

1 **An exact, unifying framework for region-based association**  
2 **testing in family-based designs, including higher criticism**  
3 **approaches, SKATs, multivariate and burden tests**

4

5 **Julian Hecker<sup>1,2\*</sup>, F. William Townes<sup>3</sup>, Priyadarshini Kachroo<sup>1</sup>, Jessica Lasky-Su<sup>1</sup>, John Ziniti<sup>1</sup>, Michael**  
6 **H. Cho<sup>1</sup>, Scott T. Weiss<sup>1</sup>, Nan M. Laird<sup>2</sup>, Christoph Lange<sup>2</sup>**

7 <sup>1</sup>Channing Division of Network Medicine, Department of Medicine, Brigham and Women's

8 Hospital and Harvard Medical School, Boston, MA, USA

9 <sup>2</sup>Department of Biostatistics, Harvard T.H. Chan School of Public Health, Boston, MA, USA

10 <sup>3</sup>Department of Computer Science, Princeton University, Princeton, New Jersey, USA

11

12 \*Correspondence: [rejhe@channing.harvard.edu](mailto:rejhe@channing.harvard.edu)

13

14

15

16

17

18

19

20

21

22

23

24

25

## 1 **Abstract**

2 Analysis of rare variants in family-based studies remains a challenge. To perform a region/set-  
3 based association analysis of rare variants in family-based studies, we propose a general  
4 methodological framework that integrates higher criticism, maximum, SKATs, and burden  
5 approaches into the family-based association testing (FBAT) framework. Using the haplotype  
6 algorithm for FBATs to compute the conditional genotype distribution under the null  
7 hypothesis of Mendelian transmissions, virtually any association test statistics can be  
8 implemented in our approach and simulation-based or exact p-values can be computed  
9 without the need for asymptotic settings. Using simulations, we compare the features of the  
10 proposed test statistics in our framework with the existing region-based methodology for  
11 family-based studies under various scenarios. The tests of our framework outperform the  
12 existing approaches. We provide general guidelines for which scenarios, e.g., sparseness of  
13 the signals or local LD structure, which test statistic will have distinct power advantages over  
14 the others. We also illustrate our approach in an application to a whole-genome sequencing  
15 dataset with 897 asthmatic trios.

16

17

18

19

20

21

22

23

## 1 **Introduction**

2 In family-based association studies, the concept of Mendelian transmissions can be utilized to  
3 construct association tests that are robust against genetic confounding (Transmission  
4 Disequilibrium Tests (TDTs) <sup>1</sup> or Family-based Association Tests (FBATs) <sup>2</sup>). This feature of  
5 family-based association tests was fundamental to establish them as a popular tool in  
6 association mapping since the days of candidate gene studies. The robustness of the approach  
7 largely out-weighted the requirement of family data, i.e., having to recruit additional related  
8 study subjects, and reduced statistical power compared to population-based cohorts with the  
9 same sample size. Testing strategies have been proposed that allowed the incorporation of  
10 the association information at the population-level in family-based designs without  
11 compromising the robustness of the test statistic <sup>3-7</sup>.

12 As genome-wide association studies (GWAS) became a standard research tool, and the  
13 multiple testing problem had to be addressed at a genome-wide level, researchers started to  
14 emphasize statistical power in their choice for study designs. Power and study design  
15 considerations substantially contributed to making population-based designs approaches the  
16 most popular choice in association analysis.

17 Now, as whole-genome sequencing (WGS) studies are replacing chip-based GWAS, region-  
18 based rare variant analysis approaches have moved to the center of the statistical  
19 methodology development, as, even for very large sample sizes, the power of single locus  
20 association analysis will be too small when the minor allele is rare. Region-based approaches  
21 are motivated by the idea that, if we can combine "association signals" across a pre-defined  
22 region in a suitable way, a stronger genetic signal could be assessed by a suitable test statistic,  
23 and the resulting region-based association test would have increased statistical power. At the  
24 same time, the multiple testing problem could become less severe as fewer association tests

1 are computed. The major statistical challenges for region-based tests are to identify suitable  
2 ways to combine the genetic information across the pre-specified region, to  
3 incorporate/estimate the correlation between the selected rare variants in the region-based  
4 test statistic, and to select a suitable test statistic.

5 For population-based designs, based on assumptions about the alternative  
6 hypothesis/distribution of disease susceptibility loci (DSLs) and their effect directions,  
7 numerous region-based tests have been proposed, e.g., burden tests and variance component  
8 tests<sup>8,9</sup>. However, population stratification is a potential problem in population-based designs  
9 that can be even more severe in the analysis of rare variants.

10 For family-based designs, the popular burden tests and the SKAT approach have been  
11 translated to the FBAT framework<sup>10,11</sup>. As with their population-based equivalents, these two  
12 approaches estimate the correlation between the genetic loci empirically. This, especially for  
13 rare variant data, can be problematic, as the rare variant allele counts are small, and the  
14 empirical estimates can be affected by numerical instabilities. Furthermore, the application of  
15 asymptotic theory may not provide accurate results here. Other recent approaches are based  
16 on mixed models and can, theoretically, analyze unrelated and related samples<sup>12</sup>.

17 As an other alternative to the transmission-based FBAT approach, methods that compare the  
18 allele frequencies between affected and unaffected individuals within families have been  
19 suggested, e.g., Generalized Disequilibrium Test (GDT) for single variant analysis<sup>13</sup>. As they  
20 can incorporate all available phenotypic and genetic information within one pedigree, they  
21 can be more powerful than the corresponding FBAT analysis. However, they require the  
22 assumption that the allele frequencies and the genetic variance for all members of a pedigree  
23 are equal under the null hypothesis. This assumption can be violated in the presence of

1 population substructure within the founders of the families or when there are departures  
2 from Hardy Weinberg equilibrium <sup>14</sup>. None of these assumptions is required for the FBAT  
3 approach to be valid. For region-based analysis, the Rare-Variant Generalized Disequilibrium  
4 Test (RV-GDT) extension has been proposed <sup>15</sup>.

5 In this communication, we will exclusively focus on transmission-based analysis approaches  
6 to construct association tests that are robust against confounding. We propose a general  
7 framework for region-based rare variant analysis in extended pedigree/nuclear families that  
8 is based on the FBAT approach. Multiple offspring per family may be available, founder/phase  
9 information can be missing, and phenotypes can be dichotomous or quantitative. In contrast  
10 to previous approaches for region-based analysis in population- and family-based designs, the  
11 joint distribution of the rare variants in the region is obtained analytically under the null  
12 hypothesis, conditioning on the sufficient statistic, using the haplotype algorithm for FBATs  
13 <sup>16,17</sup>. Based on the conditional genotype distribution, it is straightforward to implement region-  
14 based FBATs, e.g., multivariate tests, burden tests, SKAT <sup>9,11</sup>, and higher criticism approaches  
15 <sup>18-20</sup>. As the joint distribution of the rare variants is obtained analytically and can efficiently be  
16 sampled from, the significance of the test statistics in our framework can be obtained either  
17 by simulations or the construction of the exact distribution<sup>21</sup>. This flexibility of our approach  
18 enables the implementation of virtually any region/set-based test without the need for any  
19 asymptotic assumptions or approximations. We illustrate the implementation of higher  
20 criticism approaches, maximum statistics, SKATs and burden tests in our framework.

21 For different scenarios, e.g., regions with sparse signals, varying local Linkage Disequilibrium  
22 (LD) structure, we compare our proposed FBAT framework to existing methodology, using  
23 extensive simulation studies. Our simulation results support our theoretical considerations  
24 that our testing framework provides a substantial improvement over the existing

1 methodology in terms of statistical power and robustness against population substructure.  
2 We also develop general recommendations for which choice of the test statistic is preferable  
3 for which scenario. Furthermore, we also applied our methodology framework to a whole-  
4 genome sequencing study for childhood asthma with 897 trios.

