*, 1 with a *BRCA1* or *BRCA2* pathogenic variant detected in the family, 1 with unknown molecular diagnosis). The final dataset included 1 271 controls and 1 280 cases, genotyped over 197 083 SNPs.

We looked for population structure that could produce spurious associations. A principal component analysis revealed no visual differential population structure between cases and controls (S1 Fig). Independently, we did not find evidence of genomic inflation ($\lambda = 1.05$) either, further confirming the absence of confounding population structure.

4.2 High-score subnetwork search algorithms SNP- and gene-based GWAS

4.2.1 SNP and gene association

To measure association between a genotype and the phenotype

To measure the association between genotype and susceptibility to breast cancer, we performed a per-SNP 1 d.f. χ^2 allelic test using PLINK v1.90 [50]. To obtain significant SNPs, we performed a Bonferroni correction to keep the family-wise error rate below 5%. The threshold used was $\frac{0.05}{197083} = 2.54 \times 10^{-7}$.

Then, we used VEGAS2 [44] to compute the gene-level association score from the P -values of the SNPs mapped to them. Specifically, for each gene we only used the 10% of SNPs mapped to it with lowest P -values. We mapped More specifically, we mapped SNPs to genes through their genomic coordinates: all SNPs located within the

boundaries of a gene, ± 50 kb, were mapped to that gene. We computed VEGAS2 scores for each gene using only the 10% of SNPs with the lowest P-values among all those mapped to it. We used the 62 193 genes described in GENCODE 31 [51], although only 54 612 could be mapped to at least one SNP. Out of those, we focused exclusively on the 32 767 that had a gene symbol. Out of the 197 083 SNPs remaining after quality control, 164 037 were mapped to at least one of these genes.

We used such mapping to compare the outputs of methods that produce SNP lists to those that produce gene lists, and vice versa. For the former, we considered any gene that can be mapped to any of the selected SNPs as selected as well. For the latter, we considered all the SNPs that can be mapped to that gene as selected by the method. We also performed a Bonferroni correction to obtain significant genes; in this case, the threshold of significance was $\frac{0.05}{32767} = 1.53 \times 10^{-6}$.

4.3 Network methods

4.3.1 Mathematical notations

In this article, we use undirected, vertex-weighted networks, or graphs, $G = (V, E, w)$. $V = \{v_1, \dots, v_n\}$ refers to the vertices, with weights $w : V \rightarrow \mathbb{R}$. 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.

In addition to a weight, nodes have other properties. Nodes can be described by properties provided by the topology of the graph. In this article we focus on two of those: 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 two matrices that describe two different properties of a graph. Both matrices are square, and have as many rows and columns as nodes are in the network. The element (i, j) represents a relationship between v_i and v_j . The adjacency matrix W_G contains a 1 when the corresponding nodes are connected, and 0 otherwise; the its diagonal is zero. The degree matrix D_G is a diagonal matrix which contains the degree of the different nodes.

4.3.2 Networks

Gene network The statistical frameworks mathematical formulations of the different network methods are compatible with any type of network (protein interactions, gene coexpression, regulatory, etc.). Yet biological network (e.g., from protein interactions or gene co-expression). Here, we used protein-protein interaction networks (PPIN) for all of them network methods except SConES, as they PPINs are interpretable, well characterized, and they well-characterized, and the methods were designed to run efficiently on networks of their size. We built our PPIN from both binary and co-complex interactions stored in the HINT database (release April 2019) [41]. Unless otherwise specified, we used only interactions coming from high-throughput experiments, leaving out targeted studies that might bias the topology of the network. Out of the 146 722 interactions from high-throughput experiments that HINT stores, we were able to map 142 541 to a pair of gene symbols, involving 13 619 genes. 12 880 of those mapped to a genotyped SNP after quality

control, involving 127 604 interactions. The scoring function for the nodes changed from method to method (Section 4.3.3). 741
742

Additionally, we compared the results of the aforementioned obtained on this PPIN 743
with those obtained on another a PPIN built using interactions coming from both 744
high-throughput and targeted studies. In that case, out of the 179 332 interactions in 745
HINT, 173 797 mapped to a pair of gene symbols. Out of those, 13 735 mapped to a 746
genotyped SNP after quality control, involving 156 190 interactions. 747

SNP networks SConES [18] is was the only network method designed to handle 748
SNP networks. As in gene networks, two SNPs are were connected in a SNP network 749
when there is was evidence of shared functionality between two SNPs them. Azencott et 750
al. [18] proposed three ways of building these such networks: connecting the SNPs 751
consecutive in the genomic sequence (“GS network”); interconnecting all the SNPs 752
mapped to the same gene, on top of GS (“GM network”); and interconnecting all SNPs 753
mapped to two genes for which a protein-protein interaction exists, on top of GM (“GI 754
network”). We focused on the GI network using the PPIN described above, as it fits 755
fitted the scope of this work better, using the PPIN described above. However, at 756
different stages of this work we also used, we also compared GI to GS and GM for 757
comparison to understand how including the PPIN affects SConES’ output. For the 758
GM network, we used the mapping described in Section ??4.3.5. In all three the node 759
scores are the association scores of the individual SNPs with the phenotype (, we 760
scored the nodes using the 1 d.f. χ^2) statistic of association. The properties of these 761
three subnetworks are available in S1 Table. 762

4.3.3 Network methods High-score subnetwork search algorithms 763

Genes that contribute to the same function are nearby in the PPIN, and can be 764
topologically related to each other in diverse ways (densely interconnected modules, 765
nodes around a hub, a path, etc.). But this is not the only aspect to model. Several 766
aspects have to be considered when developing a network method: how to score the 767
nodes, whether the affected mechanisms form a single connected component or several, 768
how to frame the problem in a computationally efficient fashion, which network to use, 769
etc. Unsurprisingly, multiple solutions have been proposed. We examined six of them: 770
five that explore the PPIN, and one which explores SNP networks. We selected 771
methods that were open source, open-source methods that had an implementation 772
available, and an and accessible documentation. Their main differences are 773
summarized. We summarize their main differences in Table 2. We scored both SNPs 774
and genes with the P-values (or transformations) computed in Section 4.2.1. 775

dmGWAS dmGWAS seeks the subgraph with the highest local density in low 776
P-values [15]. To that end, it searches candidate subnetwork solutions using a 777
greedy, “seed and extend”, heuristic: 778

1. Select a seed node i and form the subnetwork $S_i = \{i\}$. 779
2. Compute Stouffer’s Z-score Z_m for S_i as 780

$$Z_m = \frac{1}{\sqrt{k}} \sum_{j \in S_i} z_j, \quad (1)$$

where k is the number of genes in S_i , z_j is the Z score of gene j , computed 781
as $\phi^{-1}(1 - P\text{-value}_j)$, and ϕ^{-1} is the inverse normal distribution function. 782

3. Identify neighboring nodes of S_i , i.e., nodes at distance $\leq d$. 783

Table 2. Summary of the differences between the network methods.

Method	Field	Nodes	Exhaustive	Solution	Comp.	Input	Scoring	Ref.
dmGWAS	GWAS	Genes	No	-	1	Summary	$-\log_{10}(P)$	[15]
heinz	Omics	Genes	Yes	-	1	Summary	BUM	[16]
HotNet2	Omics	Genes	Yes	Module	≥ 1	Summary	Local FDR	[17]
LEAN	Omics	Genes	Yes	Star	≥ 1	Summary	$-\log_{10}(P)$	[14]
SConES	GWAS	SNPs	Yes	Module	≥ 1	Genotypes	1 d.f. χ^2	[18]
SigMod	GWAS	Genes	Yes	Module	≥ 1	Summary	$\Phi^{-1}(1 - P)$	[19]

Field: field in which the algorithm was developed. **Nodes:** the type of nodes in the network, either genes (PPIN) or SNPs.

Exhaustive: whether the method explores all the possible solutions given the selected hyperparameters are explored parameters.

