

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}, with the GENESIS study group[¶]

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

Network methods provide a comprehensive approach to uncovering the genetics of complex diseases and building hypotheses. They discover susceptibility genes by jointly consider the statistical association between genetic variation and a phenotype, measured in a genome-wide association study (GWAS), and the biological context of each gene, represented as a network. Hence, a network method could select a gene with a high P-value of association if it connects multiple low P-value genes in the network. In this work, we evaluate six network methods which identify subnetworks with high scores of GWAS association with a phenotype. This allows them to give more compelling results than standard SNP- and gene-level GWAS analyses, recovering causal subnetworks tightly related to cancer susceptibility. We applied them to GENESIS, a GWAS on French women with a family history of breast cancer and who tested negative for pathogenic variants in *BRCA1* and *BRCA2*. We critically compared these six methods, discussing the impact of their different mathematical frameworks, and the parameter choices. Additionally, we performed an in-depth benchmarking with respect to the size and predictive power of their solutions, as well as their stability and runtimes. Importantly, we found significant overlaps between the genes in five of the solution networks and the genes significantly associated in the largest GWAS on susceptibility to breast cancer. Since most of the methods produced reasonable solutions, we proposed to combine them into a consensus solution, containing the genes selected by at least two methods. This aggregation brought further insights. For instance, it contained *COPS5*, a gene related to multiple hallmarks of cancer, and 14 of its neighbors. The main drawback of network methods with regards to conventional χ^2 -based GWAS was instability i.e. network methods' outputs changed more in front of small perturbations in the input. To compensate for this instability, we proposed a stable consensus solution, formed by the most consistently selected genes in different subsamples of the data. This

stable consensus was 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 most genes. The network was organized around *CUL3*, which encodes a ubiquitin ligase related protein that regulates the protein levels of several genes involved in cancer progression. In conclusion, we showed how network address limitations of GWAS, namely their lack of statistical power and the difficulty of their interpretation. Project-agnostic implementations of each of the network methods are available at <https://github.com/hclimente/gwas-tools> to facilitate their application to other GWAS datasets.

Author summary

In genome-wide association studies (GWAS), 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, 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. The omnigenic model offers an appealing explanation: only a limited number of core genes are directly involved in the disease; but gene functions are deeply interrelated so that many other genes 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 apply and compare 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. 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 statistical test of association (e.g. based on logistic regression) is conducted between each individual 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 our understanding of the disease.

GWAS have successfully identified thousands of variants underlying many common diseases [2]. However, this experimental setting also presents intrinsic challenges. Some of them stem from the high dimensionality of the problem, as every GWAS to date studies more variants than samples are genotyped. This limits the statistical power of the experiment, as only variants with 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, 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 it is the case in GWAS [4]. Another open issue is the interpretation of the results, as the functional consequences of most common variants

are 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 strongly inter-related and influence each other's 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 includes the structure of biological data might help alleviate the aforementioned issues.

For this reason, many authors turn to network biology to handle the complex interplay of biomolecules that lead to disease [7]. As its name suggests, network biology models biology as a network, where the biomolecules under study, often genes, are nodes, and selected functional relationships are edges that link them. These relationships come from evidence that the genes jointly contribute to a biological function; for instance, their expression is 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 [8]. In fact, an examination of biological networks shows that disease genes have differential properties [8,9]: they tend to occupy central positions in the network, interconnecting different modules TODO. Therefore, studying the neighborhood of disease-associated genes is effective at identifying new ones that are involved in the same biological functions [10].

Network-based discovery methods exploit the differential properties described above to identify disease genes on GWAS data [11]. In essence, each gene is assigned a score of association with the disease, computed from the GWAS data, and biological relationships, given by a network built on prior knowledge. Then, the problem becomes finding a functionally-related set of highly-scoring genes. Multiple solutions have been proposed to this problem, often stemming from different mathematical frameworks and considerations of what the optimal solution looks like. For example, some methods restrict the problem to specific types of subnetworks. Such is the case of LEAN [12], which focuses on "star" subnetworks, i.e. instances were both a gene and its direct interactors are associated with the disease. Other algorithms, like dmGWAS [13] and heinz [14], 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 [15], SConES [16], and SigMod [17].

In this work, we analyze the application of these six network methods on GWAS data. They use different interpretations of the omnigenic model, and provide a representative view of the field. We worked on the GENESIS dataset [18], a study on familial breast cancer conducted in the French population. After a classical GWAS approach, we use these network methods to identify additional breast cancer susceptibility genes. Lastly, we carry out a comparison of the solutions obtained by the different methods, and aggregate them to obtain a consensus solution of predisposition to familial breast cancer.

2 Results

2.1 Conventional SNP- and gene-based analyses retrieve the *FGFR2* locus in the GENESIS dataset

We conducted association analyses in the GENESIS dataset (Section 4.1) at both SNP and gene levels (Section 4.2). Two genomic regions had a P-value lower than the Bonferroni threshold on chromosomes 10 and 16 (S2 FigA). The former overlaps with

gene *FGFR2*; the latter with *CASC16* the protein-coding gene *TOX3*. Variants in both *FGFR2* and *TOX3* have been repeatedly associated with breast cancer susceptibility in other case-control studies [19], *BRCA1* and *BRCA2* carrier studies [20], and in hereditary breast and ovarian cancer families negative for mutations in *BRCA1* and *BRCA2* [21]. In GENESIS only *FGFR2* was significantly associated with breast cancer at the gene-level (S2 FigB).

Closer examination reveals two other regions (3p24 and 8q24) having low, albeit not genome-wide significant, P-values. Both of them have been associated to breast cancer susceptibility in the past [22, 23]. We applied an L1-penalized logistic regression using all GENESIS genotypes as input, and the phenotype (cancer/healthy) as outcome (Section 4.5.2). The algorithm selected 100 SNPs, both from all aforementioned regions and new ones (S2 FigC). Yet, it is 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 4.3.5) are not interconnected in the protein-protein interaction network (PPIN, Section 4.3.2). Moreover, the classification performance of the model is low (sensitivity = 55%, specificity = 55%, Section 4.5). Together, these issues motivate exploring network methods, which consider not only statistical association, but also the location of each gene in a PPIN to find susceptibility 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 (Benjamini-Hochberg [BH] correction adjusted P-value < 0.05), we obtained six solutions (Fig 1): one for each of the remaining four gene-based methods, one for SConES GI (which works at the SNP level), and the consensus.

These solutions differ in many aspects, making it hard to draw joint conclusions. For starters, the overlap between the genes featured in each solution is quite small (Fig 1A). However, the methods tend 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 the solution size: the largest solution, produced by HotNet2, contains 440 genes, while heinz's contains only 4 genes. While SConES GI did not recover any protein coding gene, working with SNP networks rather than gene networks allowed it to retrieve four subnetworks in intergenic regions, and another one overlapping an RNA gene (*RNU6-420P*).

The topologies of the six solutions differ as well (Fig 1C), as measured by the median centrality and the number of connected components (Table 1). Only two methods yield more than one connected component: SConES, as described above, and HotNet2. HotNet2 produced 135 subnetworks, 115 of which have fewer than five genes. The second largest subnetwork (13 nodes) contains the two breast cancer susceptibility genes *CASP8* and *BLM*.

