

GigaScience

A multi-omics data simulator for complex disease studies and its application to evaluate multi-omics data analysis methods for disease classification

--Manuscript Draft--

Manuscript Number:	GIGA-D-18-00397	
Full Title:	A multi-omics data simulator for complex disease studies and its application to evaluate multi-omics data analysis methods for disease classification	
Article Type:	Research	
Funding Information:	Ministry of Science and Technology, Taiwan (MOST 106-2221-E-400-005-MY3)	Dr. Ren-Hua Chung
Abstract:	<p>Background</p> <p>An integrative multi-omics analysis approach that combines multiple types of omics data including genomics, epigenomics, transcriptomics, proteomics, metabolomics, and microbiomics, has become increasingly popular for understanding the pathophysiology of complex diseases. Although many multi-omics analysis methods have been developed for complex disease studies, there is no simulation tool that simulates multiple types of omics data and models their relationships with disease status. Without such a tool, it is difficult to evaluate the multi-omics analysis methods on the same scale and to estimate the sample size or power when planning a new multi-omics disease study.</p> <p>Results</p> <p>We developed a multi-omics data simulator OmicsSIMLA, which simulates genomics (i.e., SNPs and copy number variations), epigenomics (i.e., whole-genome bisulphite sequencing), transcriptomics (i.e., RNA-seq), and proteomics (i.e., normalized reverse phase protein array) data at the whole-genome level. Furthermore, the relationships between different types of omics data, such as meQTLs (SNPs influencing methylation), eQTLs (SNPs influencing gene expression), and eQTM (methylation influencing gene expression), were modeled. More importantly, the relationships between these multi-omics data and the disease status were modeled as well. We used OmicsSIMLA to simulate a multi-omics dataset for breast cancer under a hypothetical disease model, and used the data to compare the performance among existing multi-omics analysis methods in terms of disease classification accuracy and run time.</p> <p>Conclusions</p> <p>Our results demonstrated that complex disease mechanisms can be simulated by OmicsSIMLA, and a random forest-based method showed the highest prediction accuracy when the multi-omics data were properly normalized. OmicsSIMLA can be downloaded at https://omicssimla.sourceforge.io.</p>	
Corresponding Author:	Ren-Hua Chung National Health Research Institutes Zhunan, Miaoli TAIWAN	
Corresponding Author Secondary Information:		
Corresponding Author's Institution:	National Health Research Institutes	
Corresponding Author's Secondary Institution:		
First Author:	Ren-Hua Chung	
First Author Secondary Information:		
Order of Authors:	Ren-Hua Chung	

	Chen-Yu Kang
Order of Authors Secondary Information:	
Additional Information:	
Question	Response
Are you submitting this manuscript to a special series or article collection?	No
<p>Experimental design and statistics</p> <p>Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our Minimum Standards Reporting Checklist. Information essential to interpreting the data presented should be made available in the figure legends.</p> <p>Have you included all the information requested in your manuscript?</p>	Yes
<p>Resources</p> <p>A description of all resources used, including antibodies, cell lines, animals and software tools, with enough information to allow them to be uniquely identified, should be included in the Methods section. Authors are strongly encouraged to cite Research Resource Identifiers (RRIDs) for antibodies, model organisms and tools, where possible.</p> <p>Have you included the information requested as detailed in our Minimum Standards Reporting Checklist?</p>	Yes
<p>Availability of data and materials</p> <p>All datasets and code on which the conclusions of the paper rely must be either included in your submission or deposited in publicly available repositories (where available and ethically appropriate), referencing such data using a unique identifier in the references and in</p>	Yes

the “Availability of Data and Materials” section of your manuscript.

Have you have met the above requirement as detailed in our [Minimum Standards Reporting Checklist?](#)

1 **A multi-omics data simulator for complex disease studies and its**
2
3
4 **application to evaluate multi-omics data analysis methods for disease**
5
6
7 **classification**
8
9

10 Ren-Hua Chung^{1*}, Chen-Yu Kang¹
11

12
13 ¹Division of Biostatistics and Bioinformatics, Institute of Population Health Sciences,
14
15

16 National Health Research Institutes, Zhunan, Taiwan
17
18
19
20
21
22

23 *Corresponding author: Ren-Hua Chung, PhD
24
25

26 Address: No 35, Keyan Road, Zhunan, Miaoli, 350, Taiwan
27
28

29 Tel: 886-37-246-166 #36105
30
31

32 Fax: 886-37-586-467
33
34
35

36 Email: rchung@nhri.org.tw
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

Abstract

Background

An integrative multi-omics analysis approach that combines multiple types of omics data including genomics, epigenomics, transcriptomics, proteomics, metabolomics, and microbiomics, has become increasingly popular for understanding the pathophysiology of complex diseases. Although many multi-omics analysis methods have been developed for complex disease studies, there is no simulation tool that simulates multiple types of omics data and models their relationships with disease status. Without such a tool, it is difficult to evaluate the multi-omics analysis methods on the same scale and to estimate the sample size or power when planning a new multi-omics disease study.

Results

We developed a multi-omics data simulator OmicsSIMLA, which simulates genomics (i.e., SNPs and copy number variations), epigenomics (i.e., whole-genome bisulphite sequencing), transcriptomics (i.e., RNA-seq), and proteomics (i.e., normalized reverse phase protein array) data at the whole-genome level. Furthermore, the relationships between different types of omics data, such as meQTLs (SNPs influencing methylation), eQTLs (SNPs influencing gene expression), and eQTM (methylation influencing gene expression), were modeled. More importantly, the relationships between these multi-omics data and the disease status were modeled as well. We used OmicsSIMLA to simulate a multi-omics dataset for breast cancer

1 under a hypothetical disease model, and used the data to compare the performance among
2
3
4 existing multi-omics analysis methods in terms of disease classification accuracy and run
5
6
7 time.
8
9

10 **Conclusions**

11
12 Our results demonstrated that complex disease mechanisms can be simulated by
13
14 OmicsSIMLA, and a random forest-based method showed the highest prediction accuracy
15
16
17 when the multi-omics data were properly normalized. OmicsSIMLA can be downloaded at
18
19
20 <https://omicssimla.sourceforge.io>.
21
22
23
24
25
26
27
28

29 **Keywords**

30
31
32 Multi-omics data, complex disease study, simulation tool
33
34
35
36
37
38