Solution: additional properties are enforced on the solution subnetwork, other than containing high scoring, connected nodes. **Comp.:** number of connected components in the solution. **Input:** genotype data or GWAS summary statistics. **Scoring:** how SNP/gene P-values are transformed into node scores. In the case of heinz, BUM stands for beta-uniform model, used to transform the P-values; for SigMod, Φ^{-1} represents the inverse of the cumulative distribution function of the standard Normal distribution. **Ref.:** original publication featuring the algorithm.

4. Add the neighboring nodes whose inclusion increases the $Z_{m+1}Z_{m+1}$ by more than a threshold $Z_m \times (1 + r\bar{x})$.
5. Repeat 2-4 until no further enlargement is possible.
6. Add S_i to the list of subnetworks to return. Its Normalized its Z-score is normalized as

$$Z_N = \frac{Z_m - \text{mean}(Z_m(\pi))}{\text{SD}(Z_m(\pi))}, \quad (2)$$

where $Z_m(\pi)$ represents a vector containing 100 000 random subsets of the same number of genes.

DmGWAS carries out this process on every gene in the network PPIN. We used the implementation of dmGWAS in the dmGWAS 3.0 R package [52]. We Unless otherwise specified, we used the suggested hyperparameters $d = 2$ and $r = 0.1$ parameters $d = 2$ and $r = 0.1$. We used the function simpleChoose simpleChoose to select the solution subnetwork, which aggregates the top 1% subnetworks.

heinz The goal of heinz is to identify the highest-scored connected subnetwork [16]. The authors propose proposed a transformation of the genes' P-value into a score that is negative under no weak association with the phenotype, and positive when there is under a strong one. This transformation is achieved by modelling modeling the distribution of P-values by a beta-uniform model (BUM) parameterized by the desired false discovery rate (FDR). Thus formulated, the problem is NP-complete, and hence solving it would require a prohibitively long computational time. To solve it efficiently, it is re-cast as the Prize-Collecting Steiner Tree Problem (PCST), which seeks to select the connected subnetwork S that maximizes the profit $p(S)$, defined as:

$$p(S) = \sum_{v \in V_S} p(v) - \sum_{e \in E_S} c(e). \quad (3)$$

were $p(v) = w(v) - w'$ is the profit of adding a node, $c(e) = w'$ is the cost of adding an edge, and $w' = \min_{v \in V_G} w(v)$ is the smallest node weight of G . All

three are positive quantities. Heinz implements the algorithm from Ljubić et al. [53] which, in practice, is often fast and optimal, although neither is guaranteed. We used BioNet’s implementation of heinz [54, 55].

HotNet2 HotNet2 was developed to find connected subgraphs of genes frequently mutated in cancer [17]. To that end, it considers both the local topology of the network and the scores of the nodes. The former is captured by an PPIN and the nodes’ scores. An insulated heat diffusion process captures the former: at initialization, 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—a stationary state is reached, in which the temperature of the nodes does not change anymore, and results in a diffusion matrix F . F is used to compute the similarity matrix E that models exchanged heat as

$$E = F \text{diag}(w(V)), \quad (4)$$

where $\text{diag}(w(V))$ is a diagonal matrix with the node scores in its diagonal. For any two nodes i and j , E_{ij} models the amount of heat that diffuses from node j to node i , which. Hence, E_{ij} can be interpreted as a (non-symmetric) similarity between those two nodes. To obtain densely connected subnetworksolutions, HotNet2 prunes E , only preserving edges such that $w(E) > \delta$. Lastly, HotNet2 evaluates the statistical significance of the subnetworks solutions by comparing their size to the size of networks PPINs obtained by permuting the node scores. We assigned the initial node scores as in Nakka et al. [45], assigning a score of giving a 0 for to the genes with low probability of being associated to a VEGAS2 P-values of association with the disease, and $-\log_{10}(\text{P-value})$ to those likely to be. In this the GENESIS dataset, the threshold separating both was a P-value of 0.125, which was we obtained using a local FDR approach [56]. HotNet2 has two parameters: the restart probability β , and the threshold heat δ . Both parameters are set automatically by the algorithm, which is robust to their values [17]. HotNet2 is implemented in Python [57].

LEAN LEAN searches altered “star” subnetworks, that is, subnetworks composed by of one central node and all its interactors [14]. By imposing this restriction, LEAN is able to can exhaustively test all such subnetworks (one per node). For a particular star subnetwork of size m , LEAN performs three steps:

1. Rank the P-values corresponding to of the involved nodes are ranked as $p_1 \leq \dots \leq p_m$. Then,
2. Conduct 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—
3. Transform this score into a P-value through an empirical distribution obtained via a subsampling scheme, where gene sets of the same size are selected randomly, and their score computed. Lastly,

We adjust these P-values are corrected for multiple testing through a Benjamini-Hochberg correction. We used the implementation of LEAN from the LEANR R package [58].

SConES SConES searches the minimal, modular, and maximally associated subnetwork in a SNP graph [18]. Specifically, it solves the problem

$$\arg \max_{S \subseteq G} \underbrace{\sum_{v \in V_S} w(v)}_{\text{association}} - \lambda \underbrace{\sum_{v \in V_S} \sum_{u \notin V_S} W_{vu}}_{\text{connectivity}} - \underbrace{\eta |V_S|}_{\text{sparsity}}, \quad (5)$$

where λ and η are **hyperparameters** that control the sparsity and the connectivity of the model. The connectivity term penalizes disconnected solutions, with many edges between **nodes that are selected and nodes that are not selected and unselected nodes**. Given a λ and an η , **the aforementioned problem Eq 5** has a unique solution \bar{v} that SConES finds using a graph min-cut procedure. As in Azencott et al. [18], we selected λ and η by cross-validation, choosing the values that produce the most stable solution across folds. 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. We used the version on SConES implemented in the R package `martini` [59].

SigMod SigMod **aims at identifying** the highest-scoring, most densely connected **gene** subnetwork [19]. It addresses an optimization problem similar to that of SConES (Equation Eq 5), but **the connectivity term encourages connected solutions by favoring solutions where many edges connect two selected nodes, rather than penalizing disconnected ones** **with a different connectivity term that favors solutions containing many edges**:

$$\arg \max_{S \subseteq G} \underbrace{\sum_{v \in V_S} w(v)}_{\text{association}} + \lambda \underbrace{\sum_{v \in V_S} \sum_{u \in V_S} W_{vu}}_{\text{connectivity}} - \underbrace{\eta |V_S|}_{\text{sparsity}}. \quad (6)$$

As **for** 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, it was designed for PPINs**. Second, it favors **subnetworks** containing many edges **between the selected nodes**. SConES, instead, penalizes connections between **the selected and unselected nodes**. Third, it explores the grid of **hyperparameters** differently, and processes their respective solutions. Specifically, for the range of $\lambda = \lambda_{\min}, \dots, \lambda_{\max}$ for the same η , it prioritizes the solution with the largest change in size from λ_n to λ_{n+1} . **Additionally, that change needs to be larger than a user-specified threshold** **maxjump**. Such a large change implies that the network is densely interconnected. This results in one candidate solution for each η , 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 Finally, SigMod chooses the candidate solution with the highest standardized score** **is the chosen solution, and that is not larger than a user-specified threshold** **(nmax)**. We used the default parameters **maxjump** = 10 and **nmax** = 300. SigMod is implemented in an R package [60].

Consensus We built a consensus **network** solution by retaining the **nodes that were genes** selected by at least two of the six methods (using SConES GI for SConES). **It includes any edge between the selected genes in the PPIN**.

We performed all the computations in the cluster described in Section 4.8.

4.3.4 Parameter space

We used the network methods with the parameters recommended by their authors, or with the default parameters in their absence. Additionally, we explored the parameter space of the different methods to study how they alter the output.

dmGWAS We tested multiple values for r (0.0001, 0.001, 0.01, 0.05, 0.1, 0.25, 0.5, and 1) and d (1, 2, and 3).

heinz We tested multiple FDR thresholds (0.05, 0.1, 0.15, 0.2, 0.25, 0.3, 0.35, 0.4, 0.45, 0.5, 0.55, 0.6, 0.65, 0.7, 0.75, 0.8, 0.85, 0.9, 0.95, 1).