Lastly, a pathway enrichment analysis (Section 4.4) also showed similarities and differences between the solutions of the different methods. It linked different parts of SigMod's solution to four processes (S3 Table): protein translation (including mitochondrial), mRNA splicing, protein misfolding, and keratinization (BH adjusted P-values < 0.03). Interestingly, the dmGWAS solution (S4 Table) is also related to protein misfolding (*attenuation phase*, BH adjusted P-value = 0.01). However, it additionally 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 [24]. As with SigMod, the genes in HotNet2's solution (S5 Table) 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). 136
137
138
139

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

Network	# genes	# edges	# components	Betweenness	\hat{P}_{gene}	# genes in consensus
HINT HT	13 619	142 541	15	16 706	0.46	93/93
dmGWAS	194	450	1	49 115	0.19	55/93
heinz	4	3	1	113 633	0.001	4/93
HotNet2	440	374	130	7 739	0.048	63/93
LEAN	0	0	0	-	-	0/93
SConES GI	0 (1)	0	0	-	-	0/93
SigMod	142	249	11	92 603	0.008	84/93
Consensus	93	186	21	50 737	0.006	93/93
Stable consensus	68	49	32	94 854	0.005	43/93

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 betweenness of the selected genes in the PPIN. **\hat{P}_{gene} :** median VEGAS2 P-value of the selected genes. **# genes in consensus:** Number of genes in common between the method's solution, and the 93 genes in the consensus solution. 140
141
142
143
144
145
146
147

Despite their differences, there are additional common themes. All obtained solutions have 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 1D). This illustrates the trade-off between controlling for type I error and biological relevance. However, there are nuances between solutions in this regard: heinz strongly favors genes with lower P-values, while dmGWAS is less conservative (median VEGAS2 P-values 0.0012 and 0.19, respectively); SConES tends to select whole LD-blocks; and HotNet2 and SigMod are less likely to select low scoring genes. 148
149
150
151
152
153
154
155
156
157
158

Additionally, the solutions presented other desirable properties. First, five 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 solutions displayed on average a 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 [9], 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 also more central than the average SNP (S1 Table), suggesting that causal SNPs are also more central than non-associated SNPs. 159
160
161
162
163
164
165
166
167

2.3 A case study: the consensus solution

Despite their shared properties, the differences between the solutions of the different methods suggest that each of them captures different aspects of cancer susceptibility. Indeed, out of the 668 genes that are selected by at least one method, only 93 are 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 (Fig 1B), a relationship that plateaus at 2. Hence, to leverage on their strengths and compensate their respective weaknesses, we built a consensus solution that captures the genes shared among at least two solutions (Section 4.3.3). This solution (Fig 2) contains 93 genes and 168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917
918
919
910
911
912
913
914
915
916
917
918
919
920
921
922
923
924
925
926
927
928
929
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025
1026
1027
1028
1029
1020
1021
1022
1023
1024
1025
1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079
1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1110
1111
1112
1113
1114
1115
1116
1117
1118
1119
1110
1111
1112
1113
1114
1115
1116
1117
1118
1119
1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133
1134
1135
1136
1137
1138
1139
1130
1131
1132
1133
1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187
1188
1189
1180
1181
1182
1183
1184
1185
1186
1187
1188
1189
1190
1191
1192
1193
1194
1195
1196
1197
1198
1199
1190
1191
1192
1193
1194
1195
1196
1197
1198
1199
1200
1201
1202
1203
1204
1205
1206
1207
1208
1209
1200
1201
1202
1203
1204
1205
1206
1207
1208
1209
1210
1211
1212
1213
1214
1215
1216
1217
1218
1219
1210
1211
1212
1213
1214
1215
1216
1217
1218
1219
1220
1221
1222
1223
1224
1225
1226
1227
1228
1229
1220
1221
1222
1223
1224
1225
1226
1227
1228
1229
1230
1231
1232
1233
1234
1235
1236
1237
1238
1239
1230
1231
1232
1233
1234
1235
1236
1237
1238
1239
1240
1241
1242
1243
1244
1245
1246
1247
1248
1249
1240
1241
1242
1243
1244
1245
1246
1247
1248
1249
1250
1251
1252
1253
1254
1255
1256
1257
1258
1259
1250
1251
1252
1253
1254
1255
1256
1257
1258
1259
1260
1261
1262
1263
1264
1265
1266
1267
1268
1269
1260
1261
1262
1263
1264
1265
1266
1267
1268
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295
1296
1297
1298
1299
1290
1291
1292
1293
1294
1295
1296
1297
1298
1299
1300
1301
1302
1303
1304
1305
1306
1307
1308
1309
1300
1301
1302
1303
1304
1305
1306
1307
1308
1309
1310
1311
1312
1313
1314
1315
1316
1317
1318
1319
1310
1311
1312
1313
1314
1315
1316
1317
1318
1319
1320
1321
1322
1323
1324
1325
1326
13

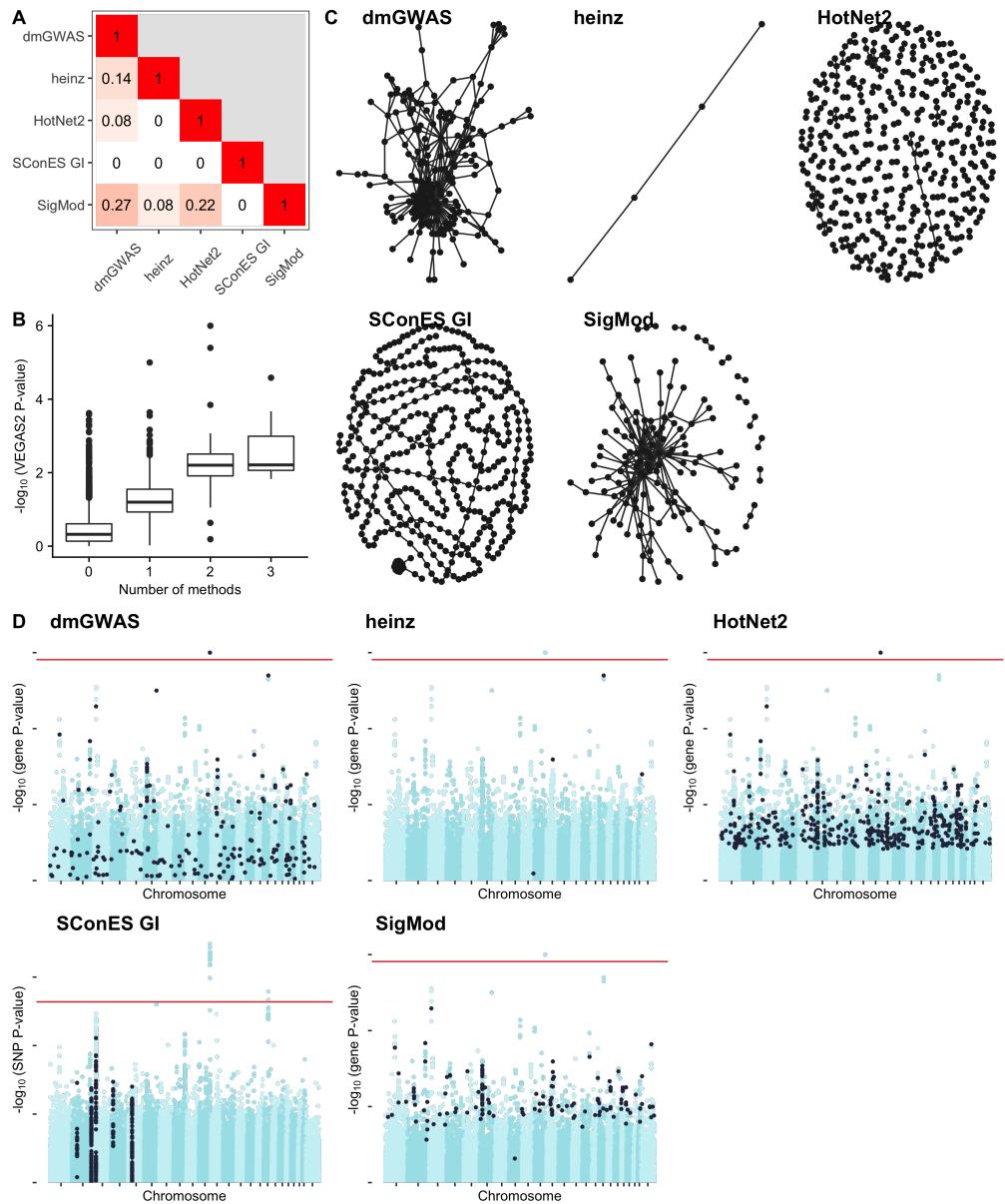


Fig 1. Overview of the solutions produced by the different network methods (Section 4.3.3) on the GENESIS dataset. As LEAN did not produce any significant solution (BH adjusted P-value < 0.05), it was excluded. Unless indicated otherwise, results refer to genes, except for SConES GI which are at the SNP-level. **(A)** Overlap between the genes selected by each of the methods, measured by Pearson correlation between indicator vectors. **(B)** 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 Bonferroni threshold is indicated by a red line (2.54×10^{-7} for SNPs, 1.53×10^{-6} for genes).

exhibits 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 solution also shows