## 5 **Methods**

6 In a family-based WGS association study, genotype data for rare variants are available for a  
7 set of marker loci that are in close physical proximity and define a genomic segment that is  
8 suitable for region-based association analysis. The genotype information may be available for  
9 multiple offspring as well as for the parents. For the  $i$ -th nuclear family, we introduce the  
10  $p \times n_i$  genotype matrix  $X_i$  and the  $n_i$  dimensional phenotype vector  $T_i$ , where  $n_i$  denotes the  
11 number of offspring in the  $i$ -th nuclear family, and  $p$  denotes the number of variants in the  
12 analysis region. We regard  $X_i$  as random while  $T_i$  is fixed in the FBAT approach. Below, we  
13 propose a possible set of test statistics that can capture the potential association between the  
14 offspring genotype data and the phenotypes under various conditions.

## 15 **Simulation-based significance testing**

16 For each region, using the haplotype algorithm for FBAT<sup>16,17</sup>, we derive the conditional  
17 distribution of offspring genotypes  $X_i$  in the  $i$ -th nuclear family under the null hypothesis,  
18 given the sufficient statistic  $S_i$  for the possible missing founder genotypes<sup>22</sup>. The sufficient  
19 statistic approach utilizes parental genotypes if they are available. The knowledge about the  
20 conditional genotype distribution allows constructing association tests that are robust against  
21 population stratification and admixture. Based on this conditional genotype distribution, it is  
22 straightforward to compute the first two moments of commonly used test statistics under the  
23 null hypothesis of no association and derive the asymptotic distribution. However, as the

1 analysis of rare variant data leads to scenarios where the application of asymptotic theory  
2 does not provide reliable approximations, simulation-based or even exact p-values<sup>21</sup> are  
3 preferred. Here, we propose to evaluate association p-values based on a sufficiently large  
4 number of simulated draws from the null distribution. This procedure can be combined with  
5 adaptive permutation/simulation-based p-value techniques. In this context, we recommend  
6 using stopping rules that are nearly optimal in terms of the required number of simulations<sup>23</sup>.  
7 This approach, therefore, also allows our testing framework to incorporate test statistics  
8 where the (asymptotic) distribution is intractable or cumbersome, e.g., maximum statistics or  
9 higher criticism approaches.

## 10 **Test statistics**

11 All test statistics under consideration are based on the following two objects. For the  $i$ -th  
12 family, we define the  $p$ -dimensional vector of Mendelian residuals  $U_i = (X_i - E[X_i|S_i])T_i$ .  
13 Also, we define the corresponding  $p \times p$  variance matrix  $V_i = Var(U_i|S_i, T_i)$ . For both  
14 objects, the moments are computed under the null hypothesis, based on the sufficient statistic  
15  $S_i$ .

16

## 1 **Burden-based approaches**

2 Burden-type FBATs can be implemented by specifying a  $p$ -dimensional weight vector  $W$  that  
3 collapses/summarizes the rare variant information of the region into a single scalar value. The  
4 specification of the weight vector  $W$  requires assumptions about the effect direction and its  
5 effect size. In this context, the contribution to the FBAT statistic of the  $i$ -th family is then given  
6 by

$$7 \quad U_i^* = W^T (X_i - E[X_i | S_i]) T_i$$

8 The corresponding FBAT-statistic for the simulation-based testing is computed by  
9  $FBAT_{burden} = (\sum_i U_i^*)^2$ . We note that for this burden test, it would be possible as well to  
10 compute an asymptotic p-value by also computing/estimating the corresponding variance.

## 11 **Variance component/SKAT approaches**

12 As an alternative to burden/collapsing association tests, SKAT/variance-component based  
13 region tests have been developed for rare variant data <sup>9</sup>. They have the advantage that they  
14 do not require any assumptions about the effect configuration at the rare variant loci under  
15 the alternative hypothesis, but they are not as powerful as burden/collapsing approaches if  
16 one is certain about the alternative hypothesis. We define the general statistic

$$17 \quad FBAT_{vc} = U^T W U$$

18 where  $U = \sum_i U_i$  and  $W$  is a fixed  $p \times p$  weight matrix. While, in the scenario of affected  
19 offspring trios and a diagonal weight matrix  $W$ , this test statistic equals the FB-SKAT statistic  
20 <sup>11</sup>, Ionita-Laza et al. assess significance based on asymptotic results that require the empirical  
21 estimation of the variance/covariance matrix for the rare variants, which, given the sparseness



1 of rare variant data, can become problematic. In our framework, the p-value of the test  
2 statistic is obtained based on simulations from the conditional genotype distribution.

3 If we set  $W = V^{-1}$ , where  $V = \sum_i V_i$ , we obtain the multivariate FBAT <sup>24</sup>. The multivariate  
4 FBAT was designed for common variants, and the asymptotic p-values also require an  
5 empirical estimate of the correlation matrix. Again, for rare variants, this can lead to unreliable  
6 results, making the implementation of the multivariate FBAT in our proposed framework  
7 preferable.

### 8 **Higher criticism and maximum statistic**

9 Besides the commonly used burden and variance component approaches, we introduce the  
10 higher criticism and maximum statistic for region-based analysis in family-based studies. Both  
11 approaches are designed to identify sparse alternatives and have been introduced to genetic  
12 association studies of unrelated individuals recently <sup>18,20</sup>.

13 Define the normalized residuals  $\frac{U_j}{\sqrt{V_{jj}}}$ ,  $j = 1, \dots, p$  and denote the corresponding association p-  
14 value based on the asymptotic marginal normal distribution by  $q_j$ . Based on the available  
15 amount of information per variant, e.g., the number of informative transmissions/families, we  
16 restrict the set of variants to a subset of variants where the marginal variance is large enough  
17 (e.g., we require at least 5 informative nuclear families). Denote the number of variants in this  
18 subset by  $p'$ . Given the ordered p-values  $q_{(1)} \leq q_{(2)} \leq \dots \leq q_{(p')}$ , we define the HC statistic  
19 as

$$20 \quad FBAT_{HC} = \max_j \frac{\frac{j}{p'} - q_{(j)}}{\sqrt{q_{(j)}(1 - q_{(j)})}}$$

1 Here, the index set  $J$  can be  $\{1, \dots, \frac{p'}{2}\}$  or  $\{1, \dots, p'\}$ , depending on assumptions about the  
2 underlying genetic architecture.

3 It is important to note that, while  $FBAT_{HC}$  contains a transformation based on the single  
4 variant asymptotic distribution, the assessment of its significance based on simulations from  
5 the conditional distribution remains a valid approach regardless of whether the assumptions  
6 that motivated the transformation hold.

7 The second approach to detect sparse signals in the tested genomic region is the MAX statistic  
8 which is simply defined as

9 
$$FBAT_{MAX} = \max_{1 \leq j \leq p} \left| \frac{U_j}{\sqrt{V_{jj}}} \right|$$

10 The theory related to the higher criticism/max statistic in the setting of unrelated case-control  
11 data and sparse signals developed in Mukherjee et al. <sup>20</sup> can be transferred to family-based  
12 studies/FBATS. In Appendix A, we derive how the theory in Mukherjee et al. <sup>20</sup> can be applied  
13 to the FBAT framework in the scenario of affected offspring trios. The corresponding  
14 optimality results for sparse signal scenarios motivate the application of these test statistics  
15 to rare variants in sequencing studies.

16 Finally, we note that it is, of course, possible to set up an omnibus statistic that is based on  
17 the maximum of multiple test statistics described above.

## 18 **Results**

19 In this section, we describe the results of two simulation studies and an analysis of a whole-  
20 genome sequencing study for childhood-asthma with 897 affected offspring trios.