39 **Introduction**

40
41
42 Complex diseases such as hypertension, type 2 diabetes, and autism are caused by multiple
43
44 genetic and environmental factors [1]. Genome-wide association studies have identified many
45
46 genetic variants (i.e., SNPs) associated with the complex diseases. However, it remains
47
48
49 difficult to understand the roles of the associated SNPs in the molecular pathophysiology of
50
51
52 the disease and how the SNPs interact with other SNPs in a biological network [2]. With the
53
54
55 advancement of high-throughput sequencing technology such as next-generation sequencing
56
57
58
59
60
61
62
63
64
65

1 (NGS) and massive parallel technology such as mass spectrometry, multiple types of omics
2
3
4 data (i.e., multi-omics data) including genomics, epigenomics, transcriptomics, proteomics,
5
6
7 metabolomics, and microbiomics are rapidly generated [3]. As a single type of data generally
8
9
10 cannot capture the complexity of molecular events causing the disease, an integrative
11
12
13 approach to combining the multi-omics data would be ideal to help elucidate the
14
15
16 pathophysiology of the disease [2].
17
18
19
20
21
22

23 Integrative methods to combine multi-omics data for disease studies have been developed
24
25
26 rapidly [4-8]. They can be generally classified into two categories: multi-staged and meta-
27
28
29 dimensional approaches [9]. The multi-staged approach aims to first identify relationships
30
31
32 between the multi-omics data, and then test the associations between the multi-omics data
33
34
35 and the phenotype. For example, Jennings et al. [7] constructed a Bayesian hierarchical
36
37
38 model consisting of two stages. The first stage partitioned gene expression into factors
39
40
41 accounted by methylation, copy number variation (CNV), and other unknown causes. These
42
43
44 factors were subsequently used as predictors for clinical outcomes in the second stage model.
45
46
47 One advantage of this approach is that the causal relationships between multi-omics data can
48
49
50 be modeled. In contrast, the meta-dimensional approach combines the multi-omics data
51
52
53 simultaneously. Raw or the transformed data from the multi-omics data are combined into a
54
55
56 single matrix for the analysis. This approach allows for a more flexible inference of the
57
58
59
60
61
62
63
64
65

1 relationships among the multi-omics data, without the assumptions of the causal relationships
2
3
4 between these data.
5
6
7
8
9

10 Although many multi-omics analysis methods for disease studies are available, they were
11
12 generally evaluated by simulations with data generated specifically to the methods. To
13
14 compare the performance among these methods, it is necessary to use the same simulated
15
16 multi-omics dataset with disease status. However, current simulation tools for disease studies
17
18 mainly focused on simulating a certain type of omics data. For example, more than 25
19
20 simulators are available for simulating genetic data with phenotypic trait, according to the
21
22 Genetic Simulation Resources website (<https://popmodels.cancercontrol.cancer.gov/gsr/>).
23
24
25 Tools such as WGBSSuite [10] and pWGBSSimla [11] can simulate whole-genome
26
27 bisulphite sequencing (WGBS) data in case-control samples. Moreover, tools such as
28
29 Polyester [12] and SimSeq [13] simulate RNA-seq data with differential gene expression
30
31 between two groups of samples. To our knowledge, there is currently no simulation tool that
32
33 is capable of simulating a variety of omics data types and modeling the complex relationships
34
35 between the data and the disease. Furthermore, sample size estimation when planning a
36
37 multi-omics study to ensure sufficient power also becomes important [3]. This also requires a
38
39 simulation tool that simulates realistic multi-omics data structures and models the
40
41 architecture of the complex disease.
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1
2
3
4 Here, we developed the multi-omics data simulator OmicsSIMLA, which simulates genomics
5
6
7 data including SNPs and CNVs, epigenomics data such as the WGBS data, transcriptomics
8
9
10 data (i.e., RNA-seq), and proteomics data such as the normalized reverse phase protein array
11
12
13 (RPPA) data at a whole-genome level. Furthermore, the relationships between different types
14
15
16 of omics data, such as meQTLs (SNPs influencing methylation), eQTLs (SNPs influencing
17
18
19 gene expression), and eQTM (methylation influencing gene expression), were modeled. More
20
21
22 importantly, the relationships between these multi-omics data and disease status were
23
24
25 modeled as well. The disease models in OmicsSIMLA are flexible so that the main effects
26
27
28 and/or interaction effects (either risk or protective) of SNPs and CNVs on the disease can be
29
30
31 specified. Differential methylation and differential gene and protein expression between cases
32
33
34 and controls can also be simulated. We demonstrated the usefulness of OmicsSIMLA by
35
36
37 simulating a multi-omics dataset for breast cancer under a hypothetical disease model, and
38
39
40 compared the performance among existing multi-omics analysis tools based on the data.
41
42
43
44
45
46
47

48 **Results**

49
50
51 Figure 1 shows the framework of OmicsSIMLA. The genomics data that can be simulated
52
53
54 include SNPs and CNVs. Genotypes at SNPs in unrelated and/or family samples are
55
56
57 simulated based on the SeqSIMLA2 algorithm [14]. CNV status (i.e., a deletion, normal, one
58
59
60

1 duplication or two duplications) on a chromosome is simulated based on the user-specified
2
3
4 chromosomal regions and CNV frequencies. Affection status of each sample is determined by
5
6
7 a logistic penetrance function conditional on the causal SNPs and CNVs, and/or the
8
9
10 interactions among the causal SNPs. The epigenomics data are the methylated and total read
11
12
13 counts at CpGs based on bisulphite sequencing, simulated using the pWGBSSimla algorithm
14
15
16 incorporating methylation profiles for 29 human cell and tissue types [11]. Allele-specific
17
18
19 methylation (ASM), in which paternal and maternal alleles have different methylation rates,
20
21
22 and differentially methylated region (DMR), where the same CpGs in the region have
23
24
25 different methylation rates among different cell types, can also be simulated. Furthermore, the
26
27
28 transcriptomics data (i.e., RNA-seq read counts) are simulated with a parametric model
29
30
31 assuming a negative-binomial distribution. Finally, the mass-action kinetic action model [15]
32
33
34 is used to simulate proteomics data at a certain time point incorporating the gene expression
35
36
37 data. Some SNPs can be specified as meQTLs and eQTLs, and some CpGs can be specified
38
39
40 as eQTM. Allele-specific expression (ASE), which alleles in a gene have different expression
41
42
43 levels, caused by cis-eQTL can also be simulated. The differential methylation, gene
44
45
46 expression, and protein expression levels between cases and controls are simulated
47
48
49 conditional on the affection status.
50
51
52
53
54
55
56
57