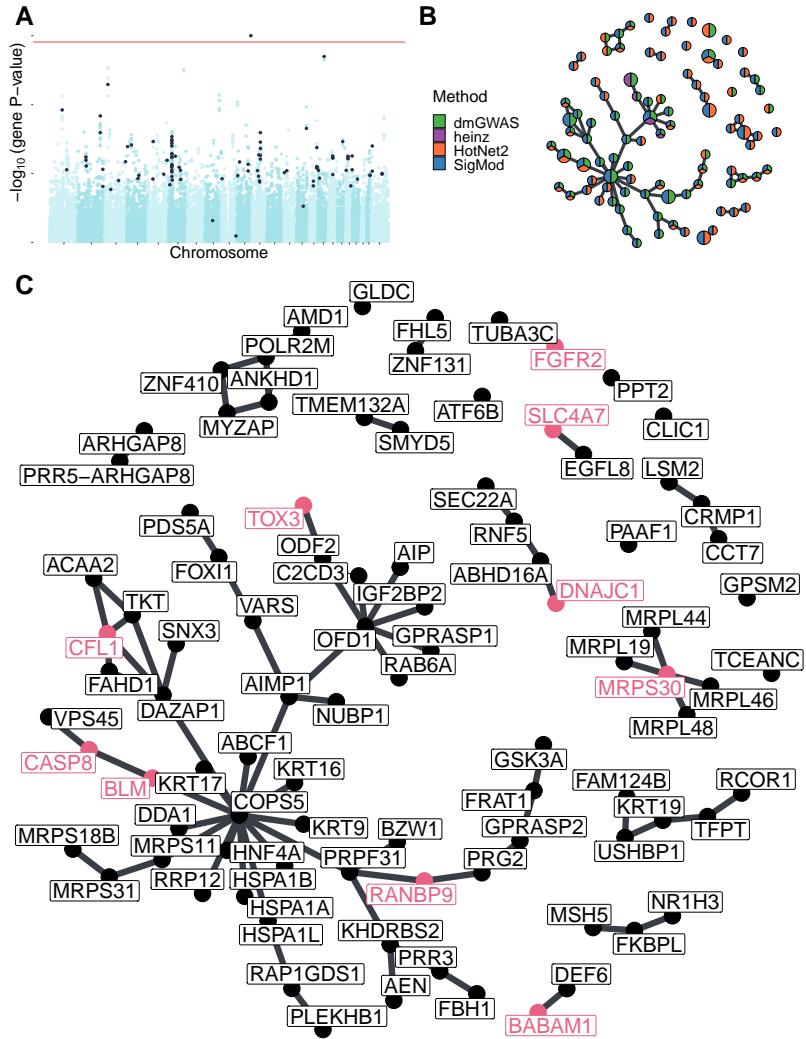


Fig 2. Consensus solution on GENESIS (Section 4.3.3) . **A** Manhattan plots of genes; in black, the ones in the consensus solution. The Bonferroni threshold is indicated by a red line (1.53×10^{-6} for genes). **B** Consensus network. Each gene is represented by a pie chart, which shows the methods that selected it. We 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 (Section 4.5.3). **C** The nodes are in the same disposition as in panel B, but every gene name is indicated. We colored in pink the names of genes that are known breast cancer susceptibility genes and/or significantly associated with breast cancer susceptibility in the BCAC dataset

similar pathways as the individual solutions (S6 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 [25]. Interestingly, increased mitochondrial translation has been found in cancer cells [26], and its inhibition proposed as a therapeutic target. With regards to the attenuation phase of heat shock response, it involves 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 to any of these three genes, 9 map to all of them, and 4 to two, making it hard to disentangle their effects. *HSPA1A* was the most strongly associated gene (VEGAS2 P-value = 8.37×10^{-4}).
178
179
180
181

Topologically, the consensus consists of a connected component composed of 49 genes, and multiple smaller subnetworks. Among the latter, 19 genes are in subnetworks containing a single gene or two connected nodes, implying that they do not have a consistently altered neighborhood, but are strongly associated themselves and hence picked by at least two methods. The large connected component contains genes that are highly central in the PPIN, a property which is weakly anti-correlated with the P-value of association to the disease (Pearson correlation coefficient = -0.26, S4 Fig). This suggests that these genes were selected because they were on the shortest path between another two highly associated genes. In view of this, we hypothesize that highly central genes might contribute to the heritability through alterations of their neighborhood, consistently with the omnigenic model of disease [6]. For instance, the most central node in the consensus solution is COPS5, a component of the COP9 signalosome which regulates multiple signaling pathways. *COPS5* is related to multiple hallmarks of cancer and is overexpressed in multiple tumors, including breast and ovarian cancer [27]. Despite its lack of association in GENESIS or in studies conducted by the Breast Cancer Association Consortium (BCAC) [19] (VEGAS2 P-value of 0.22 and 0.14 respectively), its neighbors in the consensus solution have consistently low P-values (median VEGAS2 P-value = 0.006).
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199

2.4 Network methods boost discovery

We compared the results obtained with different network methods to the European sample of BCAC, the largest GWAS to date on breast cancer (Section 4.5.3). Although BCAC case-control studies do not necessarily target cases with a familial history of breast cancer, 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 with BCAC's scale (90 times more samples than GENESIS) provides a reasonable counterfactual of what we would expect if GENESIS had a larger sample size. We computed a gene association score on BCAC (Section 4.5.3). The solutions provided by the different network methods overlap significantly with BCAC findings (Fisher's exact test P-value < 0.019). The gene-based methods achieve comparable precision (2%-25%) and recall (1.3-12.1%) at recovering BCAC-significant genes (S5 FigA). Interestingly, while SConES GI, at the SNP-level, achieves a similar recall (8.6%), it shows a much higher precision (47.3%).
200
201
202
203
204
205
206
207
208
209
210
211
212
213

2.5 Network methods share limitations

We compared the six network methods in a 5-fold subsampling setting (Section 4.5). This allowed us to measure five properties (Fig 3): size of the solution; sensitivity and specificity of an L1-penalized logistic regression classifier on the selected SNPs; stability; and 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 (Section 4.3.5). The largest solutions 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. 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 4.5.3, Fig 3B). SConES exhibits the largest area under the curve, since, when $\lambda = 0$ (network topology is disregarded), it is equivalent to a
214
215
216
217
218
219
220
221
222
223
224
225

Bonferroni correction. Between the remaining network methods, they have similar areas under the curve, with heinz being the one with the largest area.

To determine whether the selected SNPs could be used for patient classification, we computed the performance of the classifier on the *test dataset* (S5 FigB). The different classifiers displayed similarly poor 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 prediction error. Of course, although these performances are low, we do not expect to separate cases from controls well using exclusively genetic data [28].

Another desirable quality of a selection algorithm is the stability of the solution with respect to small changes in the input (Section 4.5.1). Heinz was highly stable in our benchmark, while the other methods displayed similarly low stabilities (Fig 3C).

In terms of computational runtime, the fastest method was heinz (Fig 3D), 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 down considerably gene-based methods; on this benchmark, that step took on average 1 day and 9.33 hours. Including this first step, it therefore took 5 days on average for HotNet2 to produce a result.

2.6 Network topology matters, and might lead to ambiguous results

As shown above, and despite their similarities, the network methods produced remarkably different solutions. This is due to the fact that each of them models the problem differently. Importantly, understanding which assumptions they make allows to understand the results more in depth. For instance, the fact that LEAN did not return any gene implies that there is no gene such that both itself and its environment are on average strongly associated with the disease.

In the GENESIS dataset, heinz’s solution is very conservative, providing a small solution with the lowest median P-value (Table 1). By repeatedly selecting this compact solution, heinz was the most stable method (Fig 3C). Its conservativeness stems from its preprocessing step, which models 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 a positive score after that transformation. Yet, this small solution does not provide much insight into the susceptibility mechanisms to cancer. Importantly, it ignores genes that are associated to cancer in this dataset like *FGFR2*.

On the other end of the spectrum, dmGWAS, HotNet2, and SigMod produced large solutions. dmGWAS’ solution is the lowest scoring solution on average. This is due to the greedy framework it uses, which has a bias for larger solutions [29]. It considered all nodes at distance 2 of the examined subnetwork, and accepted a weakly associated gene if it was linked to another, high scoring one. This is exacerbated when the results of successive greedy searches are aggregated, 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 solution is even harder to interpret, being composed of 440 genes divided into 135 subnetworks. Lastly, SigMod misses some of the highest scoring, breast cancer susceptibility genes in the dataset, like *FGFR2* and *TOX3*.