21

## 1 **Simulation studies**

2 We studied the performance of our proposed test statistics in two extensive simulation  
3 studies. In both studies, we compared the Type I error and power with the existing  
4 methodology for family-based region association analysis. We restricted all simulations to the  
5 scenario of trios with an affected offspring. However, it is important to note that our  
6 framework can be applied to any nuclear family and phenotype distribution. For the test  
7 statistics  $FBAT_{burden}$  and  $FBAT_{vc}$ , we applied uniform weights. In the following, we will  
8 denote the test statistics  $FBAT_{burden}$ ,  $FBAT_{vc}$ ,  $FBAT_{HC}$ , and  $FBAT_{MAX}$  by Burden, SKAT, HC,  
9 and MAX.

## 10 **Genetic regions with unphased data**

11 We extracted haplotypes for the CEU and the GBR subpopulations from the 1000 Genomes  
12 Project <sup>25</sup>, consisting of 30 and 50 consecutive rare variants with a minor allele frequency  
13 (MAF) below 3%. Based on these haplotypes, we generated genotype data for trios using  
14 Mendelian transmissions. Using a standard logistic disease model with a disease prevalence  
15 of  $\approx 10\%$ , we simulated offspring affection status and collected  $n = 1,000$  affected offspring  
16 trios. This simulation study is similar to the simulation studies described in the existing  
17 literature <sup>15,26</sup>.

18 We compared our test statistics with the GTDT <sup>26</sup> and the RV-GDT <sup>15</sup>. The GTDT <sup>26</sup> offers five  
19 different test statistics for region-based affected offspring trio analysis, designed for different  
20 modes of inheritance. The test statistics require phased haplotype data. If the phase  
21 information is not available, this information is reconstructed up to small uncertainties. We  
22 considered the test statistics GTDT-AD, GTDT-DOM, and GTDT-CH in our study. The RV-GDT <sup>15</sup>  
23 describes a generalization of the single variant GDT <sup>13</sup> for multiple variants in a genetic region.  
24 The RV-GDT can be applied to arbitrary pedigrees where affected, and unaffected samples are

1 collected; members can be missing. The test statistic compares the genotype counts between  
2 affected and unaffected members and corrects for the relatedness using the  
3 estimated/reported kinship coefficients. We note that this implies that the phenotype  
4 information for parents must be available, whereas the classical TDT/FBAT test for offspring  
5 trios does not require this information. For comparison, we included the test MAX-BF that  
6 tests if at least one single variant FBAT statistic reached the Bonferroni-corrected significance  
7 level corresponding to the number of variants  $p$  in the region. The corresponding single variant  
8 p-values were evaluated using asymptotic theory due to computational reasons.

9 To check the Type I error rates and the robustness against population stratification; we  
10 considered a null hypothesis simulation where no genetic variant is associated with the  
11 affection status and three different population admixture scenarios (Table 1). In these  
12 admixture scenarios, we generated one fixed parent based on the CEU haplotypes and the  
13 other parent based on the GBR haplotypes. The affection status of the parents differed across  
14 the three admixture scenarios. For the power analysis, we simulated six different scenarios  
15 where the number of causal variants and corresponding effect sizes differ (Table 2). In scenario  
16 5, we picked very rare and independent causal variants with a MAF below 1%, and in scenario  
17 6 we chose causal variants that are in strong LD with multiple other variants. All results are  
18 based on 1,000 replicates.

19 In Table 1, we observe that all methods control the Type I error appropriately. The only  
20 exception is the RV-GDT in the scenario of population admixture with discordant parental  
21 phenotypes (adm2 and adm3, Table 1). This is expected, as the GDT/RV-GDT test compares  
22 the frequencies between affected and unaffected family members and cannot distinguish  
23 between association and admixture in the parents. We also note that the RV-GDT computes  
24 a one-sided p-value, which explains the deflation/inflation behavior, depending on the

1 parental phenotypes. The power results in Table 2 demonstrate the advantages of non-burden  
2 tests in specific scenarios. The power results are also illustrated in Figures 1 and 2. The SKAT  
3 statistic shows the highest power in the first three scenarios and outperforms the other tests.  
4 However, the MAX and HC statistics also show substantial power. The results for scenario 4  
5 are comparable between SKAT, MAX, and HC. In scenario 5, the HC statistics achieves the  
6 highest power, which is supported by our theoretical considerations as well (see Appendix A).  
7 In the last scenario 6, all tests achieve substantial power, as expected, due to the LD structure  
8 that pushes power. The most powerful tests here are SKAT and RV-GDT, but MAX and HC test  
9 statistics achieve similar results. If we have different effect directions (scenario 2 and 4), the  
10 burden test loses power compared to the consistent effect direction scenarios 1 and 3, which  
11 is expected. The FBAT Burden test and GDT-AD have almost no power in scenarios 4 and 5. It  
12 is important to note that the GTDT-AD and the FBAT burden test are essentially based on the  
13 same test statistic idea, the only difference lies in the fact that the GTDT assigns haplotypes  
14 (with possible error) and our approach uses the robust conditional genotype distribution  
15 computed by the FBAT haplotype algorithm.

16 There is a substantial difference between the power of the MAX statistic and the MAX-BF  
17 statistic in all scenarios, because the Bonferroni correction is conservative, and the p-value of  
18 the MAX statistic is evaluated based on the joint conditional genotype distribution. We note  
19 that this difference could be even larger if the genome-wide significance levels for region- and  
20 single variant based testing are considered.

21  
22  
23  
24

## 1 **Dense genetic regions with phased data**

2 For our second set of simulation studies, we utilized the 1006 EUR population haplotypes from  
3 the 1000 Genomes Project along 1,000 consecutive rare genetic variants with MAF below 3%.  
4 In this simulation study, we consider a large number of variants in combination with a sparse  
5 signal, which means a small subset of causal variants that are not in strong LD with any other  
6 variants. We simulated affected offspring trios as described in the first simulation study but  
7 also stored the phased haplotypes for all members of the trio. In the scenario where the  
8 haplotypes are observed, the conditional distribution identified by the FBAT haplotype  
9 algorithm equals the distribution where both parents transmit one of the observed haplotypes  
10 with equal probability of 0.5. We compared the performance of the FBAT, the GTDT, and the  
11 RV-TDT BRV<sup>27</sup> statistics to demonstrate the potential advantage of non-burden tests in the  
12 presence of sparse signals. Again, we also included the MAX-BF test, where we considered the  
13 Bonferroni corrected significance level based on  $p = 1,000$  tests. As mentioned above, the  
14 knowledge about the phased haplotypes is the preferred setting for the GTDT. Also, the RV-  
15 TDT requires phased haplotypes<sup>27</sup>. We considered a null hypothesis scenario and four  
16 different power scenarios (Table 3 and Figure 3). For the null hypothesis simulation and the  
17 first two power scenarios we simulated 1,000 trios; the last two power scenarios are based on  
18 10,000 trios. The four power scenarios include causal variants that are in almost no LD with  
19 other variants, and the number of causal variants is small compared to the overall number of  
20  $p = 1,000$  variants. All results are based on 1,000 replicates, and the p-values for all test  
21 statistics were evaluated empirically based on the same 1,000 draws from the conditional  
22 haplotype distribution. The results for this simulation are also visualized in Figure 3.  
23 In Table 3, we observe that all test statistics control the Type I error appropriately. Since all p-  
24 values are evaluated empirically based on the conditional haplotype distribution by

1 simulation, this is expected. In the first power scenario 1, the MAX test statistic achieves the  
2 highest power as we simulated a sparse and rare, but strong signal in the genetic region,  
3 consisting of two rare variants. The HC test statistic also achieves substantial power, whereas  
4 all other tests (except MAX-BF) are almost powerless in this scenario. In the second scenario,  
5 where the MAF of the two causal variants is much higher, the MAX test statistics still  
6 outperforms the other tests, but also the SKAT and the HC test statistics show good  
7 performances. In scenario 3, where many very rare causal variants have a relatively small  
8 effect size, the HC is the most powerful test. This is in line with the results in Mukherjee et al.  
9 <sup>20</sup> that describe a lower detection boundary in the mild sparse regime compared to the MAX  
10 test statistics. However, the RV-TDT BRV and MAX test statistic also achieve substantial power  
11 in this scenario. The power behavior differs more in the last scenario 4, where the effects are  
12 pointing in different directions. Here, the MAX and HC statistics have a significantly increased  
13 power compared to the other tests, while the HC test statistic is the most powerful one. We  
14 note that the FBAT burden and the GTDT-AD test are very close to the nominal level in  
15 scenarios 1, 2, and 4. Both tests are equivalent since the test statistics are the same.  
16 Overall, again, the MAX test shows higher power than the MAX-BF tests, as described in the  
17 context of the first simulation study.

