

1 **RegionScan: A comprehensive R package for region-level genome-wide association testing**  
2 **with integration and visualization of multiple-variant and single-variant hypothesis testing**

3 Myriam Brossard<sup>1</sup>, Delnaz Roshandel<sup>2</sup>, Kexin Luo<sup>1</sup>, Fatemeh Yavartanoo<sup>3</sup>, Andrew D.

4 Paterson<sup>2,4</sup> Yun J. Yoo<sup>3</sup>, Shelley B. Bull<sup>1,4</sup>

5 <sup>1</sup>Lunenfeld-Tanenbaum Research Institute, Sinai Health, Toronto, Ontario, Canada;

6 <sup>2</sup>Genetics & Genome Biology Program, The Hospital for Sick Children, Toronto, Ontario,  
7 Canada;

8 <sup>3</sup>Seoul National University, Seoul, South Korea;

9 <sup>4</sup>Dalla Lana School of Public Health, University of Toronto, Toronto, Ontario, Canada.

10 **Abstract**

11 **Summary:** RegionScan is an R package for comprehensive and scalable genome-wide  
12 association testing of region-level multiple-variant and single-variant statistics and visualization  
13 of the results. It implements various state-of-the-art region-level tests to improve signal detection  
14 under heterogeneous genetic architectures and facilitates comparison of multiple-variant region-  
15 level and single-variant test results. It exploits local linkage disequilibrium (LD) structure for  
16 genomic partitioning and LD-adaptive region definition. RegionScan is compatible with VCF  
17 input file formats for genotyped and imputed variants, and options are available for analysis of  
18 multi-allelic variants and unbalanced binary phenotypes. It accommodates parallel region-level  
19 processing and analysis to improve computational time and memory efficiency and provides  
20 detailed outputs and utility functions to assist results comparison, visualization, and  
21 interpretation.

22 **Availability and implementation:** RegionScan is freely available for download on GitHub  
23 (<https://github.com/brossardMyriam/RegionScan>).

24 **Contact:** [bull@lunenfeld.ca](mailto:bull@lunenfeld.ca), [brossard@lunenfeld.ca](mailto:brossard@lunenfeld.ca).

25 **Supplementary information:** Supplementary data are available at Bioinformatics online.

## 26 **1. Introduction**

27 Compared to genome-wide single-variant testing, region-level multi-variant association analysis  
28 can better capture signals under complex genetic architectures<sup>1</sup>. Because fewer tests are  
29 conducted, multiple testing is reduced, and the genome-wide testing threshold can be relaxed.  
30 However, for comprehensive genomic analysis, region-level testing requires appropriate region  
31 definition, e.g. including intergenic, intronic, and exonic variants. It also faces analytical  
32 challenges, including high dimensionality and multi-collinearity within regions produced by  
33 complex and long-range linkage disequilibrium (LD) structure. Available region-level tests differ  
34 according to the underlying assumptions, the construction of the test statistic, and thus are  
35 sensitive to different regional genetic architectures<sup>2,3</sup>. We focus on three classes of state-of-the-  
36 art region-level tests (**Supplementary Information 1**), including multi-variant linear/logistic  
37 regression tests with and without dimension reduction<sup>3-5</sup>, variance component score tests<sup>6,7</sup>, and  
38 region-level  $\min P$  tests<sup>8-10</sup>; sensitive to heterogeneous regional architectures. Our goal is to  
39 integrate region definition with implementation of region-level and single-variant tests in one  
40 scalable R package for comprehensive genome-wide region-level association analysis and  
41 improve region discovery under heterogeneous regional architectures.

## 42 **2. Implementation and Key features**

43 We introduce the RegionScan R package for genome-wide discovery analysis and define regions  
44 using the `gpart`<sup>11</sup> R package for LD-based genomic partitioning, optimized for region-level  
45 analysis<sup>12</sup>. Although a major advantage of our approach is comprehensive analysis of the  
46 genome, including intergenic regions, RegionScan can also accommodate other user-specified  
47 region definitions.