HotNet2 We tested different thresholds to decide which genes would receive a score of 0 and which ones a score of $-\log_{10}(P\text{-value})$: 0.001, 0.01, 0.05, 0.125, 0.25, and 0.5.

LEAN We used the following significance cutoffs for LEAN's P-values (0.05, 0.1, 0.15, 0.2, 0.25, 0.3, 0.35, 0.4, 0.45, 0.5, 0.55, 0.6, 0.65, 0.7, 0.75, 0.8, 0.85, 0.9, 0.95, and 1).

SConES We used the values of λ and η that **martini** explores by default (35.54, 5.40, 0.82, 0.12, 0.02, 0.01, 4.39e-4, 6.68e-5, 1.02e-5, and 1.55e-6 in both cases).

SigMod We tested multiple values for the parameters `nmax` (10, 50, 100, 300, 700, 1000, and 10 000) and `maxjump` (5, 10, 20, 30, and 50).

4.3.5 Comparing SNP-methods to gene-methods and vice versa

In multiple steps of this article, we compared the outcome of a method that works on genes with the outcome of one that works on SNPs. For this purpose, we used the SNP-gene correspondence described in Section 4.2.1. To convert a list of SNPs into a list of genes, we included all the genes mapped to any of those SNPs. Conversely, to convert a list of genes into a list of SNPs, we included all the SNPs mapped to any of those genes.

4.4 Pathway enrichment analysis

We searched for pathways enriched in the gene **subnetworks** **solutions** produced by the above methods. We conducted ~~an a~~ hypergeometric test on pathways from Reactome [61] using the **R-package ReactomePA** function `enrichPathway` from the **ReactomePA R package** [62]. The universe of genes included any gene that we could map to a SNP in the iCOGS array (Section ?? 4.2.1). We adjusted the P-values for multiple testing as in Benjamini and Hochberg [63] (BH). ~~Pathways with an a~~: pathways with a BH adjusted P-value < 0.05 were deemed significant.

~~As the significant pathways are often overlapping and redundant, the results were manually curated afterwards.~~

4.5 Evaluation Benchmark of methods

We evaluated multiple properties (**described below**) of the different methods (**described in Sections 4.5.1 and 4.5.2**) through a 5-fold subsampling setting. We applied each method to 5 random subsets of the original **GENESIS** dataset containing 80% of the samples (*train set*). When pertinent, we evaluated the solution on the remaining 20% (*test set*). We used the 5 repetitions to estimate the average and the standard deviation of the different measures. **Every method and repetition ran in the same computational settings** (Section 4.8).

4.5.1 Properties of the solution

We compared the runtime, the number of selected genes/SNPs/features (genes or SNPs), and the stability (sensitivity of the result to small changes in the input, here, using different train sets to the choice of train set) of the different network methods. The stability was quantified Nogueira and Brown [64] proposed quantifying a method's stability using the Pearson correlation to measure the overlap between different runs as suggested by Nogueira and Brown [64] between the genes selected on different subsamples. This correlation was calculated between vectors with the length of the total number of features, containing a 0 at position i if feature i was not selected and a 1 if it was.

4.5.2 Classification accuracy of selected SNPs

A desirable solution offers good predictive power on unseen test the unseen test samples. We evaluated the predicting power of the SNPs selected by the different methods through the performance of an L1-penalized logistic regression classifier, a machine learning algorithm that searches which searches for a small subset of SNPs which provide that provides good classification accuracy. We trained the classifier exclusively on those selected SNPs to predict the at predicting the outcome (case/control). The L1 penalty helps to account for linkage disequilibrium by reducing the number of SNPs included in the model (active set), while improving the generalization of the classifier. This. The active set was a plausible, more sparse solution with comparable predictive power to the original solution. The L1 penalty was set by cross-validation, choosing the value that minimized misclassification error.

We applied each network method to each train set, and train set. Then, we trained the classifier on it as well using only on the same train set using only the selected SNPs. When the method retrieved a list of genes(all of them except SConES), we considered as selected all the SNPs mapped to any of those genes. Then, we proceeded as explained in Section 4.3.5. Lastly we evaluated the sensitivity and the specificity on the test set. The active set gave an estimate of a plausible, more sparse solution with a comparable predictive power to the original solution of the classifier on the test set. To obtain a baseline, we also trained the classifier on all the SNPs. We do of the train set.

We did not expect a linear model on selected SNPs to be able to separate cases from controls well. Indeed, the lifetime cumulative incidence of breast cancer among women with a family history of breast or ovarian cancer, and no *BRCA1/2* mutations, is only 3.9 times more than in the general population [65]. However, classification accuracy may be one additional informative criterion on which to evaluate solutions.

4.5.3 Comparison to state-of-the-art

4.6 Comparison to state-of-the-art

An alternative way to evaluate the results is comparing our results methods is by comparing their solutions to an external dataset. For that purpose, we recovered a list of used the 153 genes associated to familial breast cancer from on DisGeNET [66]. Across this article, we refer to these genes as *breast cancer susceptibility genes*.

Additionally, we used the summary statistics from the Breast Cancer Association Consortium (BCAC), a meta-analysis of case-control studies conducted in multiple countrieswhich. BCAC included 13 250 642–641 SNPs genotyped or imputed on 228 951 women of European ancestry, mostly from the general population [21]. Hence, a high proportion of breast cancer cases investigated in BCAC are sporadic (not selected according to family history), while GENESIS is a homogeneous dataset not included in BCAC and which focus on the French high-risk population attending the family

eancer clinics. Despite these differences, we expect some degree of shared genetic architecture, especially at the gene level. For that purpose, we searched associated genes. Through imputation, BCAC includes more SNPs than the iCOGS array used for GENESIS (Section 4.1). However, in all the comparisons in this paper we focused on the SNPs that passed quality control in GENESIS. Hence, we used the same Bonferroni threshold as in Section ???. We provided 4.2.1 to determine the significant SNPs in BCAC. We also computed gene-scores in the BCAC data using VEGAS2 with both, as in Section 4.1. In this case, we did use the summary statistics of all available SNPs, 13 250 641 available SNPs and the genotypes from European samples from the 1000 Genomes Project [67] to compute the LD patterns. Since these genotypes did not include chromosome X, we excluded it from this analysis. All comparisons included only the genes in common between GENESIS and BCAC, so we used a different Bonferroni threshold (1.66×10^{-6}) to call gene significance.

4.7 Code availability Network rewirings

Rewiring the PPIN while preserving the number of edges of each gene allowed to study the impact of the topology on the output of network methods. Indeed, the edges lose their biological meaning while the topology of the network is conserved. We produced 100 such rewirings by randomly swapping edges in the PPIN. We still scored the genes as described in Section 4.3.3. We only applied only four methods on the rewirings: heinz, dmGWAS, LEAN, and SigMod. We excluded HotNet2 and SConES since they took notably longer to run.

4.8 Computational resources

We ran all the computations on a Slurm cluster, running Ubuntu 16.04.2 on the nodes. The CPU models on the nodes were Intel Xeon CPU E5-2450 v2 at 2.50GHz and Intel Xeon E5-2440 at 2.40GHz. The nodes running heinz and HotNet2 had 20GB of memory; the ones running dmGWAS, LEAN, SConES, and SigMod, 60GB. For the benchmark (Section 4.5), we ran each of the methods on the same Ubuntu 16.04.2 node, with a CPU Intel Xeon E5-2450 v2 at 2.50GHz, and 60GB of memory.

4.9 Code and data availability

We developed computational pipelines for several steps of GWAS analyses, such as physically mapping SNPs to genes, computing gene scores, and performing running six different network analyses. For each of those processes, we methods. We created a pipeline with a clear interface that should work on any GWAS dataset for each of those processes. They are compiled in <https://github.com/hclimente/gwas-tools>. Although the GENESIS data is not public, the code to apply the pipelines to this dataThe code that applies them to GENESIS, as well as the code that reproduces all the analyses in this article are available at <https://github.com/hclimente/genewa>. We deposited all the produced gene subnetworks solutions on NDEEx (<http://www.ndexbio.org>), under the UUID e9b0e22a-e9b0-11e9-bb65-0ac135e8bacf.