## 18 **Real data analysis**

19 To demonstrate the applicability and the advantages of our proposed framework, we analyzed  
20 a whole-genome sequencing dataset consisting of 897 complete asthmatic trios from Costa  
21 Rica <sup>28</sup>. After standard quality control, including Mendelian error rates, we excluded all  
22 variants with a MAF above 5%. The resulting 27,345,734 non-monomorphic variants were  
23 partitioned into approximately 547,000 consecutive windows of 50 rare variants. Other  
24 partitioning approaches could be considered here <sup>29,30</sup>, but, as the focus of this data analysis

1 was to demonstrate the feasibility of our approach, we did not explore different window-  
2 strategies here.

3 For each window, we performed the Burden, SKAT, MAX, and HC test, using affection status  
4 as the phenotype. We evaluated the p-values by simulation, where we used an adaptive  
5 heuristic that increases the number of simulations if the estimated p-value is close to the  
6 minimum possible value. The smallest possible p-value was  $p = 10^{-9}$ , since the number of  
7 simulations was truncated at  $10^9$ . In Figure 4, we plotted the corresponding quantile-quantile-  
8 plot. The plot indicates that the test statistics control the Type I error rate but can also identify  
9 potential findings.

10 Based on the approximately  $4 * 547,000$  tests and a False Discovery Rate at  $\alpha = 0.05$ <sup>31</sup>, our  
11 approach identified three single significant regions on Chromosome 1, 12, and 21, as well as  
12 multiple consecutive significant regions on Chromosome 10. The significance of the three  
13 regions on Chromosomes 1, 12, and 21 was declared by the Burden test, whereas the single  
14 variant FBAT p-values within the regions were not in the range of genome-wide significance.  
15 The other regions on Chromosome 10 were identified by the MAX, HC, and SKAT tests. The  
16 lowest p-value of  $10^{-9}$  was reached by the SKAT test. For all these regions, the Burden test  
17 did not reach the magnitude of genome-wide significance. This shows the benefits of  
18 combining different test statistics to identify distinct genetic signal structures.

19

20

## 21 **Discussion**

22 In this manuscript, we propose a general framework for the region-based association analysis  
23 of sequencing datasets with family-based designs. The framework incorporates burden tests,  
24 SKAT, maximum and higher criticism approaches, and, given the flexibility of the framework,



1 any future approach can straightforwardly be implemented. In contrast to previously  
2 published approaches, the joint genotype distribution along the loci is not obtained by  
3 empirical estimates, but via the haplotype algorithm for FBATs<sup>16,17</sup>. This allows our proposed  
4 testing framework for FBATs to assess the significance of an arbitrary test statistic based on  
5 simulations or based on the exact distribution. This approach is enabled by the recent  
6 improvements in the FBAT haplotype algorithm<sup>16</sup>, which reduces the computational burden  
7 of the original approach by several magnitudes. Our simulation results illustrate that the  
8 optimal test for the region-based analysis depends on the specific genetic architecture of the  
9 disease, and any WGS analysis relying on just one single test statistic may not detect all  
10 associations contained in the data. While dense signals with consistent effect directions can  
11 be captured by burden tests, different effect directions and less dense signals can be identified  
12 by SKAT approaches. If the signal becomes more separated and sparser, the MAX and HC  
13 approaches can be the most powerful tests.

14 The proposed implementation of the simulation-based p-values requires the user to pre-select  
15 the number of simulations that FBAT performs for each test. The computational burden can  
16 be decreased by adaptive strategies<sup>23</sup>. The applications of the proposed analysis framework  
17 to simulated and real data illustrate that the theoretically expected advantages are also of  
18 practical relevance and that simulation-based p-values are not prohibitive in WGS settings. A  
19 subject of future research will be to integrate the existing FBAT approaches to multivariate  
20 phenotypes, longitudinal data, age at onset<sup>32,32-34</sup>, gene-environmental interactions, and  
21 testing strategies into the proposed framework<sup>3,6,7</sup>.

22

23

## 1 **Appendix A: Detection of sparse signals**

2 We consider the scenario of an affected offspring trio. Both parental genotypes are observed  
3 along with the  $p$  variants in the analysis region. As noted in Chen et al. <sup>26</sup>, if there is no variant  
4 where all three observed genotypes (mother, father, offspring) are heterozygous, the phase  
5 information can be recaptured from the observed unphased genotype data. However, as  
6 described in Hecker et al. <sup>14</sup>, treating inferred haplotypes as observed haplotypes can lead to  
7 misspecification.

8 Nevertheless, more specifically, if there is no variant for which both parental genotypes are  
9 heterozygous, haplotypes can be phased, and the resulting conditional genotype distribution  
10 obtained by the FBAT haplotype algorithm equals the conditional distribution where we treat  
11 the haplotypes as observed. If we restrict the genetic data to rare variants, this is true for most  
12 nuclear families. In addition, with relatively high probability, at least 1 parent has only 1 minor  
13 allele in the genetic region.

14 Let us denote the phased parental mating type for such a trio by  $G = (h_1^M, h_2^M) \times (h_1^F, h_2^F)$ . The  
15 possible offspring genotypes are denoted by  $X_1 = (h_1^M + h_1^F)$ ,  $X_2 = (h_1^M + h_2^F)$ ,  $X_3 =$   
16  $(h_2^M + h_1^F)$  and  $X_4 = (h_2^M + h_2^F)$ . We assume the following, commonly used, disease model  
17 that describes the conditional offspring genotype distribution

$$18 \quad P(X_i | T = 1, G) = \frac{\exp(\beta^T X_i)}{\sum_{j=1}^4 \exp(\beta^T X_j)}$$

19 where the  $p$  dimensional vector  $\beta$  describes the genetic effects of the variants in the region.

20 If we denote the inherited offspring haplotypes by  $(h^M, h^F)$ , this model factors into the  
21 product of the two likelihoods

$$1 \quad P(h^g = h_j^g | T = 1, G) = \frac{\exp(\beta^T h_j^g)}{\exp(\beta^T h_1^g) + \exp(\beta^T h_2^g)}, \quad j = 1, 2, \quad g = M, F$$

2 Since the haplotype data is sparse as described above, this setting matches the scenario of  
3 Weakly Correlated Designs that is described in the paper by Mukherjee et al.<sup>20</sup> about sparse  
4 binary regression (Definition 4.1). They showed that in the sparse regime, the higher  
5 criticism and the maximum statistic can identify sparse alternatives efficiently (see Theorem  
6 7.4).

7 Although this motivates the application to affected offspring trios, the statistics can be  
8 applied in all FBAT scenarios. The Type I error is preserved because we utilize a simulation-  
9 based approach.

10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21

## 1 **Acknowledgments**

2 This work was supported by Cure Alzheimer's Fund; the National Human Genome Research  
3 Institute [R01HG008976]; and the National Heart, Lung, and Blood Institute [U01HL089856,  
4 U01HL089897, P01HL120839, P01HL132825].

## 5 **Declaration of Interests**

6 The authors declare no competing interests.

## 7 **Web Resources**

8 The FBAT software is available at <https://sites.google.com/view/fbat-web-page>. A new version  
9 that implements the described methodology is currently in preparation and will be available  
10 soon.