## 48 **2.1 Capability and Scalability**

49 The main function is called *regscan* (**Supplementary Information 2**, for a detailed description  
50 and list of options). *regscan* takes four main inputs: *data* which includes genotypes (additively  
51 coded); *SNPinfo* input with variant information; *phenocov* with phenotypes (quantitative or  
52 binary) as well as covariates (if applicable) and a *REGIONinfo* input with region start/end  
53 positions, as produced with *gpart*<sup>11</sup>. Alternatively, the auxiliary function *recodeVCF* can process  
54 large VCF 4.0 files with *vcftools*<sup>13</sup> to produce a temporary subset VCF file by region,  
55 subsequently processed in R to improve memory efficiency. *regscan* also deals optionally with  
56 multi-allelic variants in addition to bi-allelic variants.

57 To improve scalability, *regscan* can process, recode and analyze each region in parallel. The  
58 processing steps for each region include variant filtering and recoding based on minor allele  
59 frequency (MAF), and an option to reduce multicollinearity by pruning variants within regions.  
60 This is followed by application of region-level tests including regression-based tests (MLC<sup>2</sup>,  
61 PC80<sup>3</sup>, LC<sup>4,5,14,15</sup>, generalized Wald tests), variance component score tests (SKAT<sup>16,17</sup>, SKAT-  
62 O<sup>7</sup>), and region-level min *P* tests (simpleM<sup>8</sup>, GATES<sup>9</sup>, MinP<sup>10,18</sup>), in addition to single-SNP  
63 tests for variants within regions. *regscan* includes an option to reduce finite-sample bias in  
64 logistic regression of unbalanced binary traits and/or variants with low minor allele counts, using  
65 a Jeffreys-prior penalized likelihood<sup>19,20</sup>. For the MLC<sup>2</sup> region-level test, variants within each  
66 region are clustered in LD bins based on pairwise correlation<sup>21</sup> for reduced-*df* region-level  
67 testing adaptive to the number of LD bins, followed by variant recoding within each bin to  
68 maximize variant pairs positive correlation (**Supplementary Information 2**, section 2.1.1); bin-  
69 level tests within regions are reported in addition to MLC region-level tests.

## 70 **2.3 Detailed Outputs and Visualization**

71 *regscan* produces six outputs detailing results for all regions analyzed (**Supplementary**  
72 **Information 2**, section 2.1.3): (1) *region-level* output with results for all regions analyzed; (2)  
73 *bin-level* output including bin-level test results for all bins within each region; (3) *variant-level*  
74 output with variant positions, LD-bin assignments, and corresponding effect sizes and *P*-values  
75 from single- and multi-variant regional regression models; within-region variance inflation factor  
76 values (VIFs) are included to facilitate identification of multi-collinearity; (4) a *list of variants*  
77 *pruned out* with reasons for exclusions; and optionally, (5) a *single-variant* output including  
78 variant-to-LD bin assignments for all the variants (available before pruning) and (6) a *covariate*  
79 output with covariate effects and *P*-values extracted from multi-variant regional regression  
80 models.

81 Utility functions are implemented to visualize comparisons between region and/or variant-level  
82 test results. For example, *MiamiPlot* produces a genome-wide comparison of  $-\log_{10}$  *P*-values for  
83 a pair of tests; *LocusPlot* displays results of several region-level tests in a set of contiguous  
84 regions; *QQPlot* assesses consistency of the observed distribution of a specified region-level  
85 statistic *P*-value with that expected under the null hypothesis and returns corresponding genomic  
86 inflation factors. *regscan* also produces optional heatmap plots within each of a set of selected  
87 region(s) to visualize correlation within region and within/across LD bins; it annotates variant  
88 positions according to the LD-bin assignment.