58 Using OmicsSIMLA, we simulated a multi-omics dataset based on hypothetical pathways for
59
60
61
62
63
64
65

1 breast cancer as described in Ritchie et al. [9] and illustrated in Figure 2. The data included a
2
3
4 deletion with a protective effect in the CYP1A1 gene, 3 common SNPs with risk effects in
5
6
7 the CYP1B1 gene, 5 rare SNPs in the COMT gene, which had interaction effects with a
8
9
10 meQTL for the XRCC1 gene, and 5 rare SNPs in the GSTM1 gene, which also had
11
12
13 interaction effects with an eQTL affecting the gene and protein expression of the XRCC3
14
15
16 gene. Finally, 5 rare SNPs in the GSTT1 gene also had interaction effects with a SNP in a
17
18
19 regulatory region. A total of 2,022 SNPs in the four genes (i.e., CYP1B1, COMT, GSTM1,
20
21
22 and GSTT1) and a regulatory region consisting of the meQTL, eQTL, and the SNP
23
24
25 interacting with GSTT1, 1 CNV in CYP1A1, 688 CpGs in XRCC1, and gene and protein
26
27
28 expression levels for 100 genes (including the expression for XRCC3 and 99 other
29
30
31 hypothetical genes in the pathways) were simulated. More details about the simulations can
32
33
34
35
36 be found in the Methods section.

37
38
39
40
41
42 We compared the performance of three multi-omics data analysis methods for disease
43
44
45 prediction using the area under the curve (AUC) measures. The three methods included the
46
47
48 random forest-based method (RFomics), a graph-based integration method (CANetwork) [5],
49
50
51 and a model-based integration method (ATHENA) [4]. The RFomics combines the
52
53
54 preprocessed multi-omics data in a single matrix for constructing the prediction model. As
55
56
57 described in the Methods section, a gene-based risk score is calculated based on SNPs for
58
59
60

1 each gene. Then the risk scores and other multi-omics data are normalized so that they can be
2
3
4 evaluated on the same scale by the RF algorithm. In contrast, CANetwork calculates a graph
5
6
7 matrix to measure the distance between samples using the composite association network
8
9
10 algorithm [16], and the prediction model is created based on the distance matrix using the
11
12
13 graph-based semi-supervised learning algorithm [17]. Finally, ATHENA creates a neural
14
15
16 network model for each type of omics data and a final integrative model is generated based
17
18
19 on these models.
20
21
22
23
24
25

26 Table 1 shows the area under the curve (AUC) for the three methods under three scenarios.
27
28
29 Scenario 1 had 500 cases and 500 controls in the training set, and 100 cases and 100 controls
30
31
32 in the validation set. Scenario 2 had the same sample sizes as those in Scenario 1, but the
33
34
35 multi-omics data had less strong effects on the disease compared to Scenario 1. The effects of
36
37
38 the multi-omics data were the same in Scenarios 3 as those in Scenario 1, but Scenarios 3 had
39
40
41 larger sample size (i.e., 1,500 cases and 1,500 controls in the training data and 500 cases and
42
43
44 500 controls in the validation data). More details of the three scenarios are provided in the
45
46
47 Methods section. Prediction models for the three methods were created based on the training
48
49
50 dataset, and their prediction accuracies were evaluated by the validation dataset. As seen in
51
52
53
54 Table 1, RFomics has the highest AUC in all 3 scenarios followed by ATHENA and
55
56
57
58 CANetwork. Table 2 shows the run time for the three methods. In Scenario 1, RFomics and
59
60
61
62
63
64
65

1 CANetwork had similar performance, while ATHENA required more than 20-times the
2
3
4 runtime of RFomics and CANetwork. In Scenario 3, CANetwork was the most efficient
5
6
7 method followed by RFomics, and ATHENA also required significantly more time than the
8
9
10 other two methods.
11

12 **Discussion**

13
14
15
16 We have developed OmicsSIMLA, which simulates multi-omics data (i.e., genomics,
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
101
102
103
104
105
106
107
108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166
167
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
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
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
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000

We used OmicsSIMLA to simulate a multi-omics dataset for breast cancer based on
hypothetical pathways. Three analysis tools were compared using the dataset. The results

1 suggest that when the different types of data were properly normalized on the same scale, the
2
3
4 RF-based method (i.e., RFomics) achieved the highest AUC. Furthermore, RFomics had
5
6
7 comparable runtime efficiency as that of CANetwork, while ATHENA was computationally
8
9
10 expensive. Therefore, RFomics can potentially be a useful analysis tool for disease prediction
11
12
13 using multi-omics data.
14
15
16
17
18
19

20 Currently, OmicsSIMLA focuses on simulating the dichotomous trait (i.e., affection status).
21
22

23 As studies for quantitative traits are also important, it is our future work to extend
24
25

26 OmicsSIMLA to simulate quantitative traits based on the classic quantitative genetics model
27
28

29 [18]. Furthermore, environmental factors and the interactions between genes and
30
31

32 environments can also play important roles in complex disease etiology. Therefore,
33
34

35 simulating exposome data such as the climate and air quality data and modeling their
36
37

38 interactions with genes are also important in the future extensions of OmicsSIMLA.
39
40
41
42
43
44

45 **Conclusions**

46
47

48 In conclusion, we developed a useful multi-omics data simulator, OmicsSIMLA, for complex
49
50

51 disease studies. Benchmark datasets can be simulated by OmicsSIMLA for evaluating
52
53

54 different multi-omics data analysis methods for disease studies. OmicsSIMLA can also be
55
56

57 used to estimate sample sizes and statistical power when designing a new multi-omics disease
58
59
60
61
62
63
64
65

1 study. As many parameters can be adjusted in OmicsSIMLA, a user-friendly web interface is
2
3
4 provided at <https://omicssimla.sourceforge.io/generateCommand.html> to conveniently specify
5
6
7 these parameters.
8
9