11

12

13

14

15

16

17

18

19

20

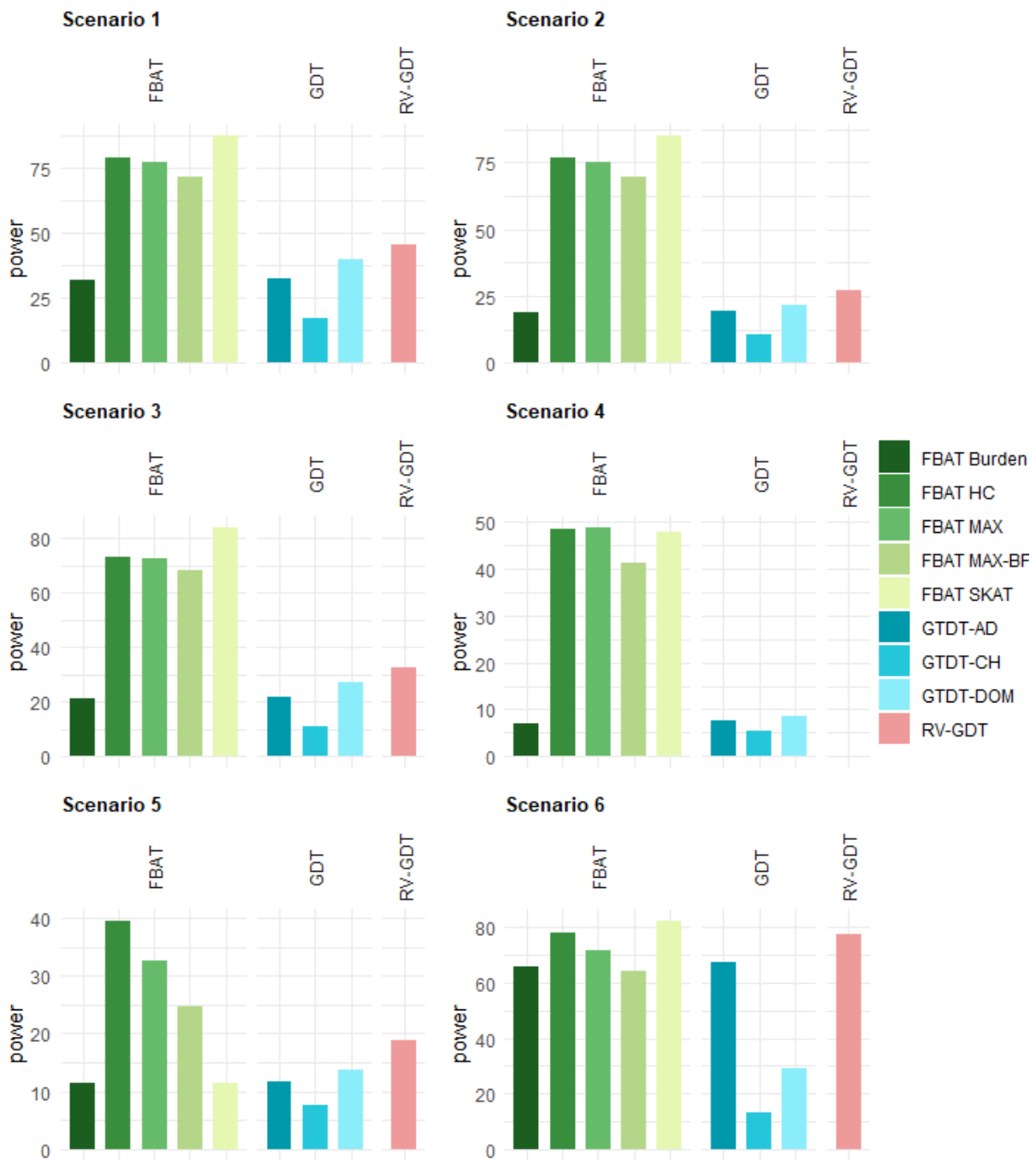
## 1 References

- 2 1. Spielman, R.S., McGinnis, R.E., and Ewens, W.J. (1993). Transmission test for linkage  
3 disequilibrium: the insulin gene region and insulin-dependent diabetes mellitus (IDDM). *Am.*  
4 *J. Hum. Genet.* *52*, 506–516.
- 5 2. Laird, N.M., and Lange, C. (2006). Family-based designs in the age of large-scale gene-  
6 association studies. *Nat. Rev. Genet.* *7*, 385–394.
- 7 3. Ionita-Laza, I., McQueen, M.B., Laird, N.M., and Lange, C. (2007). Genomewide weighted  
8 hypothesis testing in family-based association studies, with an application to a 100K scan.  
9 *Am. J. Hum. Genet.* *81*, 607–614.
- 10 4. Lange, C., Lyon, H., DeMeo, D., Raby, B., Silverman, E.K., and Weiss, S.T. (2003). A New  
11 Powerful Non-Parametric Two-Stage Approach for Testing Multiple Phenotypes in Family-  
12 Based Association Studies. *Hum. Hered.* *56*, 10–17.
- 13 5. Murphy, A., Weiss, S.T., and Lange, C. (2008). Screening and Replication using the Same  
14 Data Set: Testing Strategies for Family-Based Studies in which All Proband's Are Affected.  
15 *PLOS Genet.* *4*, e1000197.
- 16 6. Steen, K.V., McQueen, M.B., Herbert, A., Raby, B., Lyon, H., DeMeo, D.L., Murphy, A., Su,  
17 J., Datta, S., Rosenow, C., et al. (2005). Genomic screening and replication using the same  
18 data set in family-based association testing. *Nat. Genet.* *37*, 683–691.
- 19 7. Won, S., Wilk, J.B., Mathias, R.A., O'Donnell, C.J., Silverman, E.K., Barnes, K., O'Connor,  
20 G.T., Weiss, S.T., and Lange, C. (2009). On the Analysis of Genome-Wide Association Studies  
21 in Family-Based Designs: A Universal, Robust Analysis Approach and an Application to Four  
22 Genome-Wide Association Studies. *PLOS Genet.* *5*, e1000741.
- 23 8. Li, B., and Leal, S.M. (2008). Methods for detecting associations with rare variants for  
24 common diseases: application to analysis of sequence data. *Am. J. Hum. Genet.* *83*, 311–321.
- 25 9. Wu, M.C., Lee, S., Cai, T., Li, Y., Boehnke, M., and Lin, X. (2011). Rare-variant association  
26 testing for sequencing data with the sequence kernel association test. *Am. J. Hum. Genet.*  
27 *89*, 82–93.
- 28 10. De, G., Yip, W.-K., Ionita-Laza, I., and Laird, N. (2013). Rare Variant Analysis for Family-  
29 Based Design. *PLOS ONE* *8*, e48495.
- 30 11. Ionita-Laza, I., Lee, S., Makarov, V., Buxbaum, J.D., and Lin, X. (2013). Family-based  
31 association tests for sequence data, and comparisons with population-based association  
32 tests. *Eur. J. Hum. Genet. EJHG* *21*, 1158–1162.
- 33 12. Zhou, W., Nielsen, J.B., Fritsche, L.G., LeFaive, J., Taliun, S.A.G., Bi, W., Gabrielsen, M.E.,  
34 Daly, M.J., Neale, B.M., Hveem, K., et al. (2019). Scalable generalized linear mixed model for  
35 region-based association tests in large biobanks and cohorts. *BioRxiv* 583278.
- 36 13. Chen, W.-M., Manichaikul, A., and Rich, S.S. (2009). A Generalized Family-Based  
37 Association Test for Dichotomous Traits. *Am. J. Hum. Genet.* *85*, 364–376.