## 89 **3 Usage case**

90 In **Supplementary Information 3**, we demonstrate RegionScan capabilities and computational  
91 efficiency by genome-wide analysis of 1,340 individuals with type 1 diabetes (T1D) from the

92 DCCT/EDIC<sup>22,23</sup> genetic study for LDL-cholesterol (LDL-c, measured at baseline). **Fig. 1** gives  
93 an overview of RegionScan capabilities based on the results for chr19. In **Supplementary**  
94 **Information 4**, we report computation time by sample size and region size based on a realistic  
95 test dataset.

## 96 **Conclusion**

97 *RegionScan* is a flexible and versatile R package designed for scalable and comprehensive  
98 genome-wide region-level analysis that leverages region definition adaptive to local LD structure  
99 (or any other user-provided region definition). It implements multiple region-level tests sensitive  
100 to heterogeneous genetic architecture, including LD-bin reduced-*df* region-level tests, facilitates  
101 comparisons of region-level and single-variant test results, and includes options to deal with high  
102 dimensionality and multi-collinearity arising from improving resolution of  
103 genotyped/sequenced/imputed genetic data. Modular design is flexible for future developments.

## 104 **Funding**

105 This project was supported by: CIHR Project Grant (#PJT-159463), CIHR STAGE fellowship  
106 (MB #GET-101831) and NSERC Discovery Grant (SBB #RGPIN-05896). **Conflict of Interest:**  
107 none declared.

## 108 **Software, code and data availability**

109 The R package RegionScan (<https://github.com/brossardMyriam/RegionScan>) is available on  
110 GitHub and includes a vignette on how to install and run RegionScan in a realistic artificial  
111 dataset provided. The DCCT/EDIC data are available to authorized users at

112 <https://repository.niddk.nih.gov/studies/edic/>  
113 and [https://www.ncbi.nlm.nih.gov/projects/gap/cgi-bin/study.cgi?study\\_id=phs000086.v3.p1](https://www.ncbi.nlm.nih.gov/projects/gap/cgi-bin/study.cgi?study_id=phs000086.v3.p1)  
114 (IRB #07-0208-E). Data analysis, software development, and computation time estimation were  
115 performed on the Hospital for Sick Children High-performance Computing Facility, the  
116 Lunenfeld-Tanenbaum Research Institute High-performance Computing platform, and the  
117 Niagara supercomputer (with support from the Canada Foundation for Innovation under the  
118 auspices of Compute Canada, the Government of Ontario, Ontario Research Fund - Research  
119 Excellence, and the University of Toronto). For hardware specifications on Niagara, see  
120 [https://docs.computeCanada.ca/wiki/Niagara#Niagara\\_hardware\\_specifications](https://docs.computeCanada.ca/wiki/Niagara#Niagara_hardware_specifications)  
121 and [https://docs.scinet.utoronto.ca/index.php/Niagara\\_Quickstart](https://docs.scinet.utoronto.ca/index.php/Niagara_Quickstart).

## 122 **Acknowledgements**

123 This study uses data provided by the Diabetes Control and Complications Trial / Epidemiology  
124 of Diabetes Interventions and Complications (DCCT/EDIC) Research Group which is sponsored  
125 through research contracts from the National Institute of Diabetes, Endocrinology and Metabolic  
126 Diseases of the National Institute of Diabetes and Digestive Kidney Diseases (NIDDK) and the  
127 National Institutes of Health (NIH). The authors are grateful to the subjects in the DCCT/EDIC  
128 cohort for their long-term participation. A complete list of the individuals and institutions  
129 participating in the DCCT/EDIC Research Group can be found in **Supplementary Information**  
130 **3**. This project was supported by: CIHR Project Grant (#PJT-159463), CIHR STAGE fellowship  
131 (MB #GET-101831).

