# Metadata of the chapter that will be visualized online

| | | |
|---|---|---|
| Chapter Title | Gene-Gene and Gene-Environment Interactions | |
| Copyright Year | 2018 | |
| Copyright Holder | Springer Science+Business Media, LLC, part of Springer Nature | |
| Corresponding Author | Family Name | **DeWan** |
| | Particle | |
| | Given Name | **Andrew T.** |
| | Suffix | |
| | Division | Department of Chronic Disease Epidemiology |
| | Organization | Yale School of Public Health |
| | Address | New Haven, CT, USA |
| | Email | andrew.dewan@yale.edu |
| Abstract | Identifying gene–gene and gene–environment interactions may help us to better describe the genetic architecture for complex traits. While advances have been made in identifying genetic variants associated with complex traits through more dense panels of genetic variants and larger sample sizes, genome-wide interaction analyses are still limited in power to detect interactions with small effect sizes, rare frequencies, and higher order interactions. This chapter outlines methods for detecting both gene-gene and gene-environment interactions both through explicit tests for interactions (i.e., ones in which the interaction is tested directly) and non-explicit tests (i.e., ones in which an interaction is allowed for in the test, but does not test for the interaction directly) as well as approaches for increasing power by reducing the search space. Issues relating to multiple test correction, replication, and the reporting of interaction result in publications. | |
| Keywords (separated by '-') | Interaction - Epistasis - Environment - GWAS - Power - Replication | |

# Chapter 7

## Gene-Gene and Gene-Environment Interactions

### Andrew T. DeWan

#### Abstract

Identifying gene–gene and gene–environment interactions may help us to better describe the genetic architecture for complex traits. While advances have been made in identifying genetic variants associated with complex traits through more dense panels of genetic variants and larger sample sizes, genome-wide interaction analyses are still limited in power to detect interactions with small effect sizes, rare frequencies, and higher order interactions. This chapter outlines methods for detecting both gene-gene and gene-environment interactions both through explicit tests for interactions (i.e., ones in which the interaction is tested directly) and non-explicit tests (i.e., ones in which an interaction is allowed for in the test, but does not test for the interaction directly) as well as approaches for increasing power by reducing the search space. Issues relating to multiple test correction, replication, and the reporting of interaction result in publications.

**Key words** Interaction, Epistasis, Environment, GWAS, Power, Replication

## 1 Introduction

An interaction results when the effect of one factor is only evident in the presence of another. These factors could be genetic markers and/or environmental exposures. Much has been written on the topic of gene-gene (also known as epistasis) and gene-environment interactions with several comprehensive reviews of study designs and methods for analyzing both gene-gene [1–3] and gene-environment interactions [4, 5]. But why are we interested in studying interactions? Both gene-gene and gene-environment and potentially gene-gene-environment interactions allow us to better describe the underlying genetic architecture of a particular trait and as such we can begin to fill in the missing heritability [6] for a particular phenotype.

Biological interactions were originally defined as the situation when the phenotypic effect of one gene was only evident in the presence of a second gene [7]. In contrast, a statistical interaction is defined as a departure from a linear model combining two or more genetic factors (or a genetic factor and environmental factor) [8].

Andrew T. DeWan

Another way to think about this is that biological interactions are observed at the individual level and statistical interactions at the population level, but this does not imply that observing evidence of one will lead to observing evidence of another [9]. The traditional method to test for statistical interactions is to use a regression-based model that includes main effect variables representing each genetic factor and an interaction term (the product of the main effect variables), and then testing for the significance of the interaction term after adjusting for the main effects. However, several other methods exist to test for statistical interactions and these will be discussed as well. 34

We now have the ability to conduct very dense genome-wide association studies with up to five million genetic markers genotyped at one time and millions more imputed using large reference panels from population-based sequencing projects. The combinatorial problem can be immense if you attempt to look at all possible marker combinations with little power to detect significant interactions after accounting for all interactions tested. Given the exponential improvements in computational power and the relative ease of parallel computing, the computational hurdles of examining all pairwise gene-gene interactions are not insurmountable, and exhaustive searches of higher order interactions will follow. However, the immense power constraints on currently available sample sizes, while improving with studies routinely examining 10,000 to more than 100,000 subjects, are still underpowered to detect interactions with modest effect sizes ($OR < 1.2$) and low frequency variants (minor allele frequency [MAF] $< 0.01$). Though marker pruning using linkage disequilibrium can reduce this problem somewhat, there are several other data reduction approaches that will be discussed based on prior biological knowledge or statistical evidence for association. These can mitigate the multiple testing burden but depend heavily on the quality of this prior information.