Summary statistics for SNPs and genes are available at <https://github.com/hclimente/genewa>. We cannot share genotype data publicly for confidentiality reasons, but are available from GENESIS. Interested researchers can contact nadine.andrieu(at)curie.fr.

5 Supporting information

S1 Table. Summary statistics on the results of SConES on the three SNP-SNP interaction networks. Summary statistics on the results of SConES on the three SNP networks (Section 4.3.2).	1029
The first row within each block contains the summary statistics on the whole network.	1030
	1031
	1032
S2 Table. 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. Pathway enrichment analyses of the genes in SigMod's solution.	1033
	1034
	1035
	1036
S3 Table. Pathway enrichment analyses of the genes in dmGWAS' solution.	1037
	1038
S4 Table. Pathway enrichment analyses of the genes in HotNet2's solution.	1039
	1040
S5 Table. Pathway enrichment analyses of the genes in the consensus' solution.	1041
	1042
S1 Fig. GENESIS shows no differential population structure between cases and controls. 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.	1043
	1044
	1045
	1046
S2 Fig. Association in GENESIS. The red line represents the Bonferroni threshold. Association in GENESIS. The red lines represent the Bonferroni thresholds. (A) SNP association, measured from the outcome of a 1 d.f. χ^2 allelic test (Section 4.2.1). Significant SNPs that are within a coding gene, or within 50 kilobases of its boundaries, are annotated. The Bonferroni threshold is 2.54×10^{-7} . (B) Gene association, measured by P-value of VEGAS2 [44] using the 10% of SNPs with the lowest P-values (Section 4.2.1). The Bonferroni threshold is 1.53×10^{-6} . (C) SNP association as in panel (A). The SNPs in black are selected by a L1-penalized logistic regression (Section 4.5.2, $\lambda = 0.03$).	1047
	1048
	1049
	1050
	1051
	1052
	1053
	1054
	1055
S3 Fig. Relationship between the \log_{10} of the betweenness centrality and the $-\log_{10}$ of the VEGAS2 P-value of the genes in the consensus solution. The blue line represents a fitted generalized linear model.	1056
	1057
	1058
S4 Fig. Additional benchmarks of the network methods. (A) Precision and recall of the evaluated methods with respect to Bonferroni-significant SNPs/genes in BCAC. For reference, we added a gray line with a slope of 1. This panel is identical to Fig 2. (B) Sensitivity and specificity on the test set of the L1-penalized logistic regression trained on the features selected by each of the methods. The performance of the classifier trained on all SNPs is also displayed. Points are the average over the 5 runs; the error bars represent the standard error of the mean.	1059
	1060
	1061
	1062
	1063
	1064
	1065
S5 Fig. Pearson correlation between the different solution subnetworks. Pearson correlation between the different solutions. (A) Correlation between selected SNPs. (B) Correlation between selected genes. In general, the solutions display a very low overlap.	1066
	1067
	1068
	1069

S4-S6 Fig. Consensus subnetwork on GENESIS (Section 4.3.3). Number of times a gene was selected by either dmGWAS, heinz, LEAN, or SigMod in 100 rewirings of the PPIN (A) Section 4.7. Each node is represented by a pie chart, which shows the methods that selected it. We labeled (and enlarged) the two most central genes (*COPS5* and *OFD1*) and those genes that are known breast cancer susceptibility genes and/or significantly associated with breast cancer susceptibility in the BCAC dataset. The latter ones are also colored in pink. This panel is identical to Fig 2. (B) Same network, but every gene name is indicated. 4B, split by method.

S5-S7 Fig. Genes on the consensus network. Breast cancer susceptibility genes are colored in pink; the rest are colored in grey. **Size of the solutions obtained under different parameters.** (A) Number of methods selecting every gene in the subnetwork. (B) VEGAS2 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. **Size of the solution produced by different parameter values**, expressed as a percentage of the maximum solution size for the method, or the highest tested value for the parameter, respectively. The size of the solution is the median among all the solution sizes for the same parameter. (D) Relationship between the \log_{10} of the betweenness centrality and the \log_{10} of the VEGAS2 P-value of the genes in the consensus network. The blue line represents a fitted generalized linear model.

S6 Fig. Stable consensus subnetwork on GENESIS. Boxplot of the solution sizes of the methods under the explored parameters (Section 2.8, 4.3.4).

S8 Fig. (A) Pearson correlation between the solutions obtained under different parameters, computed as in Section 4.5.1. Each node is represented by a pie chart, which shows the methods that selected it. We labeled (and enlarged) the two most central genes (*CUL3*) and those genes that are known breast cancer susceptibility genes and Grey tiles represent the cases where we could not compute the Pearson correlation because the two vectors were either all ones (all genes or significantly associated with breast cancer susceptibility in the BCAC dataset. The latter ones are also colored in pink. This panel is identical to Fig 5. (B) Same network, but every gene name is indicated. SNPs were selected) or zeros (no genes/SNPs were selected).

S7-S9 Fig. Biotypes of genes from the annotation that are not present in the HINT protein-protein interaction network.

Biotypes of genes from the annotation that are not present in the HINT PPIN.

S8-S10 Fig. Comparison of benchmark on high-throughput (HT) interactions to benchmark on both high-throughput and literature curated interactions (HT+LC). **Comparison of the benchmark on high-throughput (HT) interactions to the benchmark on both high-throughput and literature curated interactions (HT+LC).** Grey lines represent no change in the statistic between the benchmarks (1 for ratios mean(HT) / mean(HT + LC), 0 for differences mean(HT) - mean(HT + LC)). (A) Ratios of the selected features between both benchmarks and of the active set (Section 4.5.2). (B) Shifts in sensitivity and specificity. (C) Shift in

Pearson correlation between benchmarks. (D) Ratio between the runtimes of the benchmarks. For gene network-based gene-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 119 980 seconds (1 day and 9.33 hours) that VEGAS2 took on average to compute the gene scores from SNP summary statistics. In general, adding additional interactions slightly improves improved the stability of the solution, but increases. However, it increased the solution size, has and the required runtime, and had mixed effects on the sensitivity and specificity, and impacts negatively the required runtime of the algorithms.

S9-S11 Fig. Overview of the solutions produced by the SConES on the GS and GM networks (Section 4.3.2) on the GENESIS dataset. Overview of the solutions produced by the SConES on the GS and GM networks (Section 4.3.2) on the GENESIS dataset. (A) Manhattan plots of SNPs (Section 4.2.1); in black, the method's solution. The red line indicates the Bonferroni threshold (2.54×10^{-7}) is indicated by a red line. (B) Precision and recall of the evaluated methods with respect to Bonferroni-significant SNPs (SConES) or genes (other methods) in BCAC. For reference, we added a gray line with a slope of 1. (C) Solution networks.

Acknowledgments

We wish to thank Om Kulkarni for helpful discussion on gene-based GWAS and PPIN databases, and the genetic epidemiology platform (the PIGE, Plateforme d'Investigation en Génétique et Epidemiologie; S. Eon-Marchais, M. Marcou, D. Le Gal, L. Toulemonde, J. Beauvallet, N. Mebirouk, E. Cavaciuti), the biological resource eentre and center (S. Mazoyer, F. Damiola, L. Barjhoux, C. Verny-Pierre, V. Sornin). We wish to pay tribute to Olga M. Sinilnikova, one of the initiators and principal investigators of GENESIS, and who died prematurely on June 30, 2014.