10 11 12 **Methods**

13 14 **Simulation of DNA sequences**

15
16
17 The SeqSIMLA2 package [14] is integrated in OmicsSIMLA to generate DNA sequences in
18
19
20 unrelated/related individuals. Similar to SeqSIMLA2, OmicsSIMLA expects a set of external
21
22
23 reference sequences (i.e., haplotypes) generated by an external sequence generator, such as
24
25
26 COSI [19] or HAPGEN2 [20] that has been widely adopted in genetics studies. Generally, a
27
28
29 set of 10,000 or more reference sequences are expected. Optional files consisting of
30
31
32 recombination rate information and pedigree structures are also accepted. A gene dropping
33
34
35 algorithm assuming random mating with crossovers is performed based on the reference
36
37
38 sequences, recombination rates, and pedigree structures to generate haplotypes in each
39
40
41 individual.
42
43
44
45
46

47 48 **Simulation of CNVs**

49
50
51 For the simulation of CNVs, we considered four CNV states including deletion (D),
52
53
54 normal (N), one duplication (U), and two duplications (UU) on a chromosome.
55
56
57 Therefore, there are 10 types of CNV states on the two chromosomes in an individual, as
58
59
60

1 shown in Supplementary Table S1, and the total copy numbers on the two chromosomes
 2
 3
 4 range from 0 to 6. The user will provide frequencies and ranges of the four CNV states.
 5

6
 7 During meiosis, we use the single-copy crossover model, assuming all crossovers
 8
 9
 10 occurred between CNVs [21].
 11

12 **Simulation of affection status**

13
 14 Genetic variants, including SNPs and CNVs, are used to determine the affection status of an
 15
 16
 17 individual based on a logistic penetrance function as follows:
 18
 19
 20

$$21 \text{logit}(P(\textit{affected})) = \beta_0 + \sum_{i \in \Omega} \beta_{C_{i1}} C_{i1} + \sum_{i \in \Omega} \beta_{C_{i2}} C_{i2} + \sum_{j \in \Psi} \beta_{G_j} G_j + \sum_{m,n \in \Upsilon} \beta_{mn} G_{mn}$$

22
 23 where $P(\textit{affected})$ is the probability of being affected, β_0 determines the baseline
 24
 25
 26 prevalence, Ω , Ψ , and Υ are sets of causal CNVs, SNPs with main effects, and SNPs
 27
 28
 29 with interaction effects, respectively, specified by the user, C_{i1} and C_{i2} are the CNV states for
 30
 31
 32 the first and second haplotypes at CNV i , respectively, G_j is the genotype coding at SNP j ,
 33
 34
 35 and G_{mn} is the genotype coding at SNPs m and n . C_{i1} and C_{i2} have values of -1, 0, 1, and 2 for
 36
 37
 38 CNV states D , N , U , and UU , respectively, where N is the baseline state. The coding of G_j is
 39
 40
 41 based on a dominant, additive or recessive model, and the coding of G_{mn} is based on several
 42
 43
 44 interaction models. If SNP j is in a CNV region, allelic CNV [22] is considered in the coding
 45
 46
 47 of G_j . More details of the coding of G_j and G_{mn} are provided in Supplementary methods. The
 48
 49
 50 parameters β_C and β_G are the effect sizes of the main effects for CNVs and SNPs,
 51
 52
 53 respectively, and β_{mn} determines the effect size of the interaction effect between SNPs m
 54
 55
 56
 57
 58
 59
 60
 61
 62
 63
 64
 65

1 and n . These parameters are specified by the user.
2
3

4 **Simulation of DNA methylation data** 5

6
7 The pWGBSSimla package [11] is integrated into OmicsSIMLA to generate the WGBS data.
8

9
10 The pWGBSSimla algorithm simulates data using methylation profiles generated based on 41
11

12 WGBS datasets for 29 human cell and tissue types. The profiles contain the information for
13

14 each CpG, such as its distance to the next site, methylation rate, methylation status (i.e.,
15

16 methylated, unmethylated, and fuzzily methylated), and read counts for each type of
17
18

19 methylation status. CpGs and the distances between the CpGs are first determined based on
20
21

22 the profiles, and then the total read count and methylated read count are simulated for each
23
24

25 CpG. Methylation level at a CpG influenced by a meQTL is simulated based on a genotype-
26
27

28 specific methylation probability, which is the methylation rate of the CpG in the profiles
29
30

31 multiplied by a ratio following an exponential distribution. Furthermore, ASMs are simulated
32
33

34 based on father- and mother-specific methylation rates for paternal and maternal alleles,
35
36

37 respectively. Finally, a DMR is generated by simulating the same genomic region using
38
39

40 profiles for different cell or tissue types. More details of the pWGBSSimla algorithm can be
41
42

43 found in Chung and Kang [11].
44
45

46 **Simulation of RNA-seq data** 47

48 We implemented a parametric simulation procedure for simulating the RNA-seq data
49
50

51 similar to that described in Benidt and Nettleton [13]. A negative binomial (NB)
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1 distribution with mean μ_{ij} and dispersion parameter ω_i is used to simulate the read
 2
 3
 4 count for gene i in individual j . The mean is calculated as $\mu_{ij} = \lambda_i c_j$, where λ_i is the
 5
 6
 7 common mean for gene i and c_j is the individual-specific normalization factor for
 8
 9
 10 individual j . The parameters λ , c , and ω for all genes were estimated using the R
 11
 12
 13 package edgeR [23] based on a whole genome RNA-seq dataset consisting of 103 normal
 14
 15
 16 tissues in patients with breast cancer from The Cancer Genome Atlas (TCGA) project
 17
 18
 19 [24]. The parameters λ_i and ω_i are randomly sampled with replacement from λ and
 20
 21
 22 ω . If more than 103 samples are simulated, we use the smoothed bootstrap procedure
 23
 24
 25 [25] to calculate c_j^* for individual j , and μ_{ij} is calculated as $\lambda_i c_j^*$. More details of the
 26
 27
 28 calculation of c_j^* are provided in Supplementary methods. The user can specify n
 29
 30
 31 differentially expressed (DE) genes between cases and controls and their fold changes,
 32
 33
 34 and the read count for DE gene i in individual j is simulated based on a NB distribution
 35
 36
 37 with mean $f_i \mu_{ij}$ and dispersion parameter ω_i , where f_i is the fold change for gene i .
 38
 39
 40
 41