132

133 **Fig. 1.** Overview of *RegionScan* for genome-wide region-level association analysis of LDL-c at  
134 baseline in 1,340 individuals from the DCCT/EDIC<sup>22,23</sup> Genetics Study of the Usage case study.  
135 Details of analysis are described in **Supplementary Information 3**. To facilitate visualization  
136 of the results, we illustrate results on chr19 which exhibits genome-wide region-level association  
137 signals at the region- and variant-levels. Panel (A) illustrates a comparison between region-level  
138 association results based, for example, on the MLC test (top panel) with single-SNP results  
139 (bottom panel) for 89,001 regions analyzed genome-wide; the dotted lines indicate the genome-  
140 wide Bonferroni-corrected significance levels:  $5.6E-7$  for region-level tests (top panel) and  $5E-8$   
141 (bottom panel) for single-SNP tests. Panel (B) illustrates partitioning results for 13 regions in  
142 chr19: 45,257,201-45,436,657bp; gene positions are shown in GRCh37. The blocks are delimited  
143 by triangles. Panel (C) illustrates comparison of results for multiple region-level tests for the  
144 same 13 regions as illustrated in Panel (B); changes in grey shading facilitate visualization of  
145 region boundaries. Panel (D) shows the LD bins constructed for the MLC test within the top  
146 LDL-c associated region, overlapping *APOE* gene (chr19: region #1690, chr19:45,385,759-  
147 45,415,935). The left panel shows the heatmap of the SNP correlation matrix (with SNPs ordered  
148 by position within LD bin, and LD bins ordered by number of SNPs assigned); the right panel  
149 shows the SNP positions (*X* axis) along the LD bins (*Y* axis).