- 1 14. Hecker, J., Laird, N., and Lange, C. (2019). A comparison of popular TDT-generalizations  
2 for family-based association analysis. *Genet. Epidemiol.* *43*, 300–317.
- 3 15. He, Z., Zhang, D., Renton, A.E., Li, B., Zhao, L., Wang, G.T., Goate, A.M., Mayeux, R., and  
4 Leal, S.M. (2017). The Rare-Variant Generalized Disequilibrium Test for Association Analysis  
5 of Nuclear and Extended Pedigrees with Application to Alzheimer Disease WGS Data. *Am. J.*  
6 *Hum. Genet.* *100*, 193–204.
- 7 16. Hecker, J., Xu, X., Townes, F.W., Fier, H.L., Corcoran, C., Laird, N., and Lange, C. (2017).  
8 Family-based tests for associating haplotypes with general phenotype data. *Genet.*  
9 *Epidemiol.* *42*, 123–126.
- 10 17. Horvath, S., Xu, X., Lake, S.L., Silverman, E.K., Weiss, S.T., and Laird, N.M. (2004). Family-  
11 based tests for associating haplotypes with general phenotype data: application to asthma  
12 genetics. *Genet. Epidemiol.* *26*, 61–69.
- 13 18. Barnett, I., Mukherjee, R., and Lin, X. (2017). The Generalized Higher Criticism for Testing  
14 SNP-Set Effects in Genetic Association Studies. *J. Am. Stat. Assoc.* *112*, 64–76.
- 15 19. Donoho, D., and Jin, J. (2004). Higher criticism for detecting sparse heterogeneous  
16 mixtures. *Ann. Stat.* *32*, 962–994.
- 17 20. Mukherjee, R., Pillai, N.S., and Lin, X. (2015). HYPOTHESIS TESTING FOR HIGH-  
18 DIMENSIONAL SPARSE BINARY REGRESSION. *Ann. Stat.* *43*, 352–381.
- 19 21. Schneiter, K., Degnan, J.H., Corcoran, C., Xu, X., and Laird, N. (2007). EFBAT: exact family-  
20 based association tests. *BMC Genet.* *8*, 86.
- 21 22. Rabinowitz, D., and Laird, N. (2000). A unified approach to adjusting association tests for  
22 population admixture with arbitrary pedigree structure and arbitrary missing marker  
23 information. *Hum. Hered.* *50*, 211–223.
- 24 23. Hecker, J., Ruczinski, I., Cho, M., Silverman, E., Coull, B., and Lange, C. (2017). A flexible  
25 and nearly optimal sequential testing approach to randomized testing: QUICK-STOP. *Genet.*  
26 *Epidemiol.* (accepted).
- 27 24. Rakovski, C.S., Xu, X., Lazarus, R., Blacker, D., and Laird, N.M. (2007). A new multimarker  
28 test for family-based association studies. *Genet. Epidemiol.* *31*, 9–17.
- 29 25. 1000 Genomes Project Consortium, Auton, A., Brooks, L.D., Durbin, R.M., Garrison, E.P.,  
30 Kang, H.M., Korbel, J.O., Marchini, J.L., McCarthy, S., McVean, G.A., et al. (2015). A global  
31 reference for human genetic variation. *Nature* *526*, 68–74.
- 32 26. Chen, R., Wei, Q., Zhan, X., Zhong, X., Sutcliffe, J.S., Cox, N.J., Cook, E.H., Li, C., Chen, W.,  
33 and Li, B. (2015). A haplotype-based framework for group-wise transmission/disequilibrium  
34 tests for rare variant association analysis. *Bioinforma. Oxf. Engl.* *31*, 1452–1459.
- 35 27. He, Z., O’Roak, B.J., Smith, J.D., Wang, G., Hooker, S., Santos-Cortez, R.L.P., Li, B., Kan, M.,  
36 Krumm, N., Nickerson, D.A., et al. (2014). Rare-variant extensions of the transmission

- 1 disequilibrium test: application to autism exome sequence data. *Am. J. Hum. Genet.* *94*, 33–  
2 46.
- 3 28. (2016). NHLBI TOPMed: The Genetic Epidemiology of Asthma in Costa Rica.
- 4 29. Fier, H.L., Prokopenko, D., Hecker, J., Cho, M.H., Silverman, E.K., Weiss, S.T., Tanzi, R.E.,  
5 and Lange, C. (2017). On the association analysis of genome-sequencing data: A spatial  
6 clustering approach for partitioning the entire genome into nonoverlapping windows. *Genet.*  
7 *Epidemiol.* *41*, 332–340.
- 8 30. He, Z., Xu, B., Buxbaum, J., and Ionita-Laza, I. (2019). A genome-wide scan statistic  
9 framework for whole-genome sequence data analysis. *Nat. Commun.* *10*, 1–11.
- 10 31. Benjamini, Y., and Hochberg, Y. (1995). Controlling the False Discovery Rate: A Practical  
11 and Powerful Approach to Multiple Testing. *J. R. Stat. Soc. Ser. B Methodol.* *57*, 289–300.
- 12 32. Ding, X., Lange, C., Xu, X., and Laird, N. (2009). New powerful approaches for family-  
13 based association tests with longitudinal measurements. *Ann. Hum. Genet.* *73*, 74–83.
- 14 33. Lange, C., Blacker, D., and Laird, N.M. (2004). Family-based association tests for survival  
15 and times-to-onset analysis. *Stat. Med.* *23*, 179–189.
- 16 34. Lange, C., Silverman, E.K., Xu, X., Weiss, S.T., and Laird, N.M. (2003). A multivariate  
17 family-based association test using generalized estimating equations: FBAT-GEE. *Biostatistics*  
18 *4*, 195–206.
- 19
- 20
- 21
- 22
- 23
- 24
- 25
- 26
- 27
- 28
- 29