42 Simulation of eQTL and allele-specific reads

43
 44
 45 We followed the procedure in the simulation study in Sun [26] to simulate eQTL and read
 46
 47
 48 counts for ASE. For eQTL l with a user-specified fold change h_l , the means for the three
 49
 50
 51 genotypes AA , Aa , and aa at the eQTL are μ_{ij} , $h_l \mu_{ij}$, and $(2h_l - 1)\mu_{ij}$, respectively, and the
 52
 53
 54 dispersion parameter is ω_i in the NB distribution for gene i influenced by the eQTL. ASE
 55
 56
 57 for a gene caused by a cis-eQTL is simulated by assuming reads were mapped to
 58
 59
 60

heterozygous SNPs (i.e., allele-specific reads) in the gene. A cis-eQTL refers to the eQTL being located in the cis-regulatory elements of the gene. Because the alleles at the cis-eQTL can be in the same haplotype as the alleles of the gene, ASE can be observed using the allele-specific reads of the gene. Furthermore, only heterozygous SNPs can be tested for cis-eQTL with the allele-specific reads. Therefore, we simulate allele-specific reads for heterozygous eQTLs. Assuming t_{ij} is the total read count for gene i in individual j , the total number of allele-specific reads is calculated as $0.005t_{ij}$, where 0.005 was estimated from real data by Sun [26]. Furthermore, also suggested by Sun [26], the number of allele-specific reads for a haplotype is simulated using a beta-binomial distribution with a mean determined by the effect size of the cis-eQTL and an overdispersion parameter of 0.1. The effect size is defined as $\log_2(\text{expression of the alternative allele at the eQTL}/\text{expression of the reference allele at the eQTL})$ [27] for a heterozygous cis-eQTL and is set to 0 for a homozygous cis-eQTL.

Simulation of eQTM

We used linear regression to model the relationship between gene expression and methylation:

$\mu_i = E(y_{ij}) = \alpha_i + \beta_i x_{ij}$, where y_{ij} and x_{ij} are the RNA-seq read count and the proportion of methylated reads, respectively, for gene i influenced by methylation in individual j . Assuming that the NB parameters for gene i are μ_i and ϕ_i , the parameter α_i is specified as μ_i , and β_i is assumed to follow a normal distribution with a mean and a standard deviation specified

1 by the user. Then the gene expression of gene i is simulated by an NB distribution with
2
3
4 parameters of μ_i and ϕ_i .
5
6

7 **Protein expression simulation**

8
9

10 We assumed that the protein expression level for protein k at a time point t in sample j
11
12 follows a normal distribution with a mean η_{kjt} and a standard deviation τ_k after
13
14 normalization. We used the mass-action kinetic action model [15] to simulate protein
15
16 expression at a certain time point. The mean $\eta_{k,j,t+1}$ for the protein expression at time $t+1$
17
18 was determined as follows:
19
20
21
22
23
24

$$25 \eta_{k,j,t+1} = \eta_{kjt} + (x_{kjt} \kappa_{jt}^s - \eta_{kjt} \kappa_{jt}^d),$$

26
27
28

29 where x_{kjt} is the normalized gene expression for the gene encoding protein k , and κ_{jt}^s and
30
31 κ_{jt}^d are the protein synthesis and degradation rates, respectively, in individual j at time t . The
32
33 normalized gene expression x_{kjt} is calculated using the median absolute deviation (MAD)
34
35 scale normalization [28] based on the RNA-seq data simulated from the previous section.
36
37
38
39

40 Similar to the simulation study in Teo et al. [15], κ_{jt}^d is fixed to be 1, and κ_{jt}^s with a default
41
42 value of 1 can be changed by the user. A vector of standard deviations τ were estimated
43
44
45 from the level 4 protein expression data of primary tumor tissue in 874 breast cancer patients
46
47
48 from the TCGA project downloaded from the cancer proteome atlas (TCPA) [29] website.
49
50
51

52 The level 4 data consist of protein expression data for 224 proteins that have been normalized
53
54
55 across the samples as well as across the proteins, and a replication-based method was used to
56
57
58
59
60

1 account for differences in protein expression among different batches. The parameter τ_j is
2
3
4 then randomly sampled with replacement from τ .
5
6

7 **A random-forest based method for integrating multi-omics data for disease studies**

8
9

10 Multi-omics data can have different data types (e.g., discrete data for SNP genotypes,
11
12 categorical data for CNV statuses, and continuous data for proportions of methylated reads,
13
14 RNA-seq read counts, and normalized protein expression) and different variations (e.g., three
15
16 possible values of 0, 1, and 2 for minor allele counts at SNPs, and real numbers ranging
17
18 between 0 and 1 for the proportions of methylated reads). When developing a method for
19
20 integrating these data, it is important to account for the properties of different data types so
21
22 that the analysis results would not be biased toward certain variables [9]. We developed a
23
24 preprocessing algorithm for the multi-omics data. A gene-based risk score, which is a
25
26 weighted sum of the numbers of risk alleles at SNPs in the gene, for each individual is
27
28 constructed. The weights are the effect sizes of the risk alleles at the SNPs. More details for
29
30 calculating the risk score are provided in Supplementary methods. Then each variable from
31
32 different omics data, including the gene-based risk scores, CNV statuses of genes,
33
34 methylation proportions at CpGs, gene and protein expression levels, is normalized so that it
35
36 has a mean 0 and a standard deviation of 1. The normalized variables are then used in RF for
37
38 classification.
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56