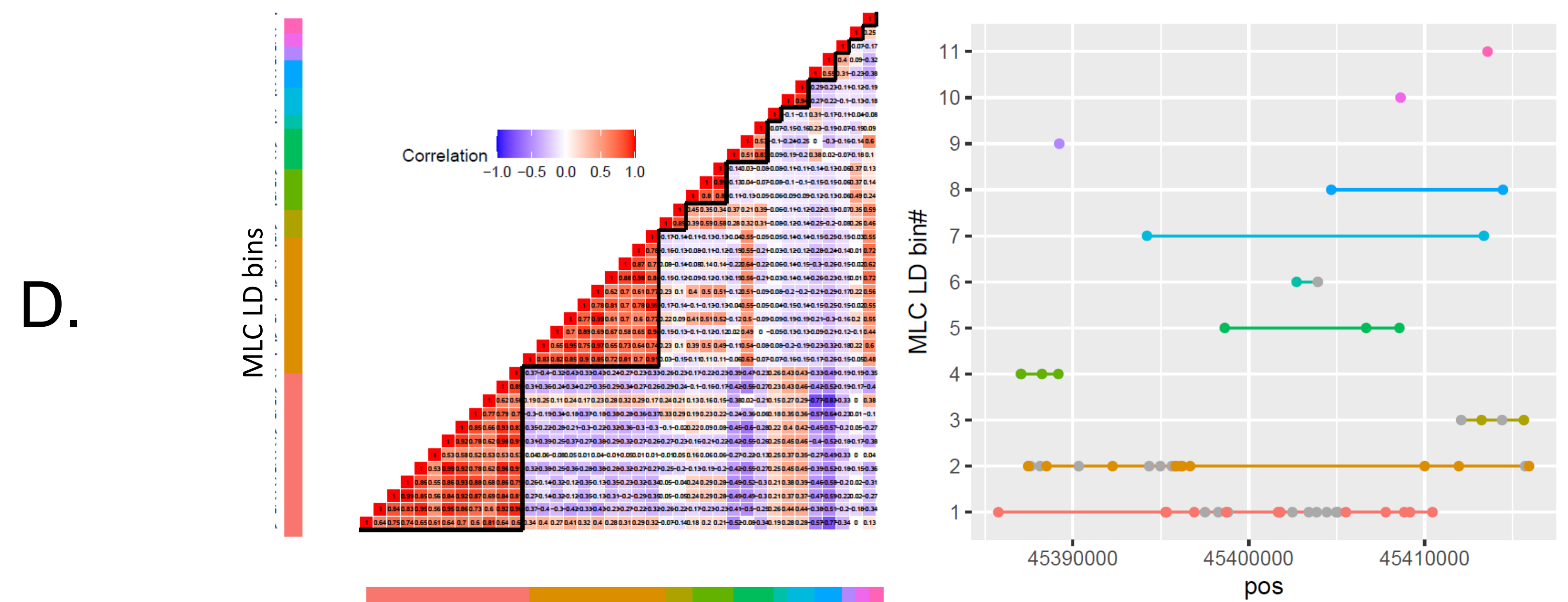
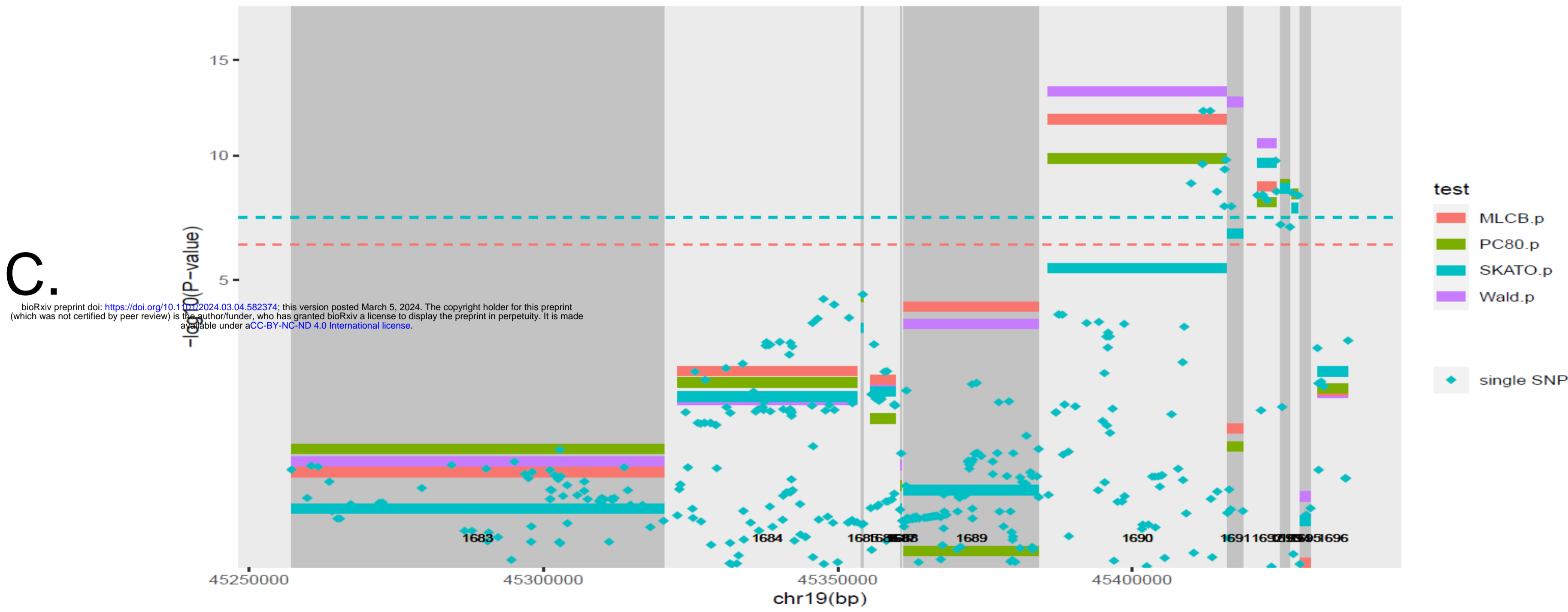