We thank all the GENESIS collaborating cancer clinics clinics (Clinique Sainte Catherine, Avignon; H. Dreyfus; Hôpital Saint Jacques, Besançon; M-A. Collonge-Rame; Institut Bergonié, Bordeaux; M. Longy, A. Floquet, E. Barouk-Simonet; CHU, Brest; S. Audebert; Centre François Baclesse, Caen; P. Berthet; Hôpital Dieu, Chambéry; S. Fert-Ferrer; Centre Jean Perrin, Clermont-Ferrand; Y-J. Bignon; Hôpital Pasteur, Colmar; J-M. Limacher; Hôpital d'Enfants CHU – Centre Georges François Leclerc, Dijon; L. Faivre-Olivier; CHU, Fort de France; O. Bera; CHU Albert Michallon, Grenoble; D. Leroux; Hôpital Flaubert, Le Havre; V. Layet; Centre Oscar Lambret, Lille; P. Vennin, C. Adenis; Hôpital Jeanne de Flandre, Lille; S. Lejeune-Dumoulin, S. Manouvrier-Hanu; CHRU Dupuytren, Limoges; L. Venat-Bouvet; Centre Léon Bérard, Lyon; C. Lasset, V. Bonadona; Hôpital Edouard Herriot, Lyon; S. Giraud; Institut Paoli-Calmettes, Marseille; F. Eisinger, L. Huiart; Centre Val d'Aurelle – Paul Lamarque, Montpellier; I. Coupier; CHU Arnaud de Villeneuve, Montpellier; I. Coupier, P. Pujol; Centre René Gauducheau, Nantes; C. Delnatte; Centre Catherine de Sienne, Nantes; A. Lortholary; Centre Antoine Lacassagne, Nice; M. Frénay, V. Mari; Hôpital Caremeau, Nîmes; J. Chiesa; Réseau Oncogénétique Poitou Charente, Niort; P. Gestal; Institut Curie, Paris; D. Stoppa-Lyonnet, M. Gauthier-Villars, B. Buecher, A. de Pauw, C. Abadie, M. Belotti; Hôpital Saint-Louis, Paris; O. Cohen-Haguenauer; Centre Viggo-Petersen, Paris; F. Cornélis; Hôpital Tenon, Paris; A. Fajac; GH Pitié Salpêtrière et Hôpital Beaujon, Paris; C. Colas, F. Soubrier, P. Hammel, A. Fajac; Institut Jean Godinot, Reims; C. Penet, T. D. Nguyen; Polyclinique Courlancy, Reims; L. Demange*, C. Penet; Centre Eugène Marquis, Rennes; C. Dugast*, Centre Henri Becquerel, Rouen; A. Chevrier, T. Frebourg, J. Tinat, I. Tennevret, A. Rossi; Hôpital René

Huguenin/Institut Curie, Saint Cloud; C. Noguès, L. Demange*, E. Mouret-Fourme; CHU, Saint-Etienne; F. Prieur; Centre Paul Strauss, Strasbourg; J-P. Fricker, H. Schuster; Hôpital Civil, Strasbourg; O. Caron, C. Maugard; Institut Claudius Regaud, Toulouse; L. Gladieff, V. Feillel; Hôpital Bretonneau, Tours; I. Mortemousque; Centre Alexis Vautrin, Vandoeuvre-les-Nancy; E. Luporsi; Hôpital de Bravois, Vandoeuvre-les-Nancy; P. Jonveaux; Gustave Roussy, Villejuif; A. Chompret*, O. Caron). *Deceased prematurely

1166
1167
1168
1169
1170
1171
1172

Author contributions

1173

Conceptualization Héctor Climente-González, Christine Lonjou, Chloé-Agathe Azencott.

1174
1175

Data curation Christine Lonjou, GENESIS Study collaborators.

1176

Formal Analysis Héctor Climente-González, Christine Lonjou.

1177

Funding acquisition Dominique Stoppa-Lyonnet, Nadine Andrieu, Chloé-Agathe Azencott.

1178
1179

Investigation Héctor Climente-González, Christine Lonjou.

1180

Methodology Héctor Climente-González, Christine Lonjou, Chloé-Agathe Azencott.

1181

Project administration Chloé-Agathe Azencott.

1182

Resources GENESIS Study collaborators, Dominique Stoppa-Lyonnet, Nadine Andrieu.

1183

1184

Software Héctor Climente-González, Christine Lonjou.

1185

Supervision Christine Lonjou, Fabienne Lesueur, Nadine Andrieu, Chloé-Agathe Azencott.

1186

1187

Validation Christine Lonjou, Fabienne Lesueur.

1188

Visualization Héctor Climente-González.

1189

Writing – original draft Héctor Climente-González.

1190

Writing – review & editing Héctor Climente-González, Christine Lonjou, Fabienne Lesueur, Nadine Andrieu, Chloé-Agathe Azencott.

1191

1192

References

1. Bush WS, Moore JH. Chapter 11: Genome-Wide Association Studies. PLoS Computational Biology. 2012;8(12):e1002822. doi:10.1371/journal.pcbi.1002822.
2. Buniello A, MacArthur JA, Cerezo M, Harris LW, Hayhurst J, Malangone C, et al. The NHGRI-EBI GWAS Catalog of published genome-wide association studies, targeted arrays and summary statistics 2019. Nucleic Acids Research. 2019;47(D1):D1005–D1012. doi:10.1093/nar/gky1120.
3. Visscher PM, Wray NR, Zhang Q, Sklar P, McCarthy MI, Brown MA, et al. 10 Years of GWAS Discovery: Biology, Function, and Translation. The American Journal of Human Genetics. 2017;101(1):5–22. doi:10.1016/j.ajhg.2017.06.005.

4. Wang MH, Cordell HJ, Van Steen K. Statistical methods for genome-wide association studies. *Seminars in Cancer Biology*. 2018;doi:10.1016/j.semancer.2018.04.008.
5. Barton NH, Etheridge AM, Véber A. The infinitesimal model: Definition, derivation, and implications. *Theoretical Population Biology*. 2017;118:50–73. doi:10.1016/j.tpb.2017.06.001.
6. Boyle EA, Li YI, Pritchard JK. An Expanded View of Complex Traits: From Polygenic to Omnipigenic. *Cell*. 2017;169(7):1177–1186. doi:10.1016/j.cell.2017.05.038.
7. Furlong LI. Human diseases through the lens of network biology. *Trends in Genetics*. 2013;29(3):150–159. doi:10.1016/j.tig.2012.11.004.
8. Leiserson MD, Eldridge, Jonathan V, Ramachandran, Sohini, Raphael, Benjamin J. Network analysis of GWAS data. *Current Opinion in Genetics & Development*. 2013;23(6):602–610.
9. Barabási AL, Gulbahce N, Loscalzo J. Network medicine: a network-based approach to human disease. *Nature Reviews Genetics*. 2011;12(1):56–68. doi:10.1038/nrg2918.
10. Piñero J, Berenstein A, Gonzalez-Perez A, Chernomoretz A, Furlong LI. Uncovering disease mechanisms through network biology in the era of Next Generation Sequencing. *Scientific Reports*. 2016;6(1):24570. doi:10.1038/srep24570.
11. Huang JK, Carlin DE, Yu MK, Zhang W, Kreisberg JF, Tamayo P, et al. Systematic Evaluation of Molecular Networks for Discovery of Disease Genes. *Cell Systems*. 2018;6(4):484–495.e5. doi:10.1016/j.cels.2018.03.001.
12. Cai JJ, Borenstein E, Petrov DA. Broker Genes in Human Disease. *Genome Biology and Evolution*. 2010;2:815–825. doi:10.1093/gbe/evq064.
13. Azencott CA. Network-Guided Biomarker Discovery. In: Machine Learning for Health Informatics. vol. 9605. Cham: Springer International Publishing; 2016. p. 319–336. Available from: http://link.springer.com/10.1007/978-3-319-50478-0_16.
14. Gwinner F, Boulday G, Vandiedonck C, Arnould M, Cardoso C, Nikolayeva I, et al. Network-based analysis of omics data: The LEAN method. *Bioinformatics*. 2016; p. btw676. doi:10.1093/bioinformatics/btw676.
15. Jia P, Zheng S, Long J, Zheng W, Zhao Z. dmGWAS: dense module searching for genome-wide association studies in protein–protein interaction networks. *Bioinformatics*. 2011;27(1):95–102. doi:10.1093/bioinformatics/btq615.
16. Dittrich MT, Klau GW, Rosenwald A, Dandekar T, Muller T. Identifying functional modules in protein-protein interaction networks: an integrated exact approach. *Bioinformatics*. 2008;24(13):i223–i231. doi:10.1093/bioinformatics/btn161.
17. Leiserson MDM, Vandin F, Wu HT, Dobson JR, Eldridge JV, Thomas JL, et al. Pan-cancer network analysis identifies combinations of rare somatic mutations across pathways and protein complexes. *Nature Genetics*. 2015;47(2):106–114. doi:10.1038/ng.3168.