57 **Simulation studies**

58
59
60
61
62
63
64
65

1 We used OmicsSIMLA to evaluate the performance of the proposed RF-based method,
2
3
4 compared with CANetwork and ATHENA. A hypothetical disease model for breast cancer
5
6
7 involving multi-omics data [9] was simulated, as shown in Figure 2. To be more specific, a
8
9
10 deletion with a frequency of 20%, which had a protective effect with an odds ratio (OR) of
11
12
13 0.67, in the CYP1A1 gene and 3 common variants, which had main effects (ORs = 1.5) with
14
15
16 minor allele frequencies (MAFs) > 10%, in the CYP1B1 gene were simulated. We also
17
18
19 simulated 5 rare variants with MAFs < 3% in the COMT gene, which had interaction effects
20
21
22 (ORs = 5) with a meQTL for the XRCC1 gene. The CpG in XRCC1 influenced by the
23
24
25 meQTL caused a difference in methylation rates of 10% between cases and controls.
26
27
28
29 Furthermore, we simulated 5 rare variants in the GSTM1 gene, which had interaction effects
30
31
32 (ORs = 5) with a cis-eQTL for the XRCC3 gene, and 5 rare variants in the GSTT1 gene,
33
34
35 which had interaction effects (ORs = 5) with a SNP located in the same region as that of the
36
37
38 meQTL and eQTL. The eQTL caused a fold change of 1.5 in the XRCC3 gene expression
39
40
41 compared to the reference genotype, and a fold change of 1.5 was simulated for the
42
43
44 differential gene expression of XRCC3 between cases and controls. In summary, the total
45
46
47 variables consisted of 200, 687, 264, and 176 SNPs in the CYP1B1, COMT, GSTM1, and
48
49
50 GSTT1 genes, respectively, and 695 SNPs harboring the meQTL, eQTL, and the SNP
51
52
53 interacting with GSTT1 in the regulatory region, a variable for CNV status in CYP1A1,
54
55
56
57 methylation levels at 688 CpGs in XRCC1, and gene and protein expression levels for 100
58
59
60

1 genes and their encoded proteins. More details for generating the reference sequences in the
2
3
4 genes and the simulations for each omics data type are provided in Supplementary methods.
5
6
7

8
9
10 We simulated a training dataset consisting of 500 cases and 500 controls as well as a
11
12 validation dataset consisting of 100 cases and 100 controls. The training dataset was used by
13
14 RFomics, CANetwork, or ATHENA to construct a prediction model. The validation dataset
15
16 was then used to calculate the AUC based on the prediction model. Note that a 5-fold cross-
17
18 validation was performed in ATHENA, and a best model based on the testing dataset (i.e.,
19
20 one of the five random 20% of the training dataset) was created for each cross-validation. The
21
22 model with the highest AUC based on the testing dataset was selected and applied to the
23
24 validation dataset. This simulation scenario was referred to as Scenario 1. We also simulated
25
26 a scenario with less strong genetic effects (Scenario 2) and a scenario with larger sample size
27
28 (Scenario 3). More details about Scenarios 2 and 3 are provided in Supplementary methods.
29
30 For each scenario, 1,000 batches of training and validation datasets were simulated, and the
31
32 AUC for each algorithm was averaged over the 1,000 batches.
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50

51 **Availability of supporting source code and requirements**

52 Project name: OmicsSIMLA
53
54

55 Project home page: <https://omicssimla.sourceforge.io>
56
57
58
59
60

1 Operating system: Linux
2

3
4 Programming language: C++
5

6
7 Other requirements: C++11 compiler and Eigen and boost libraries if directly compiling the
8
9
10 source code.
11

12
13 License: GPL-3.0
14
15
16
17
18
19

20 **Availability of supporting data**

21

22
23 The simulated datasets supporting the conclusions of this article are available from the
24
25
26 OmicsSIMLA website (<https://omicssimla.sourceforge.io/download.html>).
27
28
29
30
31

32 **Declarations**

33

34 **Funding**

35

36
37
38
39 This work has been supported by a grant from the Ministry of Science and Technology
40
41
42 (MOST 106-2221-E-400-005-MY3) in Taiwan.
43
44

45 **Authors' contributions**

46

47
48 RHC and CYK both designed the framework of the simulation tool and implemented the
49
50
51 software. RHC designed the simulation study and CYK performed the simulation analysis.
52
53

54
55 Both authors read and approved the final manuscript.
56

57 **Competing interests**

58
59
60
61
62
63
64
65

The authors declare that they have no competing interests.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

References

1. Timpson NJ, Greenwood CMT, Soranzo N, Lawson DJ and Richards JB. Genetic architecture: the shape of the genetic contribution to human traits and disease. *Nature reviews Genetics*. 2018;19 2:110-24. doi:10.1038/nrg.2017.101.
2. Karczewski KJ and Snyder MP. Integrative omics for health and disease. *Nature reviews Genetics*. 2018;19 5:299-310. doi:10.1038/nrg.2018.4.
3. Hasin Y, Seldin M and Luskis A. Multi-omics approaches to disease. *Genome biology*. 2017;18 1:83. doi:10.1186/s13059-017-1215-1.
4. Holzinger ER, Dudek SM, Frase AT, Pendergrass SA and Ritchie MD. ATHENA: the analysis tool for heritable and environmental network associations. *Bioinformatics*. 2014;30 5:698-705. doi:10.1093/bioinformatics/btt572.
5. Yan KK, Zhao H and Pang H. A comparison of graph- and kernel-based -omics data integration algorithms for classifying complex traits. *BMC bioinformatics*. 2017;18 1:539. doi:10.1186/s12859-017-1982-4.
6. Ruffalo M, Koyuturk M and Sharan R. Network-Based Integration of Disparate Omic Data To Identify "Silent Players" in Cancer. *PLoS computational biology*. 2015;11 12:e1004595. doi:10.1371/journal.pcbi.1004595.
7. Jennings EM, Morris JS, Carroll RJ, Manyam GC and Baladandayuthapani V. Bayesian methods for expression-based integration of various types of genomics data. *EURASIP J Bioinform Syst Biol*. 2013;2013 1:13. doi:10.1186/1687-4153-2013-13.
8. Tyekuceva S, Marchionni L, Karchin R and Parmigiani G. Integrating diverse genomic data using gene sets. *Genome biology*. 2011;12 10:R105. doi:10.1186/gb-2011-12-10-r105.
9. Ritchie MD, Holzinger ER, Li R, Pendergrass SA and Kim D. Methods of integrating data to uncover genotype-phenotype interactions. *Nature reviews Genetics*. 2015;16 2:85-97. doi:10.1038/nrg3868.
10. Rackham OJ, Dellaportas P, Petretto E and Bottolo L. WGBSSuite: simulating whole-genome bisulphite sequencing data and benchmarking differential DNA methylation analysis tools. *Bioinformatics*. 2015;31 14:2371-3. doi:10.1093/bioinformatics/btv114.
11. Chung R-H and Kang C-Y. pWGBSSimla: a profile-based whole-genome bisulphite sequencing data simulator incorporating methylation QTLs, allele-specific methylations and differentially methylated regions. *bioRxiv*. 2018; doi:10.1101/390633.
12. Frazee AC, Jaffe AE, Langmead B and Leek JT. Polyester: simulating RNA-seq datasets with differential transcript expression. *Bioinformatics*. 2015;31 17:2778-84. doi:10.1093/bioinformatics/btv272.
13. Benidt S and Nettleton D. SimSeq: a nonparametric approach to simulation of RNA-