Another peculiarity of network methods is their relationship to degree centrality. We studied random rewirings of the PPIN while preserving node centrality (Section 4.5.4). In this setting, network methods favored central genes, as they often connect high

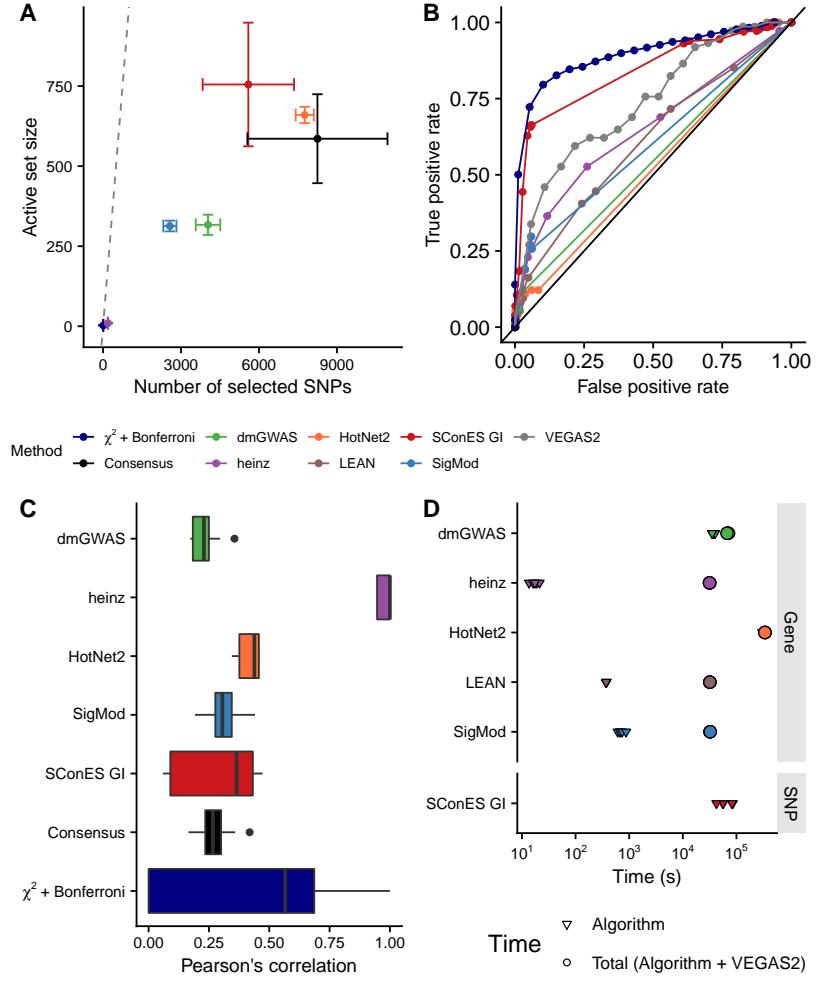


Fig 3. Comparison of network-based GWAS methods on GENESIS. Each method was run 5 times on a random subset containing 80% of the samples, and tested on the remaining samples (Section 4.5). As LEAN did not select any gene, it was excluded from all panels except **D**. **(A)** Number of SNPs selected by each method and number of SNPs in the active set found by the classifier (Section 4.5.2). Points are the average over the 5 runs; the error bars 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 (i.e. the number of SNPs in the solution). For reference, the active set of Lasso using all the SNPs included, on average, 154 117.4 SNPs. **(B)** True positive rate and true negative rate, using significant SNPs (for SConES and $\chi^2 + \text{Bonferroni}$) and genes (for the remaining methods) in BCAC (Section 4.5.3) as true positives, of multiple parameter combinations for different methods (Section 2.7). **(C)** Pairwise Pearson correlations of the solutions 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)** Runtime of the evaluated methods, by type of network used (PPIN or SNP). For gene-based methods, inverted triangles represent the runtime of the algorithm itself, and circles the total time, which includes the algorithm themselves and the additional 119 980 seconds (1 day and 9.33 hours) that VEGAS2 took on average to compute the gene scores from SNP summary statistics.

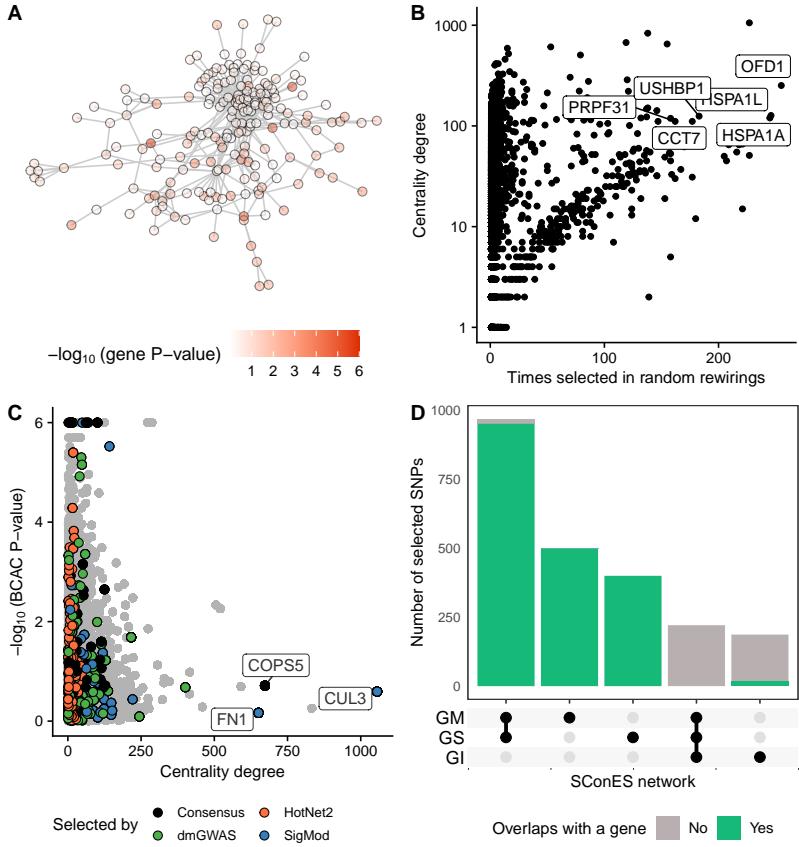


Fig 4. Drawbacks encountered when using network methods. (A) DmGWAS solution, with the genes colored according to the $-\log_{10}$ of their P-value. (B) Centrality degree and $-\log_{10}$ of the VEGAS2 P-value in BCAC for each of the nodes in the PPIN. We highlighted the genes selected by each method, and the ones selected by more than one (“Consensus”). We labeled the three most central genes that were picked by any method. (C) 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. (D) Overlap between the solutions of SConES GS, GM or GI in the different genomic regions. SNPs that were not selected in the studied network, but were selected in another one, are displayed in background color.

scoring nodes (Fig 4B). This is despite the fact that highly central genes often had no association to breast cancer susceptibility (Fig 4C). This was especially the case of SigMod, which selected three highly central, unassociated genes in both the PPIN and in many of the random rewirings: *COPS5*, *CUL3* and *FN1*. As we showed in Section 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 [30], and it is negatively correlated with poor prognosis in other cancer types [31,32].

By virtue of using a SNP subnetwork, SConES analyzes each SNP in their functional context. It can therefore select SNPs located in genes not included in the PPIN, as well as SNPs in non-coding regions or in non-interacting genes. We compared the solution of SConES in the GI network (using PPIN information), to the one using only positional information (GS network), or positional and gene annotations (GM network).

Importantly, SConES produces similar results on the GS and GM networks (S3 Fig). While the solutions on those two greatly overlap with SConES GI's, they contain additional gene-coding segments (Fig 4C). In fact the formers' solutions are chromosome regions related to breast cancer, like 3p24 (*SLC4A7/NEK10* [33]), 5p12 (*FGF10, MRPS30* [25]), 10q26 (*FGFR2*), and 16q12 (*TOX3*). On top of those SConES GS selects region 8q24 (*POU5F1B* [34]).

2.7 Different parameters produce similarly-sized solutions

We explored the parameter space of the different methods by running them under different combinations of parameters (Section 4.3.4). In agreement with their formulations (Section 4.3.3), larger values of certain parameters produce less astringent solutions (S6 FigA): for HotNet2 and heinz, we examined the threshold to decide each gene has a positive score or zero; for dmGWAS, the d parameter controls how far neighbors could be added; node involved how far the search from the explored subnetwork; for SigMod $nmax$ specifies the maximum size of the solution; and for LEAN, it was the P-value threshold to consider a solution significant. Two parameters had the opposite effect (larger is more stringent): SigMod's maxjump , which sets the threshold to consider an increment in λ "large enough"; and SConES η , where higher values produce more stringent 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 an additional 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 size of the solutions was relatively robust to the choice of parameters (S6 FigB).