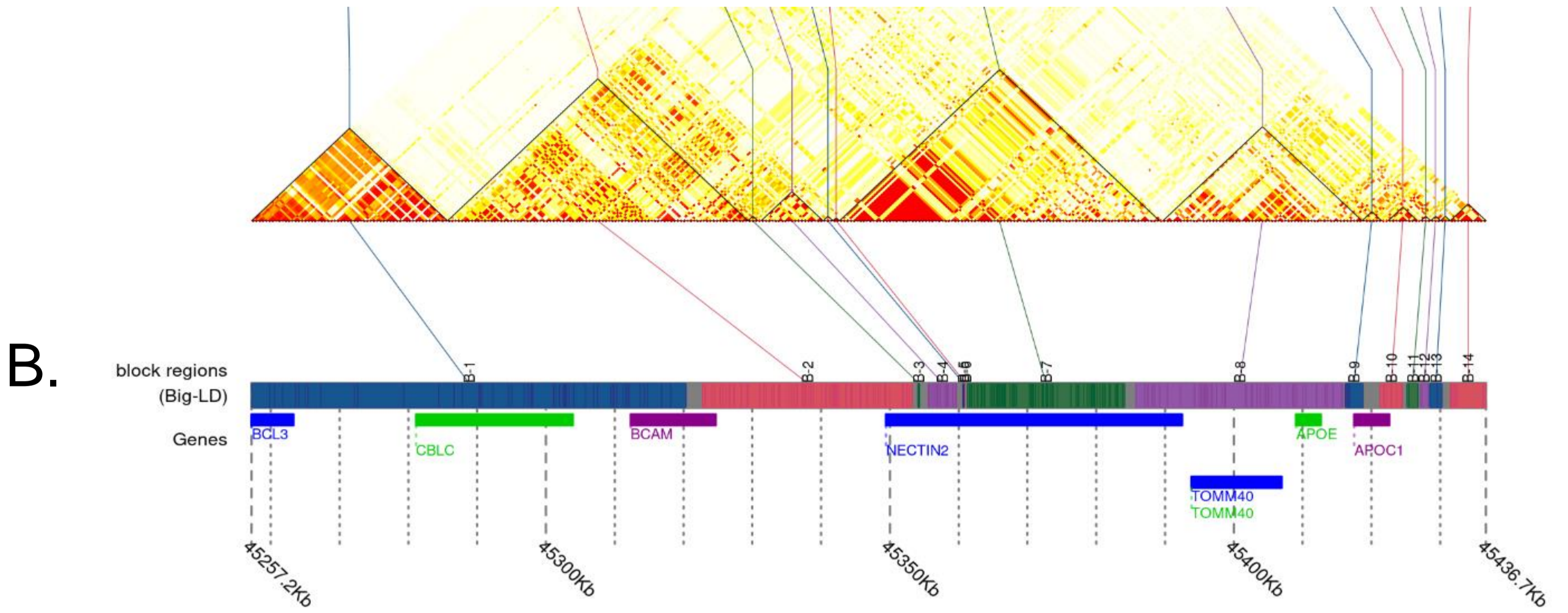
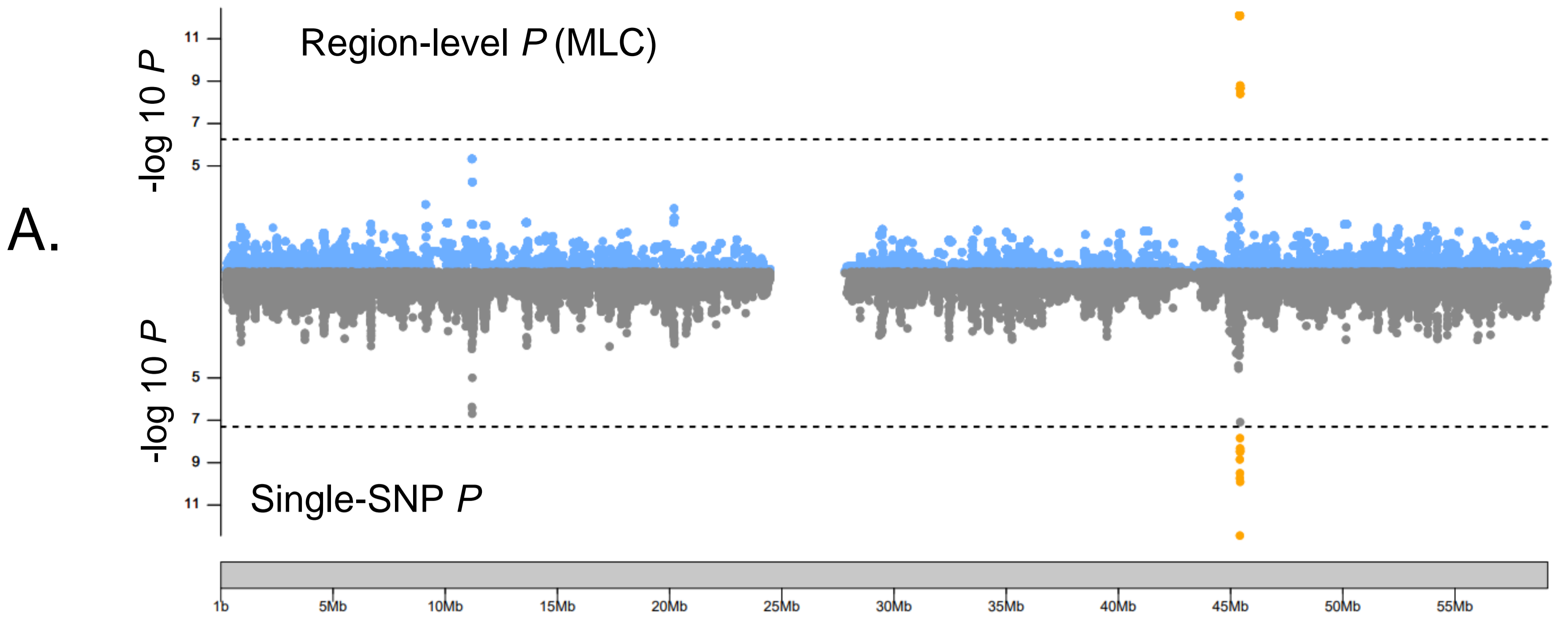
## 150 References