- 1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
- sequence datasets. *Bioinformatics*. 2015;31 13:2131-40.
doi:10.1093/bioinformatics/btv124.
14. Chung RH, Tsai WY, Hsieh CH, Hung KY, Hsiung CA and Hauser ER. SeqSIMLA2: simulating correlated quantitative traits accounting for shared environmental effects in user-specified pedigree structure. *Genetic epidemiology*. 2015;39 1:20-4.
doi:10.1002/gepi.21850.
15. Teo G, Vogel C, Ghosh D, Kim S and Choi H. *A Mass-Action-Based Model for Gene Expression Regulation in Dynamic Systems*. Cambridge University Press; 2015.
16. Mostafavi S, Ray D, Warde-Farley D, Grouios C and Morris Q. GeneMANIA: a real-time multiple association network integration algorithm for predicting gene function. *Genome biology*. 2008;9 Suppl 1:S4. doi:10.1186/gb-2008-9-s1-s4.
17. Tsuda K, Shin H and Scholkopf B. Fast protein classification with multiple networks. *Bioinformatics*. 2005;21 Suppl 2:ii59-65. doi:10.1093/bioinformatics/bti1110.
18. Falconer DS and Mackay TF. *Quantitative genetics*. San Francisco: Benjamin Cummings; 1996.
19. Schaffner SF, Foo C, Gabriel S, Reich D, Daly MJ and Altshuler D. Calibrating a coalescent simulation of human genome sequence variation. *Genome research*. 2005;15 11:1576-83. doi:10.1101/gr.3709305.
20. Su Z, Marchini J and Donnelly P. HAPGEN2: simulation of multiple disease SNPs. *Bioinformatics*. 2011;27 16:2304-5. doi:10.1093/bioinformatics/btr341.
21. Hartasanchez DA, Valles-Codina O, Braso-Vives M and Navarro A. Interplay of interlocus gene conversion and crossover in segmental duplications under a neutral scenario. *G3*. 2014;4 8:1479-89. doi:10.1534/g3.114.012435.
22. Usher CL and McCarroll SA. Complex and multi-allelic copy number variation in human disease. *Briefings in functional genomics*. 2015;14 5:329-38.
doi:10.1093/bfgp/elv028.
23. Robinson MD, McCarthy DJ and Smyth GK. edgeR: a Bioconductor package for differential expression analysis of digital gene expression data. *Bioinformatics*. 2010;26 1:139-40. doi:10.1093/bioinformatics/btp616.
24. Cancer Genome Atlas Research N. Comprehensive genomic characterization defines human glioblastoma genes and core pathways. *Nature*. 2008;455 7216:1061-8.
doi:10.1038/nature07385.
25. Efron B and Tibshirani RJ. *An Introduction to the Bootstrap*. Chapman and Hall/CRC; 1993.
26. Sun W. A statistical framework for eQTL mapping using RNA-seq data. *Biometrics*. 2012;68 1:1-11. doi:10.1111/j.1541-0420.2011.01654.x.
27. Mohammadi P, Castel SE, Brown AA and Lappalainen T. Quantifying the regulatory effect size of cis-acting genetic variation using allelic fold change. *Genome research*.

2017;27 11:1872-84. doi:10.1101/gr.216747.116.

28. Fundel K, Kuffner R, Aigner T and Zimmer R. Normalization and gene p-value estimation: issues in microarray data processing. *Bioinform Biol Insights*. 2008;2:291-305.
29. Li J, Lu Y, Akbani R, Ju Z, Roebuck PL, Liu W, et al. TCPA: a resource for cancer functional proteomics data. *Nature methods*. 2013;10 11:1046-7. doi:10.1038/nmeth.2650.

Figures

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

Figure 1. Simulation framework of OmicsSIMLA. The black solid lines represent the relationships among different types of omics data. The black dotted lines represent the causal effects of genomics data to the disease. The red dotted lines represent the retrospective simulations of the methylation, gene expression and protein expression levels conditional on the disease status.

Figure 2. Hypothetical pathways involved in breast cancer. The brown solid lines represent the main effects of SNPs and CNVs on the disease, while the green solid lines represent the interaction effects of SNPs on the disease. The black solid lines represent the regulatory effects of the meQTL and eQTL on methylation and gene expression, respectively. The red dotted lines represent the retrospective simulations of the methylation, gene expression and protein expression levels conditional on the disease status.

Tables

Table 1. Area under the curve (AUC) for RFomics, CANetwork, and ATHENA under different scenarios