## 1 Figures



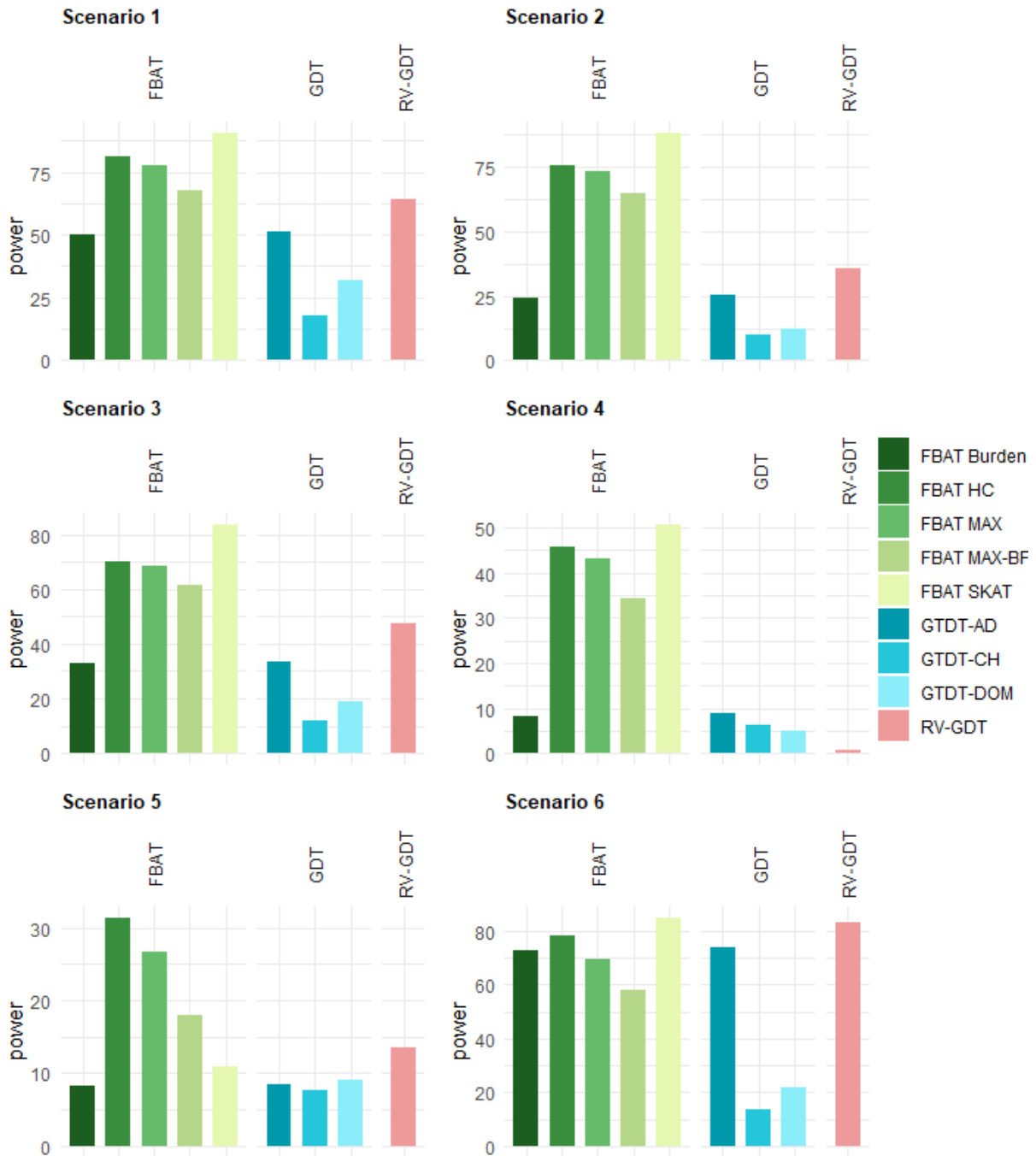
2

3 Figure 1. Power results in six different scenarios for genetic regions consisting of 30 variants at a significance  
 4 level of  $\alpha = 0.05$ . All results based on 1,000 replicates.

5

6



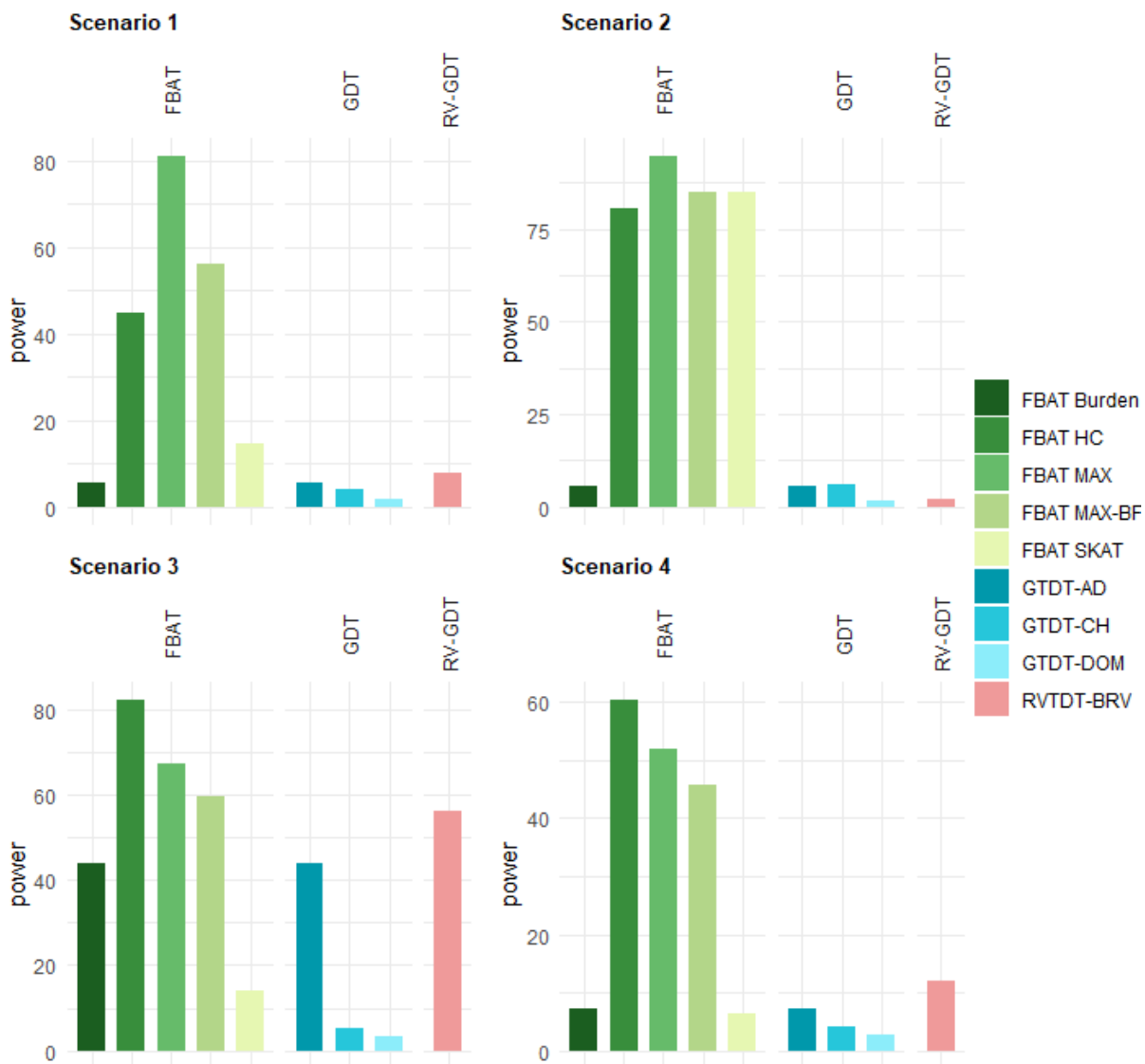


1

2 Figure 2. Power results in six different scenarios for genetic regions consisting of 50 variants at a significance  
 3 level of  $\alpha = 0.05$ . All results based on 1,000 replicates.

4

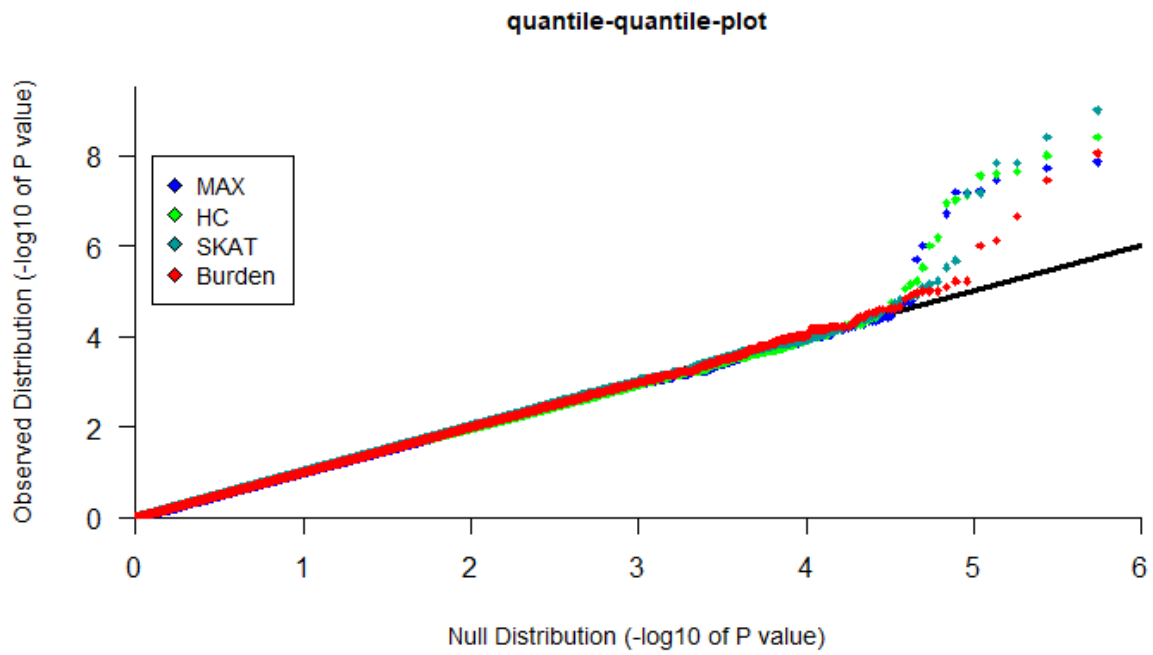
5



1

2 Figure 3. Power results for four different scenarios for genetic regions consisting of 1,000 variants at a  
 3 significance level of  $\alpha = 0.05$ . All results based on 1,000 replicates.