Similar to gene-gene interactions, gene-environment interactions arise when the effect of a genetic factor on a phenotype is dependent on the presence or absence of an environmental factor. Statistically, this can be tested in a similar fashion as for gene-gene interactions traditionally done using a regression-based model with main effects for the genetic factor and environmental factor and an interaction term and testing for the significance of the interaction term. This environmental factor may be one traditionally thought of as an environmental exposure such as smoking, indoor $NO_2$ levels, or sun exposure, but these could also be other potentially genetically influenced phenotypes such as obesity, blood glucose levels, or birthweight that may be influencing the effect of the genetic factor depending on the value of these secondary phenotypes.

As outlined below, direct assessment of interactions (termed as "explicit" test) is often less powerful than joint tests that include

either an interaction term or allow for interaction but do not test for the significance of the interaction directly (termed as "non-explicit" test). This classification of tests for interactions will be used throughout this chapter as different approaches and tests are discussed.

## 2 Materials

### 2.1 Data

#### 2.1.1 Genotype Data

The core of both gene-gene and gene-environment interaction analyses is obtaining high-quality genotype data. One may use data generated specifically for a particular project, but there are also many outstanding datasets available for analysis from a number of online repositories.

If samples are being genotyped for the specific project there are many options available. There are whole genome microarray panels of markers based on relatively even coverage across the genome, panels of markers that maximize the coverage for specific race/ethnic groups, and panels that allow custom markers to be added-on to existing panels to increase coverage of specific genes of interest or previously associated markers, for example. There are also panels of markers that target specific regions of the genome such as the exonic regions, cancer-associated genes, and metabolic genes. The selection of the right microarray will depend on your specific study hypotheses, type of study, and budget. Lower throughput genotyping can be done for single variants to hundreds of markers simultaneously depending again on the hypotheses and goals of the project.

Genotyping known panels of markers is not the only choice for interaction studies. One could choose to utilize sequencing-based approaches to genotype unknown variants and/or low frequency and rare variants. Whole genome or whole exome sequencing could be utilized for hypothesis-free analyses or targeted sequencing if specific genes or pathways are hypothesized to be involved in the interaction(s).

If secondary data analysis is an option, there are a plethora of datasets with genome-wide data available as well as extensive phenotype data. Two such repositories are the database of genotype and phenotype (dbGaP) maintained by NCBI primarily of studies conducted in the United States (https://www.ncbi.nlm.nih.gov/gap) and the European Genome-phenome Archive (EGA) maintained by EBI primarily of studies conducted in Europe (https://www.ebi.ac.uk/ega/home). These databases contain hundreds of datasets accessible through an application to the respective data access committees.

#### 2.1.2 Environmental Data

The environmental factors that could be considered in a study of gene-environment interactions are extremely broad. These include

Andrew T. DeWan

chemical (e.g., polychlorinated biphenyl (PCB)), physical (e.g., airborne particulate matter), biological (e.g., viral infections), and lifestyle (e.g., physical activity). The measurement of each of these individual environmental factors is going to vary widely depending on the environmental factor of interest. This could range from measure $NO_2$ levels in the air via chemiluminescence, measuring radiation exposure via a dosimeter, conducting a daily food diary to estimate saturated fat intake, or reviewing charts to collect data on BMI history. The discussion in this chapter will focus on environmental exposures at one time point, but there is some evidence that longitudinal environmental data may increase the power to detect gene-environment interactions for common diseases [10].

*2.1.3 Biological/Functional Data*

For analyses that are pursuing a hypothesis-driven approach and/or filtering based on biological information, one may want to utilize prior biological or functional data. There are a variety of databases and programs that can be accessed to provide this type of information. One could simply curate information from publications in the scientific literature through systemic reviews of publications in databases such as PubMed. Data on functionality of variants can be obtained more systematically from databases such as ENCODE (Encyclopedia of DNA Elements) [11] that contains a comprehensive list of functional elements at both the RNA and protein levels and is available for viewing or downloading from the UCSC Genome Browser (www.genome.ucsc.edu/ENCODE). Direct annotation of variants could be conducted using a program such as Annovar [12] to annotate variants as to their respective genes, coding vs. noncoding, and predicted functional consequence. An alternative annotation program more directly related to filtering variants for interaction analyses is Biofilter which allows for the annotation of variants based on previous association studies and biological knowledge, filtering variants based on specific biological hypotheses, and building sets of testable variant interactions based on implication indices compiled from available data [13].

*2.1.4 Previous Statistical Data*

For analyses filtering variants based on prior statistical knowledge, data from one's own GWAS or single-variant association study could be used, results mined from previous publications or, alternatively, association results obtained from databases such as the GWAS Catalog (http://www.ebi.ac.uk/gwas/). While results from previous publications or the GWAS Catalog are a convenient and useful resource, they have the disadvantage of being biased toward reporting only genome-wide significant results and other nominally significant results will likely not be available and should be kept in mind when planning the analysis approach.