- 151 1. Neale, B. M. & Sham, P. C. The future of association studies: gene-based analysis and  
152 replication. *The American Journal of Human Genetics* **75**, 353–362 (2004).
- 153 2. Yoo, Y. J., Sun, L., Poirier, J. G., Paterson, A. D. & Bull, S. B. Multiple linear  
154 combination (MLC) regression tests for common variants adapted to linkage  
155 disequilibrium structure. *Genet Epidemiol* **41**, 108–121 (2017).
- 156 3. Gauderman, W. J., Murcray, C., Gilliland, F. & Conti, D. V. Testing association between  
157 disease and multiple SNPs in a candidate gene. *Genet Epidemiol* **31**, 383–395 (2007).
- 158 4. LI, Q. H. & LAGAKOS, S. W. On the Relationship between Directional and Omnibus  
159 Statistical Tests. *Scandinavian Journal of Statistics* **33**, 239–246 (2006).
- 160 5. Yoo, Y. J., Sun, L., Poirier, J. G., Paterson, A. D. & Bull, S. B. Multiple linear  
161 combination (MLC) regression tests for common variants adapted to linkage  
162 disequilibrium structure. *Genet Epidemiol* **41**, 108–121 (2017).
- 163 6. Ionita-Laza, I., Lee, S., Makarov, V., Buxbaum, J. D. & Lin, X. Sequence kernel  
164 association tests for the combined effect of rare and common variants. *Am J Hum Genet*  
165 **92**, 841–853 (2013).
- 166 7. Lee, S. *et al.* Optimal Unified Approach for Rare-Variant Association Testing with  
167 Application to Small-Sample Case-Control Whole-Exome Sequencing Studies. 224–237  
168 (2012) doi:10.1016/j.ajhg.2012.06.007.
- 169 8. Gao, X., Starmer, J. & Martin, E. R. A multiple testing correction method for genetic  
170 association studies using correlated single nucleotide polymorphisms. *Genet Epidemiol*  
171 **32**, 361–369 (2008).
- 172 9. Li, M. X., Gui, H. S., Kwan, J. S. H. & Sham, P. C. GATES: A rapid and powerful gene-  
173 based association test using extended Simes procedure. *Am J Hum Genet* **88**, 283–293  
174 (2011).
- 175 10. James, S. Approximate multinormal probabilities applied to correlated multiple endpoints  
176 in clinical trials. *Stat Med* **10**, 1123–1135 (1991).
- 177 11. Kim, S. A. *et al.* gpart: human genome partitioning and visualization of high-density SNP  
178 data by identifying haplotype blocks. *Bioinformatics* **35**, 4419–4421 (2019).
- 179 12. Kim, S. A., Cho, C.-S., Kim, S.-R., Bull, S. B. & Yoo, Y. J. A new haplotype block  
180 detection method for dense genome sequencing data based on interval graph modeling of  
181 clusters of highly correlated SNPs. *Bioinformatics* **34**, 388–397 (2018).
- 182 13. Danecek, P. *et al.* The variant call format and VCFtools. *Bioinformatics* **27**, 2156–2158  
183 (2011).