4



1

2 Figure 4. Real data analysis of 897 asthmatic offspring trios in a WGS study. Quantile-quantile plot for Burden,  
3 SKAT, MAX, and HC test statistics based on approximately 547,000 windows of 50 consecutive rare variants.

4

5

6

7

8

9

10

11

12

13

14

15

## 1 Tables

	FBAT					GTD			RV-GDT
	Burden	SKAT	MAX	HC	MAX-BF	GTD-AD	GTD-DOM	GTD-CH	RV-GDT
<b><math>p = 30</math></b>									
null	4.7%	5.4%	5.7%	5.0%	2.9%	5.2%	5.6%	5.2%	4.4%
adm1	5.7%	5.5%	4.4%	5.3%	3.0%	6.2%	5.3%	5.4%	5.7%
adm2	4.8%	4.7%	3.9%	4.1%	2.3%	5.1%	6.2%	5.0%	0.0%
adm3	4.2%	4.6%	3.2%	4.0%	2.0%	4.5%	5.1%	4.9%	99.6%
<b><math>p = 50</math></b>									
null	3.1%	3.6%	3.7%	4.0%	2.0%	3.4%	3.5%	3.8%	4.4%
adm1	4.3%	5.2%	4.0%	3.5%	1.2%	4.7%	4.8%	5.2%	4.0%
adm2	4.1%	4.4%	4.8%	4.7%	2.4%	4.0%	5.3%	4.9%	0.0%
adm3	4.5%	5.4%	6.4%	5.1%	2.7%	4.7%	4.5%	5.2%	35.4%

2 Table 1. Type I errors at a significance level of 5% for the FBAT, GTD and RV-GDT statistics. We considered four  
3 scenarios, separately for  $p = 30$  and  $p = 50$  variants. All results based on 1,000 replicates.  
4 null: no association between genetic variants and phenotype, no population admixture  
5 adm1: 1 parent generated from CEU haplotypes, 1 parent generated from GBR haplotypes. Parents unaffected,  
6 offspring affected.  
7 adm2: 1 parent generated from CEU haplotypes, 1 parent generated from GBR haplotypes. CEU parent  
8 affected, GBR parent unaffected, offspring affected.  
9 adm3: 1 parent generated from CEU haplotypes, 1 parent generated from GBR haplotypes. CEU parent  
10 unaffected, GBR parent affected, offspring affected.  
11

12

13

14

15

16

17

18

19

	FBAT					GTD			RV-GDT	
	Burden	SKAT	MAX	HC	MAX-BF	GTD-AD	GTD-DOM	GTD-CH	RV-GDT	
<b><i>p</i> = 30</b>	1	32.2%	87.9%	77.6%	79.2%	72.0%	32.7%	39.7%	17.2%	45.5%
	2	19.2%	85.4%	75.1%	76.9%	69.7%	19.7%	21.5%	10.7%	27.4%
	3	20.9%	83.7%	72.5%	73.1%	68.2%	21.9%	27.0%	11.0%	32.7%
	4	7.1%	48.0%	48.8%	48.7%	41.4%	7.7%	8.7%	5.3%	0.17%
	5	11.3%	11.5%	32.6%	39.5%	24.8%	11.6%	13.7%	7.5%	18.8%
	6	65.9%	82.3%	71.5%	77.8%	64.2%	67.5%	29.1%	13.4%	77.2%
<b><i>p</i> = 50</b>	1	50.2%	90.8%	77.8%	80.9%	67.9%	51.3%	31.8%	17.9%	64.1%
	2	24.0%	88.3%	73.5%	75.3%	64.9%	25.2%	12.1%	9.9%	35.6%
	3	32.9%	83.4%	68.6%	70.0%	61.3%	33.4%	18.9%	11.7%	47.6%
	4	8.2%	50.6%	43.3%	45.8%	34.3%	8.7%	4.9%	6.4%	0.7%
	5	8.3%	11.0%	26.8%	31.4%	17.9%	8.5%	9.0%	7.6%	13.5%
	6	72.6%	84.9%	69.3%	78.5%	58.0%	73.7%	21.9%	13.5%	83.3%

1 Table 2. Power estimates at a significance level of 5% for the FBAT, GTD and RV-GDT statistics. We considered  
2 six scenarios, separately for  $p = 30$  and  $p = 50$  variants. All results based on 1,000 replicates.  
3 scenario 1: three causal variants, effect sizes  $0.4|\log_{10}(MAF)|$ , same direction  
4 scenario 2: three causal variants, effect sizes  $0.4|\log_{10}(MAF)|$ , different direction  
5 scenario 3: two causal variants, effect sizes  $0.4|\log_{10}(MAF)|$ , same direction  
6 scenario 4: two causal variants, effect sizes  $0.4|\log_{10}(MAF)|$ , different direction  
7 scenario 5: four very rare causal variants, effect size 1.0, same direction  
8 scenario 6: three causal variants, effect sizes  $0.4|\log_{10}(MAF)|$ , same direction, in strong LD with other  
9 variants.

10  
11  
12  
13  
14  
15  
16  
17  
18  
19

1  
2  
3

		FBAT					GTDT			RV-TDT
		Burden	SKAT	MAX	HC	MAX-BF	GTDT-AD	GTDT-DOM	GTDT-CH	RV-TDT BRV
$\alpha = 0.05$	null	5.2%	4.5%	3.9%	5.0%	0.7%	5.2%	2.5%	5.9%	5.5%
	1	5.6%	14.5%	81.2%	45.0%	56.0%	5.6%	1.7%	4.1%	8.0%
	2	5.6%	85.0%	94.9%	80.8%	85.2%	5.6%	1.5%	6%	2.1%
	3	43.8%	14.2%	67.4%	82.3%	59.7%	43.8%	3.3%	5.2%	56.2%
	4	7.1%	6.3%	51.9%	60.4%	45.8%	7.1%	2.8%	4.1%	12.0%
$\alpha = 0.01$	null	1.3%	1.5%	0.7%	1.2%	0.0%	1.3%	0.4%	1.1%	1.3%
	1	1.3%	3.4%	61.2%	39.9%	33.0%	1.3%	0.2%	0.8%	2.3%
	2	0.9%	54.6%	87.6%	77.7%	71.0%	0.9%	0.2%	1.3%	0.5%
	3	23.0%	2.7%	39.2%	48.3%	32.0%	23.0%	0.5%	1.3%	32.2%
	4	1.8%	0.9%	27.3%	31.9%	20.8%	1.8%	0.5%	0.8%	3.1%

4 Table 3. Type I error and power estimates at significance levels of 1% and 5%, all results based on 1,000  
5 replicates. MAX-BF refers to the test that at least one single variant test statistic reached the Bonferroni-  
6 corrected significance level based on  $p = 1,000$  tests.  
7 scenario 1: two causal variants,  $MAF \sim 0.2\%$ , almost no LD with other variants, effect size 1.8  
8 scenario 2: two causal variants,  $MAF \sim 2\%$ , almost no LD with other variants, effect size 0.7  
9 scenario 3: 16 causal variants,  $MAF \sim 0.1\%$ , almost no LD with other variants, effect size 0.7  
10 scenario 4: 16 causal variants,  $MAF \sim 0.1\%$ , almost no LD with other variants, effect size 0.7, effect direction  
11 alternates  
12