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 (S6 FigC). This showed that dmGWAS and SigMod are robust to certain parameters: dmGWAS' output was mostly determined by the parameter d , rather than r ; SigMod's by $nmax$, rather than maxjump .

SConES presented an interesting case in terms of feature selection: most of the explored combinations of parameters led to trivial solutions (either all the SNPs, or none of them were included) (S6 FigA). To explore a more meaningful parameter space, we selected the parameters in two rounds. First, we explored the whole sample space. Then, we focused in 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, and allowed to find more interesting solutions.

2.8 Building a stable consensus network preserves global network properties

Most of the network methods, including the consensus, were highly unstable, raising questions about the reliability of the results. We built a new, *stable consensus* 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 stability-based consensus was composed of 68 genes (Fig 5). This network shares most of the properties of the consensus: breast cancer susceptibility genes are overrepresented ($P\text{-value} = 3 \times 10^{-4}$), 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.

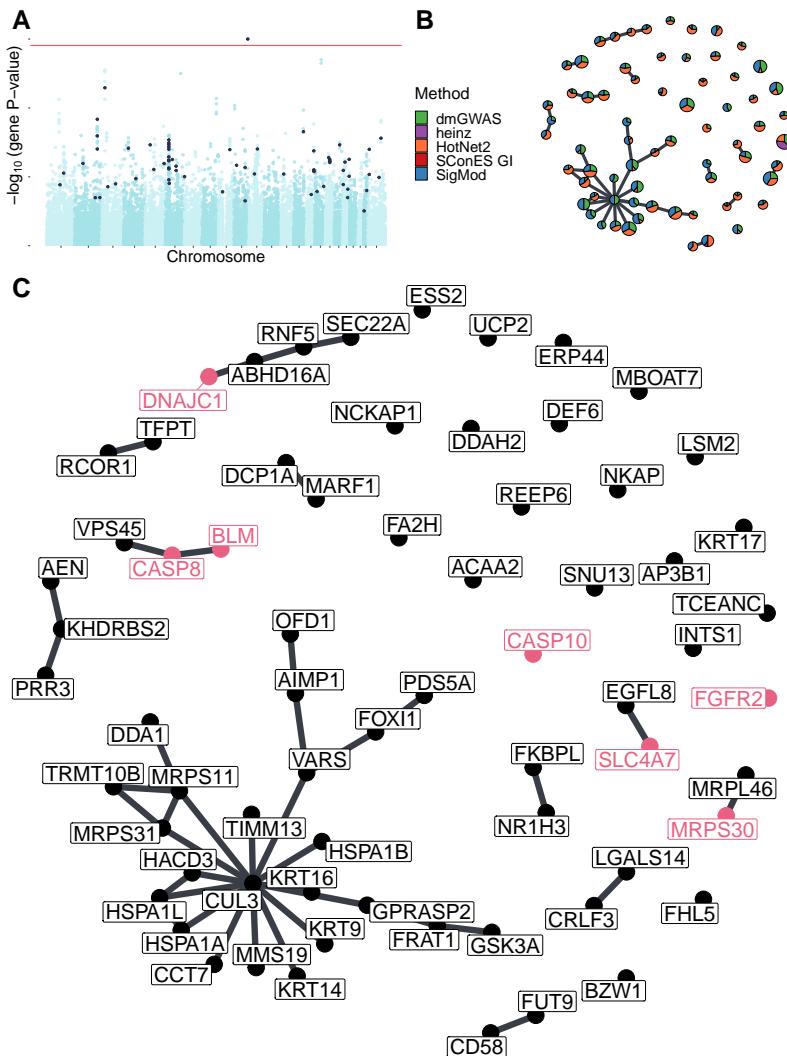


Fig 5. Stable consensus solution on GENESIS (Section 2.8). **A** Manhattan plots of genes; in black, the ones in the stable consensus solution. The Bonferroni threshold is indicated by a red line (1.53×10^{-6} for genes). **B** Stable consensus network. Each gene is represented by a pie chart, which shows the methods that selected it. We enlarged the most central gene (*CUL3*) and those genes that are known breast cancer susceptibility genes and/or significantly associated with breast cancer susceptibility in the BCAC dataset (Section 4.5.3). **C** The nodes are in the same disposition as in panel B, but every gene name is indicated. We colored in pink the names of genes that are known breast cancer susceptibility genes and/or significantly associated with breast cancer susceptibility in the BCAC dataset

Despite these similarities, the consensus and the stable consensus solutions include different genes. In the stable consensus network, the most central gene is *CUL3*, which is absent from the previous consensus solution 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 [35], and its overexpression was
345
also linked to increased sensitivity to carcinogens [36].
346

3 Discussion

In recent years, the ability of GWAS to unravel the mechanisms leading to complex
348 diseases has been called into question [6]. First, the omnigenic model proposes that gene
349 functions are interwoven with each other in a dense co-function network. The practical
350 consequence is that larger and larger GWAS will lead to the discovery of an
351 uninformative wide-spread pleiotropy. Second, discovery in GWAS is hindered by a
352 conservative statistical framework. Network methods elegantly address these two issues
353 by using both association scores and an interaction network to take into consideration
354 the biological context of each of the genes and SNPs. Based on what could be
355 considered diverse interpretations of the omnigenic model, several methods for
356 network-guided discovery have been proposed in recent years. In this article we
357 evaluated the relevance of six of these methods by applying them to the study of
358 GENESIS, a GWAS dataset on familial breast cancer.
359

DmGWAS, Heinz, HotNet2, SConES and SigMod all yield interesting solutions,
360 which include (but are not limited to) known breast cancer susceptibility genes. In
361 general, the selected genes and SNPs were more central than average, in accordance
362 with the observation that disease genes are more relatively central [9]. However, very
363 central nodes are also more likely to be connecting any given random pair of nodes,
364 making them more likely to be selected by network methods (Section 2.6). Yet, across
365 this article we show that highly central genes that were selected (*COPS5*, *CUL3* and
366 *FN1*) could plausibly be involved in breast cancer susceptibility. Despite these
367 similarities, the solutions obtained were notably different. At one end of the spectrum,
368 SConES and heinz preferred small, highly associated solutions, at the expense of not
369 shedding much light on the etiology of the disease. On the other end, SigMod and
370 dmGWAS gravitate towards larger, less associated solutions which provide a wide
371 overview of the biological context. While this deepens our understanding of the disease
372 and provides biological hypotheses, they require further analyses. For instance, a user
373 might need to examine the centrality of the selected genes, and discern the extent to
374 which that property is driving the selection of each gene. HotNet2 balances both
375 approaches at the expense of producing the largest solution: a constellation of many,
376 highly associated, small subnetworks. Additionally, all solutions share two drawbacks.
377 First, they are all equally bad at discriminating cases from controls. Yet, the
378 classification accuracy of network methods is similar to that of a classifier trained on
379 the entire genome, which suggests that cases and controls are difficult to separate in the
380 GENESIS dataset. This may be due to limited statistical power, which reduces the
381 ability to identify relevant SNPs; but in any event, we do not expect to be able to
382 separate people who have or will developed cancer from others on the sole basis of their
383 genomes, ignoring all environmental factors and chance events. Second, all methods are
384 remarkably unstable, yielding different solutions for slightly different inputs. This might
385 partly be caused by the instability of the P-values themselves in low statistical power
386 settings [37]. Hence, heinz’s conservative transformation of P-values, which favors only
387 the most extreme ones, leads to improved stability. Another source of instability might
388 be the redundancy inherent to biological networks, since they are subject to an
389 evolutionary pressure to avoid single points of failure [38]. Hence, biological networks
390 will often have multiple paths connecting two high-scoring nodes.
391

To overcome these limitations while exploiting the strengths of the individual
392 methods, we proposed combining them into a consensus solution. We use a
393 straightforward strategy of including any node that was recovered by multiple methods.
394

We thus proposed two networks: a consensus solution, meant to address the heterogeneity of the solutions in the full dataset, and a stable consensus solution, which in addition addressed the instability of the methods. They both included the majority of the strongly associated smaller solutions and captured genes and broader mechanisms related to cancer, thus synthetizing 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 are related to cancer at the expression level and interact with genes with consistently high association scores. Importantly, while the consensus approach was as unstable as the individual network-guided methods, the stable consensus network retained the ability to provide compelling hypotheses and had better stability. This supports that instability might be caused by redundant but equivalent biological mechanisms, and hence validates 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 [19]. Network methods show modest precision, but much higher recall at recovering BCAC hits (Section 2.4). While precision might be desirable when a subset of good markers is required (for instance, for diagnosis), higher recall is desirable in exploratory settings. Nonetheless, BCAC is not an ideal ground truth. First, the studied populations are non-overlapping: BCAC focuses on a pan-European cohort, while GENESIS targets the French population specifically. Second, the study designs differ: 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 focuses on the French high-risk population attending the family cancer clinics. Despite these differences, we expect some degree of shared genetic architecture, especially at the gene level. 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 low 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 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.