18. Azencott CA, Grimm D, Sugiyama M, Kawahara Y, Borgwardt KM. Efficient network-guided multi-locus association mapping with graph cuts. *Bioinformatics*. 2013;29(13):i171–i179. doi:10.1093/bioinformatics/btt238.
19. Liu Y, Brossard M, Roqueiro D, Margaritte-Jeannin P, Sarnowski C, Bouzigon E, et al. SigMod: an exact and efficient method to identify a strongly interconnected disease-associated module in a gene network. *Bioinformatics*. 2017; p. btx004. doi:10.1093/bioinformatics/btx004.
20. Sinilnikova OM, Dondon MG, Eon-Marchais S, Damiola F, Barjhoux L, Marcou M, et al. GENESIS: a French national resource to study the missing heritability of breast cancer. *BMC Cancer*. 2016;16(1):13. doi:10.1186/s12885-015-2028-9.
21. Michailidou K, Lindström S, Dennis J, Beesley J, Hui S, Kar S, et al. Association analysis identifies 65 new breast cancer risk loci. *Nature*. 2017;551(7678):92–94. doi:10.1038/nature24284.
22. Mulligan AM, Couch FJ, Barrowdale D, Domchek SM, Eccles D, et al. Common breast cancer susceptibility alleles are associated with tumour subtypes in BRCA1 and BRCA2 mutation carriers: results from the Consortium of Investigators of Modifiers of BRCA1/2. *Breast Cancer Research*. 2011;13(6). doi:10.1186/bcr3052.
23. Rinella ES, Shao Y, Yackowski L, Pramanik S, Oratz R, Schnabel F, et al. Genetic variants associated with breast cancer risk for Ashkenazi Jewish women with strong family histories but no identifiable BRCA1/2 mutation. *Human Genetics*. 2013;132(5):523–536. doi:10.1007/s00439-013-1269-4.
24. Brisbin AG, Asmann YW, Song H, Tsai YY, Aakre JA, Yang P, et al. Meta-analysis of 8q24 for seven cancers reveals a locus between NOV and ENPP2 associated with cancer development. *BMC Medical Genetics*. 2011;12(1):156. doi:10.1186/1471-2350-12-156.
25. SEARCH, The GENICA Consortium, kConFab, Australian Ovarian Cancer Study Group, Ahmed S, Thomas G, et al. Newly discovered breast cancer susceptibility loci on 3p24 and 17q23.2. *Nature Genetics*. 2009;41(5):585–590. doi:10.1038/ng.354.
26. Nielsen FC, van Overeem Hansen T, Sørensen CS. Hereditary breast and ovarian cancer: new genes in confined pathways. *Nature Reviews Cancer*. 2016;16(9):599–612. doi:10.1038/nrc.2016.72.
27. Quigley DA, Fiorito E, Nord S, Van Loo P, Alnaes GG, Fleischer T, et al. The 5p12 breast cancer susceptibility locus affects MRPS30 expression in estrogen-receptor positive tumors. *Molecular Oncology*. 2014;8(2):273–284. doi:10.1016/j.molonc.2013.11.008.
28. Yu M, Li R, Zhang J. Repositioning of antibiotic levofloxacin as a mitochondrial biogenesis inhibitor to target breast cancer. *Biochemical and Biophysical Research Communications*. 2016;471(4):639–645. doi:10.1016/j.bbrc.2016.02.072.
29. Liu G, Claret FX, Zhou F, Pan Y. Jab1/COPS5 as a Novel Biomarker for Diagnosis, Prognosis, Therapy Prediction and Therapeutic Tools for Human Cancer. *Frontiers in Pharmacology*. 2018;9:135. doi:10.3389/fphar.2018.00135.
30. de los Campos G, Vazquez AI, Hsu S, Lello L. Complex-Trait Prediction in the Era of Big Data. *Trends in Genetics*. 2018;34(10):746–754. doi:10.1016/j.tig.2018.07.004.

31. Nikolayeva I, Guitart Pla O, Schwikowski B. Network module identification—A widespread theoretical bias and best practices. *Methods*. 2018;132:19–25. doi:10.1016/j.ymeth.2017.08.008.
32. Ioachim E, Charchanti A, Briassoulis E, Karavasilis V, Tsanou H, Arvanitis DL, et al. Immunohistochemical expression of extracellular matrix components tenascin, fibronectin, collagen type IV and laminin in breast cancer: their prognostic value and role in tumour invasion and progression. *European Journal of Cancer*. 2002;38(18):2362–2370. doi:10.1016/s0959-8049(02)00210-1.
33. Yi W, Xiao E, Ding R, Luo P, Yang Y. High expression of fibronectin is associated with poor prognosis, cell proliferation and malignancy via the NF- κ B/p53-apoptosis signaling pathway in colorectal cancer. *Oncology Reports*. 2016;36(6):3145–3153. doi:10.3892/or.2016.5177.
34. Sponzillo M, Rosignolo F, Celano M, Maggisano V, Pecce V, Rose RFD, et al. Fibronectin-1 expression is increased in aggressive thyroid cancer and favors the migration and invasion of cancer cells. *Molecular and Cellular Endocrinology*. 2016;431:123–132. doi:10.1016/j.mce.2016.05.007.
35. Ahmed S, Thomas G, Ghousaini M, Healey CS, Humphreys MK, Platte R, et al. Newly discovered breast cancer susceptibility loci on 3p24 and 17q23.2. *Nature Genetics*. 2009;41(5):585–590. doi:10.1038/ng.354.
36. Breyer J, Dorset D, Clark T, Bradley K, Wahlfors T, McReynolds K, et al. An Expressed Retrogene of the Master Embryonic Stem Cell Gene POU5F1 Is Associated with Prostate Cancer Susceptibility. *The American Journal of Human Genetics*. 2014;94(3):395–404. doi:10.1016/j.ajhg.2014.01.019.
37. Chen HY, Chen RH. Cullin 3 Ubiquitin Ligases in Cancer Biology: Functions and Therapeutic Implications. *Frontiers in Oncology*. 2016;6. doi:10.3389/fonc.2016.00113.
38. Loignon M, Miao W, Hu L, Bier A, Bismar TA, Scrivens PJ, et al. Cul3 overexpression depletes Nrf2 in breast cancer and is associated with sensitivity to carcinogens, to oxidative stress, and to chemotherapy. *Molecular Cancer Therapeutics*. 2009;8(8):2432–2440. doi:10.1158/1535-7163.mct-08-1186.
39. Halsey LG, Curran-Everett D, Vowler SL, Drummond GB. The fickle P value generates irreproducible results. *Nature Methods*. 2015;12(3):179–185. doi:10.1038/nmeth.3288.
40. Wagner A, Wright J. Alternative routes and mutational robustness in complex regulatory networks. *Biosystems*. 2007;88(1-2):163–172. doi:10.1016/j.biosystems.2006.06.002.
41. Das J, Yu H. HINT: High-quality protein interactomes and their applications in understanding human disease. *BMC Systems Biology*. 2012;6(1):92. doi:10.1186/1752-0509-6-92.
42. Gallagher MD, Chen-Plotkin AS. The Post-GWAS Era: From Association to Function. *The American Journal of Human Genetics*. 2018;102(5):717–730. doi:10.1016/j.ajhg.2018.04.002.
43. Xie T. Analysis of the Gene-Dense Major Histocompatibility Complex Class III Region and Its Comparison to Mouse. *Genome Research*. 2003;13(12):2621–2636. doi:10.1101/gr.1736803.