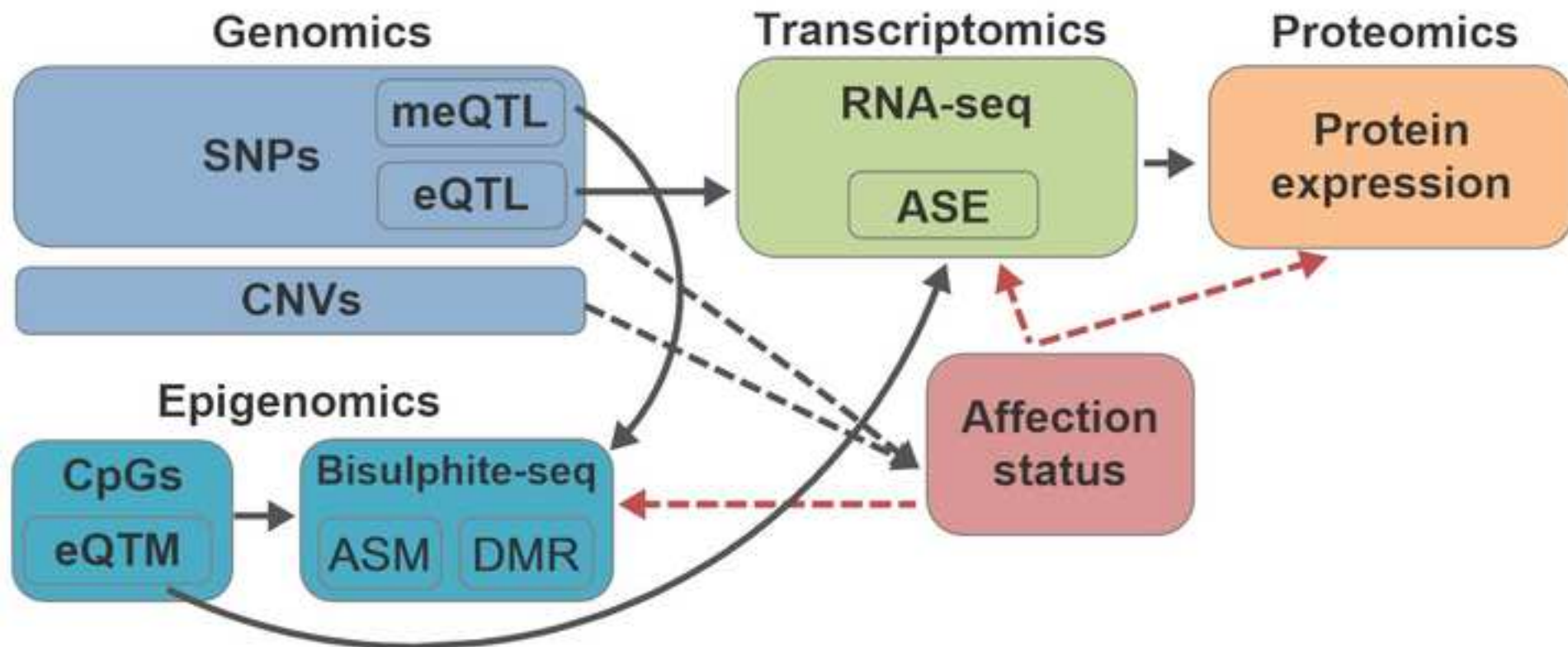
	RFomics	CANetwork	ATHENA
Scenario 1	0.861 (0.026) ¹	0.596 (0.042)	0.831 (0.042)
Scenario 2	0.566 (0.041)	0.529 (0.029)	0.559 (0.068)
Scenario 3	0.876 (0.012)	0.649 (0.019)	0.835 (0.031)

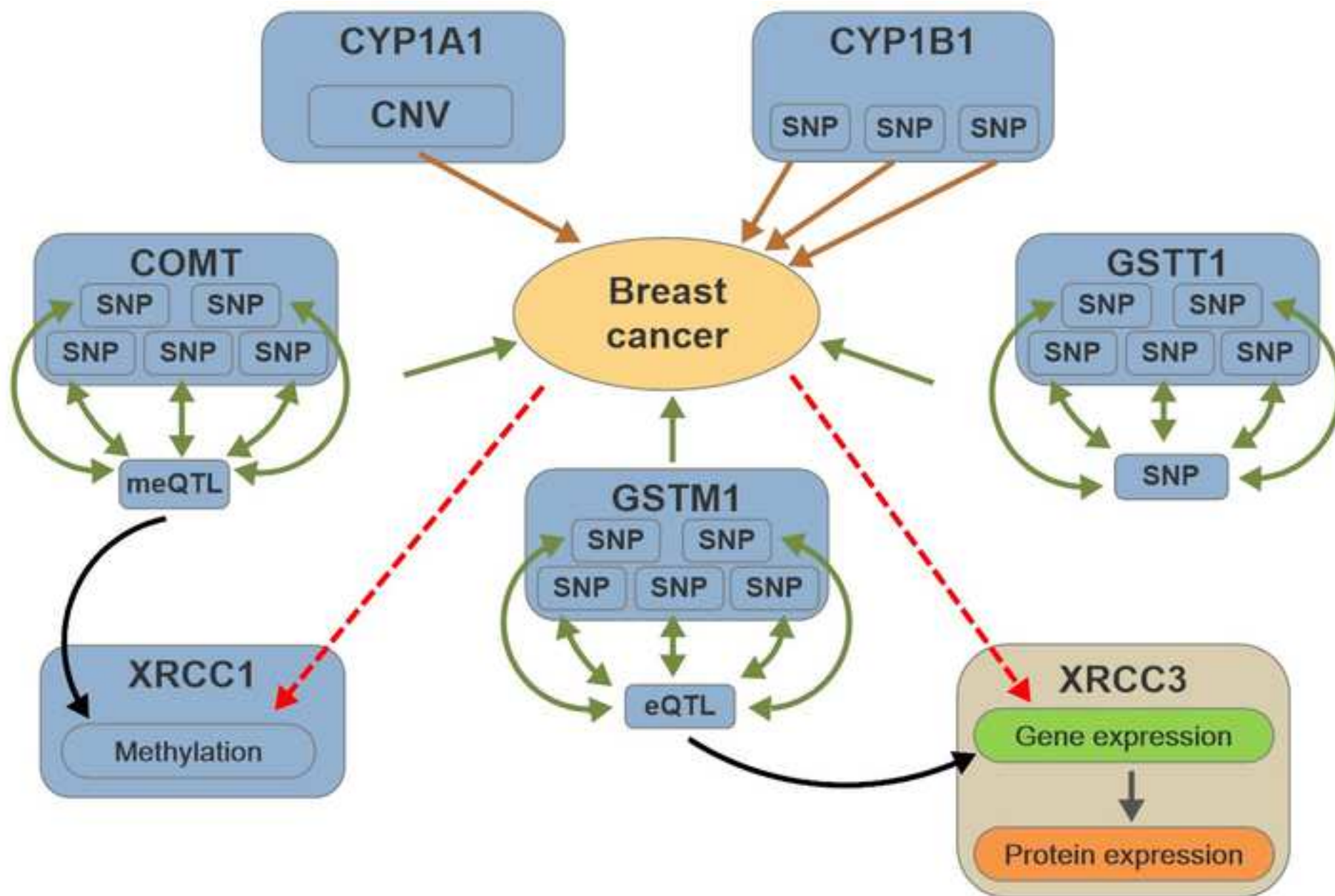
¹The mean AUC and its standard error estimated based on 1,000 batches

Table 2. Run time (in seconds) for RFomics, CANetwork, and ATHENA under Scenarios 1 and 3

	RFomics	CANetwork	ATHENA
Scenario 1	37.78	40.54	823.14
Scenario 3	143.91	94.16	2195.95

¹The mean time (in seconds) was estimated based on 100 batches







Click here to access/download
Supplementary Material
Supplementary_material.docx