The strength of network-based analyses comes from leveraging prior knowledge to boost discovery. In consequence, they show their shortcomings on understudied genes, especially those not in the network. Out of the 32 767 genes to which we can map the genotyped SNPs, 60.7% (19 887) are not in the PPIN. The majority of those (14 660) are non-coding genes, mainly lncRNA, miRNA, and snRNA (S7 Fig). Yet, RNA genes like *CASC16* are associated to breast cancer (Section 2.1), reminding us of the importance of using networks beyond coding genes. In addition, even protein-coding genes linked to breast cancer susceptibility [33], like *NEK10* (P -value 1.6×10^{-5} , overlapping with *SLC4A7*) or *POU5F1B*, were absent from the PPIN. However, on average protein-coding genes absent from the PPIN are less associated with breast cancer susceptibility (Wilcoxon rank-sum P -value = 2.79×10^{-8} , median P -values of 0.43 and 0.47). This cannot be due to well-known genes having more known interactions because we are only using interactions from high-throughput experiments. As disease genes tend to be more central [9], we hypothesize that it is due to interactions between central genes being more likely. It is worth noting that network methods that do not use PPIs, like SConES GS and GM, did recover SNPs in *NEK10* and *CASC16*. Moreover, both SConES GM and GI recovered intergenic regions, which might contain key regulatory elements [39] and, yet, 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, and other strategies (e.g. eQTLs, SNPs associated to the expression of a gene) 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 [40, 41], and they might bias discovery by reducing the average shortest path length between two random nodes. On the other hand, Huang et al. [10] 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), 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 (S8 Fig). 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 methods is the computation of the gene score. In this work we used VEGAS2 [42] 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 [43]), and other algorithms like PEGASUS [43], SKAT [44] or COMBAT [45] might have more statistical power.

How to handle linkage disequilibrium (LD) is often a concern among GWAS practitioners. 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. With regards to SConES, fewer SNPs would lead to simpler SNP networks and, possibly, shorter runtimes. However, as mentioned in Section 2.6, LD patterns seem to drive SConES’ solutions, and an LD pruning step could potentially alter them. In Section 2.3 we highlight ambiguities that appear when genes overlap or are in LD. In fact, the presented case is paradigmatic, since all three genes are located in the HLA region, the most gene-dense region of the genome [46]. 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. LD also affects SConES, since it penalizes selecting a SNP and not its neighbors, via a nonzero parameter η in Equation 5. Due to LD, nearby SNPs’ P-values are correlated; and since SNP networks are determined by positional information, nearby SNPs are likely to be linked. Hence, SConES will tend to select LD-blocks formed by low P-value SNPs. This might explain why SConES produces similar results on the GS and GM networks, heavily affected by LD (Section 2.6). 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. We hypothesize that this is caused by the absence of joint association of a gene and a majority of its neighbors. This is supported by the lack of results from LEAN as well. Yet, a different combination of parameters could lead to a more informative SConES’ solution (e.g. a lower λ in Equation 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 under studied — hindering the

accuracy with which the GS network captures the genome structure. A stringent LD pruning might address such problems.

To produce the two consensus solutions, we had to face practical challenges due to the differences in interfaces, preprocessing steps, and unexpected behaviors of the various methods. To make it easier for other to apply these methods to new datasets and aggregate their solutions, we built six nextflow pipelines [47] with a consistent interface and, whenever possible, parallelized computation. They are available on GitHub: <https://github.com/hclimente/gwas-tools> (Section 4.1). Importantly, those methods that had a permissive license were compiled into a Docker image for easier use, which is available on Docker Hub hclimente/gwas-tools.

4 Materials and methods

4.1 GENESIS dataset, preprocessing and quality control

The GENE Sisters (GENESIS) study was designed to investigate risk factors for familial breast cancer in the French population [18]. Index cases are patients with infiltrating mammary or ductal adenocarcinoma, who had a sister with breast cancer, and who have been tested negative for *BRCA1* and *BRCA2* pathogenic variants. Controls are unaffected colleagues and/or friends of the cases, born around the year of birth of their corresponding case (± 3 years). We focused on the 2 577 samples of European ancestry, of which 1 279 were controls and 1 298 were cases. The genotyping was performed using the iCOGS array, a custom Illumina array designed to study the genetic susceptibility to hormone-related cancers [48]. It contains 211 155 SNPs, including SNPs putatively associated with breast, ovarian, and prostate cancers, SNPs associated with survival after diagnosis, and SNPs associated to other cancer-related traits, as well as candidate functional variants in selected genes and pathways.

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 subset of 20 duplicated SNPs in *FGFR2* were also removed. In addition, we removed the samples with more than 10% missing genotypes. After controlling for relatedness, 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 SNP- and gene-based GWAS

To measure association between a genotype and susceptibility to breast cancer, we performed a per-SNP 1 d.f. χ^2 allelic test using PLINK v1.90 [49]. In order to obtain significant SNPs, we performed Bonferroni correction, to keep family-wise error rate below 5%. The threshold used was $\frac{0.05}{197083} = 2.54 \times 10^{-7}$.

Then, we used VEGAS2 [42] to compute the gene-level association score from the P-values of the SNPs mapped to them. 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. For each gene, we computed VEGAS2 scores using only the 10% of SNPs with lowest P-values among all those that were mapped to it. We used the 62 193 genes described in GENCODE 31 [50], 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 also used 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.

Nodes can be described by properties provided by the topology of the graph. 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.

We also use two matrices that describe two different properties of a graph. These matrices are square, and have as many rows and columns as nodes in the network. The element (i, j) hence 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; 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 mathematical formulations of the different network methods are compatible with any type of network (protein interactions, gene coexpression, regulatory, etc.). Here, we used protein-protein interaction networks (PPIN) for all of them except SConES, as PPINs are interpretable, well-characterized, and the methods were designed to run efficiently on them. 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 PPIN. 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).

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

SNP networks SConES [16] is the only network method designed to handle SNP networks. As in gene networks, two SNPs are connected in a SNP network when there is evidence of shared functionality between two SNPs. Azencott et al. [16] proposed three ways of building these networks: connecting the SNPs consecutive in the genomic sequence (“GS network”); interconnecting all the SNPs mapped to the same gene, on top of GS (“GM network”); and interconnecting all SNPs mapped to two genes for which a protein-protein interaction exists, on top of GM (“GI network”). We focused on the GI network, as it fits the scope of this work better, using the PPIN described above. However, at different stages we also compared GI to GS and GM to understand how considering the PPIN affects SConES’ output. For the GM network, we used the mapping described in Section 4.3.5. In all three the node scores are the association scores of the individual SNPs with the phenotype (1 d.f. χ^2 statistic). The properties of these three subnetworks are available in S1 Table.

4.3.3 High-score subnetwork search algorithms

Genes that contribute to the same function are nearby in the PPIN, and can be topologically related to each other in diverse ways (densely interconnected modules, nodes around a hub, a path, etc.). Several aspects have to be taken into consideration when developing a network method: how to score the nodes, whether the affected mechanisms form a single connected component or several, how to frame the problem in a computationally efficient fashion, which network to use, etc. Unsurprisingly, multiple solutions have been proposed. We examined six of them: five that explore the PPIN, and one which explores SNP networks. We selected methods that were open-source, and had an implementation available and accessible documentation. Their main differences are summarized in Table 2. We scored both SNPs and genes with the P-values (or transformations of them) computed in Section 4.2.

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	$-\log_{10}(P)$	[13]
heinz	Omics	Genes	Yes	-	1	Summary	BUM	[14]
HotNet2	Omics	Genes	Yes	Module	≥ 1	Summary	Local FDR	[15]
LEAN	Omics	Genes	Yes	Star	≥ 1	Summary	$-\log_{10}(P)$	[12]
SConES	GWAS	SNPs	Yes	Module	≥ 1	Genotypes	1 d.f. χ^2	[16]
SigMod	GWAS	Genes	Yes	Module	≥ 1	Summary	$\Phi^{-1}(1 - P)$	[17]

Field: field in which the algorithm was developed. **Nodes:** the type of nodes in the network, either genes (PPIN) or SNPs. **Exhaustive:** whether all the possible solutions given the selected hyperparameters are explored. **Solution:** additional properties are enforced on the solution, 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.