44. Mishra A, Macgregor S. VEGAS2: Software for More Flexible Gene-Based Testing. *Twin Research and Human Genetics*. 2015;18(1):86–91. doi:10.1017/thg.2014.79.
45. Nakka P, Raphael BJ, Ramachandran S. Gene and Network Analysis of Common Variants Reveals Novel Associations in Multiple Complex Diseases. *Genetics*. 2016;204(2):783–798. doi:10.1534/genetics.116.188391.
46. Ionita-Laza I, Lee S, Makarov V, Buxbaum J, Lin X. Sequence Kernel Association Tests for the Combined Effect of Rare and Common Variants. *The American Journal of Human Genetics*. 2013;92(6):841–853. doi:10.1016/j.ajhg.2013.04.015.
47. Wang M, Huang J, Liu Y, Ma L, Potash JB, Han S. COMBAT: A Combined Association Test for Genes Using Summary Statistics. *Genetics*. 2017;207(3):883–891. doi:10.1534/genetics.117.300257.
48. Di Tommaso P, Chatzou M, Floden EW, Barja PP, Palumbo E, Notredame C. Nextflow enables reproducible computational workflows. *Nature Biotechnology*. 2017;35(4):316–319. doi:10.1038/nbt.3820.
49. Sakoda LC, Jorgenson E, Witte JS. Turning of COGS moves forward findings for hormonally mediated cancers. *Nature Genetics*. 2013;45(4):345–348. doi:10.1038/ng.2587.
50. Chang CC, Chow CC, Tellier LC, Vattikuti S, Purcell SM, Lee JJ. Second-generation PLINK: rising to the challenge of larger and richer datasets. *GigaScience*. 2015;4(1):7. doi:10.1186/s13742-015-0047-8.
51. Frankish A, Diekhans M, Ferreira AM, Johnson R, Jungreis I, Loveland J, et al. GENCODE reference annotation for the human and mouse genomes. *Nucleic Acids Research*. 2019;47(D1):D766–D773. doi:10.1093/nar/gky955.
52. Wang Q, Jia P. dmGWAS 3.0; 2014. <https://bioinfo.uth.edu/dmGWAS/>.
53. Ljubić I, Weiskircher R, Pferschy U, Klau GW, Mutzel P, Fischetti M. An Algorithmic Framework for the Exact Solution of the Prize-Collecting Steiner Tree Problem. *Mathematical Programming*. 2006;105(2-3):427–449. doi:10.1007/s10107-005-0660-x.
54. Beisser D, Klau GW, Dandekar T, Muller T, Dittrich MT. BioNet: an R-Package for the functional analysis of biological networks. *Bioinformatics*. 2010;26(8):1129–1130. doi:10.1093/bioinformatics/btq089.
55. Dittrich M, Beisser D. BioNet; 2008. <https://bioconductor.org/packages/BioNet/>.
56. Scheid S, Spang R. twilight; a Bioconductor package for estimating the local false discovery rate. *Bioinformatics*. 2005;21(12):2921–2922. doi:10.1093/bioinformatics/bti436.
57. Leiserson MDM, Vandin F, Wu HT, Dobson JR, Eldridge JV, Thomas JL, et al.. HotNet2; 2018. <https://github.com/raphael-group/hotnet2>.
58. Gwinner F. LEANR; 2016. <https://cran.r-project.org/web/packages/LEANR/>.

59. Climente-González H, Azencott CA. martini; 2019. <https://www.bioconductor.org/packages/martini/>.
60. Liu Y. SigMod v2; 2018. <https://github.com/YuanlongLiu/SigMod>.
61. Jassal B, Matthews L, Viteri G, Gong C, Lorente P, Fabregat A, et al. The reactome pathway knowledgebase. Nucleic Acids Research. 2019;doi:10.1093/nar/gkz1031.
62. Yu G, He QY. ReactomePA: an R/Bioconductor package for reactome pathway analysis and visualization. Molecular BioSystems. 2016;12(2):477–479. doi:10.1039/c5mb00663e.
63. Benjamini Y, Hochberg Y. Controlling the False Discovery Rate: A Practical and Powerful Approach to Multiple Testing. Journal of the Royal Statistical Society: Series B (Methodological). 1995;57(1):289–300. doi:10.1111/j.2517-6161.1995.tb02031.x.
64. Nogueira S, Brown G. Measuring the Stability of Feature Selection. In: Machine Learning and Knowledge Discovery in Databases. vol. 9852. Cham: Springer International Publishing; 2016. p. 442–457. Available from: http://link.springer.com/10.1007/978-3-319-46227-1_28.
65. Metcalfe KA, Finch A, Poll A, Horsman D, Kim-Sing C, Scott J, et al. Breast cancer risks in women with a family history of breast or ovarian cancer who have tested negative for a BRCA1 or BRCA2 mutation. British Journal of Cancer. 2008;100(2):421–425. doi:10.1038/sj.bjc.6604830.
66. Piñero J, Bravo A, Queralt-Rosinach N, Gutiérrez-Sacristán A, Deu-Pons J, Centeno E, et al. DisGeNET: a comprehensive platform integrating information on human disease-associated genes and variants. Nucleic Acids Research. 2017;45(D1):D833–D839. doi:10.1093/nar/gkw943.
67. The 1000 Genomes Project Consortium, Gibbs RA, Boerwinkle E, Doddapaneni H, Han Y, Korchina V, et al. A global reference for human genetic variation. Nature. 2015;526(7571):68–74. doi:10.1038/nature15393.

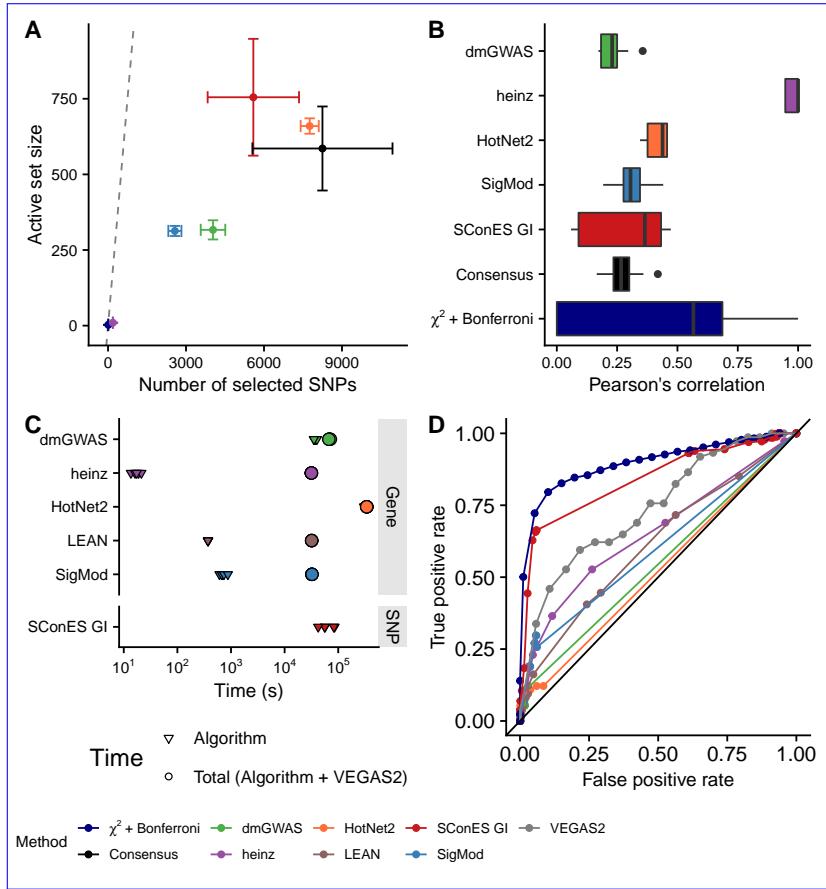


Fig 3. Comparison of network-based GWAS methods on GENESIS. **Comparison of network methods on GENESIS.** Each method was run 5 times on a random subset containing 80% of the samples, and tested on the remaining samples (Section 2.2–4.5). As LEAN did not select any gene, it was excluded from all panels except D and B. (A) Number of SNPs selected by each method and number of SNPs in the active set found (i.e., the number of SNPs selected by the classifier, Section 4.5.2). 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, indicating the upper bound of the active set. For reference, the active set of Lasso the classifier using all the SNPs as input 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 correlations of the solutions used produced by different methods. A Pearson correlation of 1 means the two solutions are the same. A Pearson correlation of 0 means that there is no SNP in common between the two solutions. (D) (C) Runtime of the evaluated methods, by type of network used (gene PPIN or SNP). For gene network-based gene-based methods, inverted triangles represent the runtime of the algorithm itself alone, and circles the total time, which includes the algorithm themselves and the additional 119 980 seconds (1 day and 9.33 hours) that VEGAS2 took on average to compute the gene scores from SNP summary statistics. (D) True positive rate and true negative rate of the methods, obtained using different parameter combinations (Section 2.7). We used as true positives BCAC-significant SNPs (for SConES and χ^2 + Bonferroni) and genes (for the remaining methods, Section 4.6). We used the whole dataset in this panel.