- 184 14. Pocock, S. J., Geller, N. L. & Tsiatis, A. A. The Analysis of Multiple Endpoints in  
185 Clinical Trials. *Biometrics* **43**, 487 (1987).
- 186 15. Stram, D. O., Wei, L. J. & Ware, J. H. Analysis of repeated ordered categorical outcomes  
187 with possibly missing observations and time-dependent covariates. *J Am Stat Assoc* **83**,  
188 631–637 (1988).
- 189 16. Kwee, L. C., Liu, D., Lin, X., Ghosh, D. & Epstein, M. P. A Powerful and Flexible  
190 Multilocus Association Test for Quantitative Traits. *The American Journal of Human*  
191 *Genetics* **82**, 386–397 (2008).
- 192 17. Wu, M. C. *et al.* Rare-variant association testing for sequencing data with the sequence  
193 kernel association test. *Am J Hum Genet* (2011) doi:10.1016/j.ajhg.2011.05.029.
- 194 18. Wang, P., Rahman, M., Jin, L. & Xiong, M. A new statistical framework for genetic  
195 pleiotropic analysis of high dimensional phenotype data. *BMC Genomics* **17**, 1–24 (2016).
- 196 19. Kosmidis, I., Kenne Pagui, E. C. & Sartori, N. Mean and median bias reduction in  
197 generalized linear models. *Stat Comput* **30**, 43–59 (2020).
- 198 20. Kosmidis, I. & Firth, D. Jeffreys-prior penalty, finiteness and shrinkage in binomial-  
199 response generalized linear models. *Biometrika* **108**, 71–82 (2021).
- 200 21. Yoo, Y. J., Kim, S. A. & Bull, S. B. Clique-Based Clustering of Correlated SNPs in a  
201 Gene Can Improve Performance of Gene-Based Multi-Bin Linear Combination Test.  
202 *Biomed Res Int* **2015**, 1–11 (2015).
- 203 22. Paterson, A. D. *et al.* A genome-wide association study identifies a novel major locus for  
204 glycemic control in type 1 diabetes, as measured by both A1C and glucose. *Diabetes* **59**,  
205 539–549 (2010).
- 206 23. Roshandel, D. *et al.* Meta-genome-wide association studies identify a locus on  
207 chromosome 1 and multiple variants in the MHC region for serum C-peptide in type 1  
208 diabetes. *Diabetologia* **61**, 1098–1111 (2018).
- 209





bioRxiv preprint doi: <https://doi.org/10.1101/2024.03.04.582374>; this version posted March 5, 2024. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-NC-ND 4.0 International license.