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

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

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

where k is the number of genes in S_i ; z_j is the Z score of gene j , computed as $\phi^{-1}(1 - P\text{-value}_j)$; and ϕ^{-1} is the inverse normal distribution function. 622
623

3. Identify neighboring nodes of S_i , i.e. nodes at distance $\leq d$. 624
4. Add the neighboring nodes whose inclusion increases Z_{m+1} by more than a threshold $Z_m \times (1 + r)$. 625
626
5. Repeat 2-4 until no further enlargement is possible. 627
6. Add S_i to the list of subnetworks to return. Normalize its Z-score as 628

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

where $Z_m(\pi)$ represents a vector containing 100 000 random subsets of the same number of genes. 629
630

DmGWAS carries out this process on every gene in the PPIN. We used the implementation of dmGWAS in the dmGWAS 3.0 R package [51]. Unless otherwise specified, we used the suggested hyperparameters $d = 2$ and $r = 0.1$. We used the function `simpleChoose` to select the solution, which aggregates the top 1% subnetworks. 631
632
633
634
635

heinz The goal of heinz is to identify the highest-scored connected subnetwork [14]. The authors propose a transformation of the genes' P-value into a score that is negative under weak association with the phenotype, and positive under a strong one. This transformation is achieved by 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: 636
637
638
639
640
641
642
643
644

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

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. [52] which, in practice is often fast and optimal, although neither is guaranteed. We used BioNet's implementation of heinz [53, 54]. 645
646
647
648
649

HotNet2 HotNet2 was developed to find connected subgraphs of genes frequently mutated in cancer [15]. To that end, it considers both the local topology of the PPIN and the scores of the nodes. The former is captured by an insulated heat diffusion process: 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 a stationary state is reached, in which the temperature of 650
651
652
653
654
655
656

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
657
658

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

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 can be interpreted as a (non-symmetric) similarity between those
two nodes. To obtain densely connected solutions, HotNet2 prunes E , only
preserving edges such that $w(E) > \delta$. Lastly, HotNet2 evaluates the statistical
significance of the solutions by comparing their size to the size of PPINs obtained
by permuting the node scores. We assigned initial node scores as in Nakka et
al. [43], assigning a score of 0 for the genes with a low probability of being
associated to the disease, and $-\log_{10}(\text{P-value})$ to those likely to be. In the
GENESIS dataset, the threshold separating both was a P-value of 0.125, which
was obtained using a local FDR approach [55]. 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 [15]. HotNet2 is
implemented in Python [56].
659
660
661
662
663
664
665
666
667
668
669
670
671
672

LEAN LEAN searches altered “star” subnetworks, that is, subnetworks composed by
one central node and all its interactors [12]. By imposing this restriction, LEAN
can exhaustively test all such subnetworks (one per node). For a particular
subnetwork of size m , the P-values corresponding to the involved nodes are ranked
as $p_1 \leq \dots \leq p_m$. Then, k binomial tests are conducted, to compute the
probability of having k out of m P-values lower or equal to p_k under the null
hypothesis. The minimum of these k P-values is the score of the subnetwork. This
score is transformed into a P-value through an empirical distribution obtained via
a subsampling scheme, where gene sets of the same size are selected randomly,
and their score computed. Lastly, P-values are corrected for multiple testing
through a Benjamini-Hochberg correction. We used the implementation of LEAN
from the LEANR R package [57].
673
674
675
676
677
678
679
680
681
682
683
684

SConES SConES searches the minimal, modular, and maximally associated
subnetwork in a SNP graph [16]. Specifically, it solves the problem
685
686

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

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. Given a λ and an η ,
the aforementioned problem has a unique solution, that SConES finds using a
graph min-cut procedure. As in Azencott et al. [16], 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 [58].
687
688
689
690
691
692
693
694
695
696

SigMod SigMod aims at identifying the highest-scoring, most densely connected
subnetwork [17]. It addresses an optimization problem similar to that of SConES
697
698

(Equation 5), but the connectivity term encourages connected solutions by favoring solutions where many edges connect two selected nodes, rather than penalizing disconnected ones.

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

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 PPINs. Second, it favors solutions containing many edges. 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 candidate solution with the highest standardized score and that is not larger than a user-specified threshold (`nmax`) is the chosen solution. SigMod is implemented in an R package [59].

Consensus We built a consensus solution by retaining the genes that were selected by at least two of the six methods (using SConES GI for SConES). Any edge between the selected genes in the PPIN was included.

We performed all the computations in the cluster described in Section 4.6.

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 $-\log_{10}(\text{P-value})$: 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 (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 740
In multiple steps of this article, we needed to compare the outcome of a method that 741
works on genes with the outcome of another method that works on SNPs. For this 742
purpose, we used the SNP-gene correspondence described in Section 4.2. To convert a 743
list of SNPs into a list of genes, we included all the genes that can be mapped to any of 744
the SNPs. Conversely, to convert a list of genes into a list of SNPs, we included all the 745
SNPs that can be mapped to any of the genes. 746

4.4 Pathway enrichment analysis 747
We searched for pathways enriched in the gene solutions produced by the above 748
methods. We conducted a hypergeometric test on pathways from Reactome [60] using 749
the function `enrichPathway` from the ReactomePA R package [61]. The universe of 750
genes included any gene that we could map to a SNP in the iCOGS array (Section 4.2). 751
We adjusted the P-values for multiple testing as in Benjamini and Hochberg [62] (BH). 752
Pathways with a BH adjusted P-value < 0.05 were deemed significant. 753

4.5 Benchmark of methods 754
We evaluated multiple properties (described below) of the different methods through a 755
5-fold subsampling setting. We applied each method to 5 random subsets of the original 756
GENESIS dataset containing 80% of the samples (*train set*). When pertinent, we 757
evaluated the solution on the remaining 20% (*test set*). We used the 5 repetitions to 758
estimate the average and the standard deviation of the different measures. Every 759
method and repetition ran on the same computational settings (Section 4.6). 760

4.5.1 Properties of the solution 761
We compared the runtime, the number of selected features (genes or SNPs), and the 762
stability (sensitivity of the result to small changes in the input, here, using different 763
train sets) of the different network methods. Nogueira and Brown [63] proposed 764
quantifying the stability of a method using the Pearson correlation between the genes 765
selected on different subsamples. This correlation was calculated between vectors with 766
the length of the total number of features, containing a 0 at position i if feature i was 767
not selected, and 1 if it was. 768

4.5.2 Classification accuracy of selected SNPs 769
A desirable solution offers good predictive power on unseen *test* samples. We evaluated 770
the predicting power of the SNPs selected by the different methods through the 771
performance of an L1-penalized logistic regression classifier, which searches for a small 772
subset of SNPs which provides good classification accuracy. We trained the classifier 773
exclusively on those selected SNPs to predict the outcome (case/control). The L1 774
penalty helps to account for linkage disequilibrium by reducing the number of SNPs 775
included in the model (*active set*) while improving the generalization of the classifier. 776
This penalty was set by cross-validation, choosing the value that minimized 777
misclassification error. We applied each network method to each *train* set, and trained 778
the classifier on the same train set using only the selected SNPs. When the method 779
retrieved a list of genes (all of them except SConES), we considered as selected all the 780
SNPs mapped to any of those genes. Then we evaluated the sensitivity and the 781
specificity on the *test set*. The active set gave an estimate of a plausible, more sparse 782
solution with comparable predictive power to the original solution. To obtain a baseline, 783
we also trained the classifier on all the SNPs. We do not expect a linear model on 784

selected SNPs to be able to separate cases from controls well. Indeed, the lifetime
785
cumulative incidence of breast cancer among women with a family history of breast or
786
ovarian cancer, and no *BRCA1/2* mutations, is only 3.9 times more than in the general
787
population [64]. However, classification accuracy may be one additional informative
788
criterion on which to evaluate solutions.
789

4.5.3 Comparison to state-of-the-art

An alternative way to evaluate the results is by comparing our results to an external
791
dataset. For that purpose, we used the 153 genes associated to familial breast cancer on
792
DisGeNET [65]. Across this article we refer to these genes as *breast cancer susceptibility*
793
genes.
794