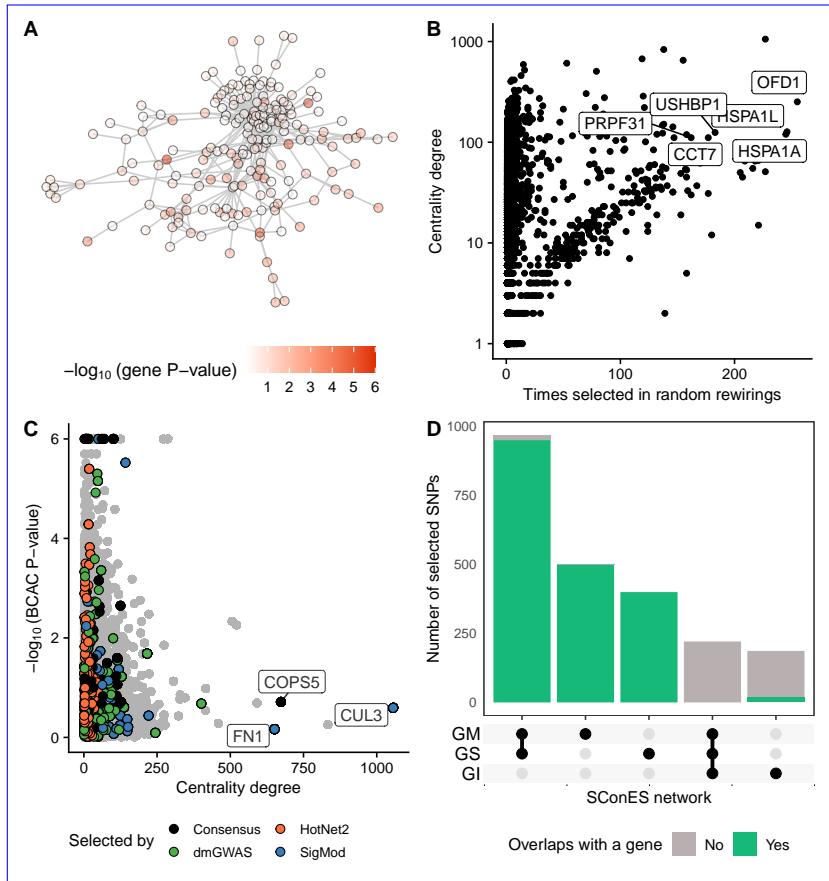


Fig 4. Drawbacks encountered when using network-guided methods. **Drawbacks encountered when using network methods.** (A) DmGWAS solution subnetwork. Genes with $\text{-}\log_{10}$ of their P-value > 0.1 are highlighted in red. (B) Number of times a gene was selected by either dmGWAS, heinz, LEAN, or SigMod in 100 rewirings of the PPIN (Section 4.4) and its centrality degree. (C) Centrality degree and $-\log_{10}$ of the VEGAS2 P-value in BCAC for each of the nodes in the PPIN. We highlighted the genes selected by each method, and the ones selected by more than one (“Consensus”). We labeled the three most central genes that were picked by any method. (E)(D) Overlap between the solutions of SConES GS, GM, or GI in the different genomic regions. Barplots are colored based on whether the SNPs that were map to a gene or not selected in the studied network, but were selected in another one, are displayed in background color (Section 4.3.5).

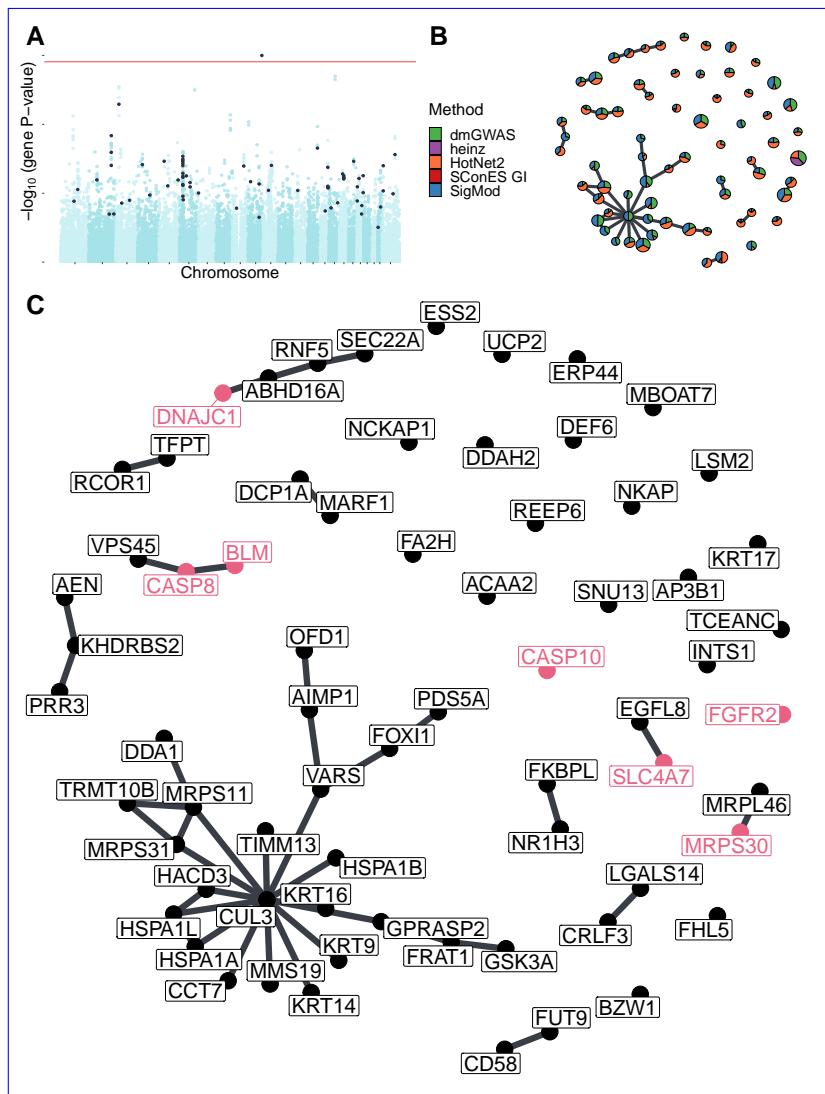


Fig 5. Stable consensus solution on GENESIS (Section 2.8).
(A) Manhattan plot of genes; in black, the ones in the stable consensus subnetwork on GENESIS solution. The red line indicates the Bonferroni threshold (Section 1.53 2.8×10^{-6} for genes). **(B)** Stable consensus network. Each node gene is represented by a pie chart, which shows the methods that selected it. We labeled (and enlarged) the most central genes gene (*CUL3*) and those genes that are, the known breast cancer susceptibility genes, and /or significantly associated with breast cancer susceptibility in the BCAC dataset BCAC-significant genes (Section 4.6). **(C)** The latter ones nodes are also in the same disposition as in panel B, but we indicated every gene name. We colored in pink. All gene the names are indicated in of known breast cancer susceptibility genes and BCAC-significant genes.