Additionally, we used the summary statistics from the Breast Cancer Association
795
Consortium (BCAC), a meta-analysis of case-control studies conducted in multiple
796
countries which included 13 250 641 SNPs genotyped or imputed on 228 951 women of
797
European ancestry mostly from the general population [19]. Through imputation,
798
BCAC includes more SNPs than the iCOGS array used for GENESIS (Section 4.1). Yet,
799
in all the comparisons in this paper, we focused on the subset of the GENESIS SNPs
800
that passed quality control (Section 4.1). Hence, we used the same Bonferroni threshold
801
as in Section 4.2 to determine the significant SNPs in BCAC. We also computed
802
gene-scores in the BCAC data using VEGAS2, as in Section 4.1. In this case, we did use
803
the summary statistics of all 13 250 641 available SNPs, and the genotypes from
804
European samples from the 1000 Genomes Project [66] to compute the LD patterns.
805
Since these genotypes did not include chromosome X, we excluded it from this analysis.
806
All comparisons included only the genes common to GENESIS and BCAC, so we used
807
the corresponding Bonferroni threshold (1.66×10^{-6}) to call gene significance.
808

4.5.4 Network rewirings

Rewiring the PPIN while preserving the number of edges of each gene allows to study
810
the impact of the topology on the output of network methods. Indeed, the edges lose
811
their biological meaning but the topology of the network is conserved. We produced 100
812
such rewirings by swapping edges in the PPIN. We scored the genes as described in
813
section 4.3.3. We only applied only four methods on the rewirings: heinz, dmGWAS,
814
LEAN and SigMod. We excluded HotNet2 and SConES, since they take notably longer
815
to run than the other methods, while taking up more computational resources.
816

4.6 Computational resources

We ran all the computations on a Slurm cluster, running Ubuntu 16.04.2 on the nodes.
818
The CPU models on the nodes were Intel Xeon CPU E5-2450 v2 at 2.50GHz and Intel
819
Xeon E5-2440 at 2.40GHz. The nodes running heinz and HotNet2 had 20GB of
820
memory; the ones running dmGWAS, LEAN, SConES, and SigMod, 60GB. For the
821
benchmark (Section 4.5), we ran each of the methods on the same Ubuntu 16.04.2 node,
822
with a CPU Intel Xeon E5-2450 v2 at 2.50GHz, and 60GB of memory.
823

4.7 Code and data availability

We developed computational pipelines for several steps of GWAS analyses, such as
825
physically mapping SNPs to genes, computing gene scores, and performing six different
826
network analyses. For each of those processes, we created a pipeline with a clear
827
interface that should work on any GWAS dataset. They are compiled in
828
<https://github.com/hclimente/gwas-tools>. Although the GENESIS data is not
829

public, the code to apply the pipelines to this data, as well as the code that reproduces all the analyses in this article are available at
830
831
<https://github.com/hclimente/genewa>. We deposited all the produced gene
832
solutions on NDEx (<http://www.ndexbio.org>), under the UUID
833
e9b0e22a-e9b0-11e9-bb65-0ac135e8bacf.
834

Summary statistics for SNPs and genes are available
835
at <https://github.com/hclimente/genewa>. We cannot share genotype data publicly
836
for confidentiality reasons, but are available from GENESIS. Interested researchers can
837
contact nadine.andrieu(at)curie.fr.
838

5 Supporting information

S1 Table. Summary statistics on the results of SConES on the three SNP-SNP interaction networks. The first row within each block contains the summary statistics on the whole network.
839
840
841
842

S2 Table. Summary statistics on the results of multiple network methods on the PPIN. The first row contains the summary statistics on the whole PPIN.
843
844

S3 Table. Pathway enrichment analyses of the genes in SigMod solution.
845

S4 Table. Pathway enrichment analyses of the genes in dmGWAS solution.
846
847

S5 Table. Pathway enrichment analyses of the genes in HotNet2 solution.
848

S6 Table. Pathway enrichment analyses of the genes in the consensus solution.
849
850

S1 Fig. 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.
851
852
853

S2 Fig. Association in GENESIS. The red line represents the Bonferroni threshold. (A) SNP association, measured from the outcome of a 1 d.f. χ^2 allelic test (Section 4.2). 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 [42] using the 10% of SNPs with the lowest P-values (Section 4.2). 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, $\lambda = 0.03$).
854
855
856
857
858
859
860
861

S3 Fig. 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.
862
863
864

S4 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.
865
866
867

S5 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.

868
869
870
871
872
873
874

S6 Fig. Parameter space of the network methods. (A) Boxplot of the solution sizes of the methods under the explored parameters (Section 4.3.4). (B) Size of SConES's with regards to each pair of parameters. (C) Pearson correlation between the solutions of the different runs.

875
876
877
878

S7 Fig. Biotypes of genes from the annotation that are not present in the HINT PPIN.

879
880

S8 Fig. 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-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 the stability of the solution, but increases the solution size, has mixed effects on the sensitivity and specificity, and impacts negatively the required runtime of the algorithms.

881
882
883
884
885
886
887
888
889
890
891
892
893
894

S9 Fig. 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); in black, the method's solution. 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.

895
896
897
898
899
900

Acknowledgments

901

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 centre (S. Mazoyer, F. Damiola, L. Barjhoux, C. Verny-Pierre, V. Sornin). We wish to pay a tribute to Olga M. Sinilnikova, who was one of the initiators and principal investigators of GENESIS and who died prematurely on June 30, 2014.

902
903
904
905
906
907
908
909
910
911
912

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. Manouvier-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. Gesta; 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. Feille; 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	913 914 915 916 917 918 919 920 921 922 923 924 925 926 927 928 929 930 931 932 933 934 935 936 937 938
---	--

Author contributions

Conceptualization Héctor Climente-González, Christine Lonjou, Chloé-Agathe Azencott.	939 940 941
Data curation Christine Lonjou, GENESIS Study collaborators.	942
Formal Analysis Héctor Climente-González, Christine Lonjou.	943
Funding acquisition Dominique Stoppa-Lyonnet, Nadine Andrieu, Chloé-Agathe Azencott.	944 945
Investigation Héctor Climente-González, Christine Lonjou.	946
Methodology Héctor Climente-González, Christine Lonjou, Chloé-Agathe Azencott.	947
Project administration Chloé-Agathe Azencott.	948
Resources GENESIS Study collaborators, Dominique Stoppa-Lyonnet, Nadine Andrieu.	949 950
Software Héctor Climente-González, Christine Lonjou.	951
Supervision Christine Lonjou, Fabienne Lesueur, Nadine Andrieu, Chloé-Agathe Azencott.	952 953
Validation Christine Lonjou, Fabienne Lesueur.	954
Visualization Héctor Climente-González.	955

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.semcan.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 Omnipathogenic. 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. 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.
9. 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.
10. 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.
11. 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.
12. 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.

13. 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.
14. 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.
15. 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.
16. 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.
17. 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.
18. 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.
19. 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.
20. 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.
21. 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.
22. 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.
23. 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.
24. 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.

25. 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.
26. 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.
27. 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.
28. 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.
29. 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.
30. 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.
31. 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.
32. 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.
33. 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.
34. 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.
35. 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.
36. 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.
37. 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.

38. 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.
39. 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.
40. Cai JJ, Borenstein E, Petrov DA. Broker Genes in Human Disease. *Genome Biology and Evolution*. 2010;2:815–825. doi:10.1093/gbe/evq064.
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. 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.
43. 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.
44. 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.
45. 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.
46. 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.
47. 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.
48. 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.
49. 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.
50. 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.
51. Wang Q, Jia P. dmGWAS 3.0; 2014. <https://bioinfo.uth.edu/dmGWAS/>.
52. 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.

53. 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.
54. Dittrich M, Beisser D. BioNet; 2008. <https://bioconductor.org/packages/BioNet/>.
55. 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.
56. Leiserson MDM, Vandin F, Wu HT, Dobson JR, Eldridge JV, Thomas JL, et al.. HotNet2; 2018. <https://github.com/raphael-group/hotnet2>.
57. Gwinner F. LEANR; 2016. <https://cran.r-project.org/web/packages/LEANR/>.
58. Clemente-González H, Azencott CA. martini; 2019. <https://www.bioconductor.org/packages/martini/>.
59. Liu Y. SigMod v2; 2018. <https://github.com/YuanlongLiu/SigMod>.
60. 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.
61. 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.
62. 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.
63. 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.
64. 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.
65. 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.
66. 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.