**2.2 Software**

Below are a listing of programs that can be used to conduct gene-gene and/or gene-environment interaction analyses with other

programs mentioned and described throughout the chapter. This is 175
not an exhaustive list and is only an example of programs that are 176
commonly used to assess interactions in genetic epidemiological 177
studies. Additional programs not directly related to the interaction 178
analysis such as for computing eigenvalues from principal compo- 179
nents analyses (e.g., EIGENSTRAT [14]) and imputation of var- 180
iants (e.g., IMPUTE2 [15]) are not listed. 181

182

*2.2.1 PLINK*

This is a suite of tools designed to conduct genome association 183
analyses, including both gene-gene and gene-environment interac- 184
tions [16]. The primary interaction analyses are based on logistic 185
and linear regression. They can accommodate both gene-gene and 186
gene-environment interactions on the genome-wide scale or on a 187
smaller number of variants by creating subsets of genetic variants to 188
test again each other or the environmental factor. The program has 189
the flexibility to conduct both explicit tests for interaction by 190
testing for the significance of the interaction term directly in the 191
regression model or a non-explicit joint test by testing the main and 192
interactions effects. There is also a faster option for conducting 193
genome-wide gene-gene interactions (*fast-epistasis*) based on the 194
Z-score for the differences in OR for SNP-SNP combinations 195
between cases and controls or for cases alone (case-only test). 196

197

*2.2.2 CASSI*

This is a software package that is specifically designed to conduct 198
genome-wide gene-gene interaction analyses in a computationally 199
efficient manner ([17]; https://www.staff.ncl.ac.uk/richard. 200
howey/cassi/index.html). This package corrects a minor error in 201
the Wu et al. statistic [18] in the calculation of the variance for 202
estimated rather than observed haplotypes and in the *fast-epistasis* 203
variance originally implemented in PLINK. 204

205

*2.2.3 BOOST*

The Boolean Operation-based Screening and Testing (BOOST) 206
program was designed to efficiently screen and then explicitly test 207
for genome-wide gene-gene interactions [19]. The screening phase 208
involves a non-iterative procedure to approximate the likelihood 209
ratio and then all variant pairs that survive this screening are sub- 210
jected to a classical likelihood ratio test in the testing phase. 211

212

*2.2.4 MDR*

The Multifactor Dimensionality Reduction (MDR) software pack- 213
age [20] is designed to conduct data mining on discrete variables 214
and can be used to detect both gene-gene and gene-environment 215
interactions and dichotomous outcomes [21]. The traditional 216
MDR approach is a non-explicit test for interaction as it is a non- 217
parametric test that combines factors that may be interacting in 218
order to best discriminate the subjects among the dichotomous 219
outcome. An extension of the MDR has been developed that 220
incorporates a permutation-based approach that can explicitly test 221
for interactions [22]. A recent extension to the MDR has 222

implemented a *t*-test approach that allows for quantitative out- 223
comes [23]. The MDR method, however, is designed primarily 224
for smaller sets of markers, but parallel computing could be utilized 225
to conduct a genome-wide analysis. 226

227

## 3 Methods
228

### 3.1 Quality Control (QC)

No analysis can be successful without high-quality data. The spe- 229
cific steps of the genotype QA/QC will depend on the type of assay 230
used to generate the genotypes. These range from single-variant 231
assays based on PCR, whole genome microarray genotyping, and 232
whole-exome and whole-genome sequencing. A brief outline of the 233
QC steps for each is outlined below. 234

235

#### 3.1.1 Single-Variant QC

The primary steps are to assess the overall performance of the 236
individual genotyping assays through examination of the variant 237
call rate (variant call rate = total number of genotype calls/total 238
number of individuals genotyped) and Hardy-Weinberg Equilib- 239
rium (HWE). While the thresholds chosen to eliminate variants can 240
be arbitrary, typically one would look for variant call rates >98% 241
(which should be examined in cases and controls separately if con- 242
ducting a case-control analysis to ensure no bias due to differences 243
in call rates between cases and controls) and HWE *p*-values $>10^{-4}$ 244
which if conducting a case-control study are assessed only in con- 245
trols. Without the benefit of genome-wide genotype data it is 246
impossible to assess the data for population stratification, but 247
adjustments can be made in the analysis (if using regression-based 248
methods) for relevant covariates that may capture potential stratifi- 249
cation such as self-reported race/ethnicity. 250

251

#### 3.1.2 Microarray QC

As with single-variant QC, one will examine both the individual 252
variant call rates and HWE to ensure that each variant probe is 253
generating high-quality genotype data with similar thresholds 254
applied as mentioned above. However, additional steps can and 255
should be taken into account. The individual subject call rates 256
should be examined first to determine if there were general pro- 257
blems with the individual array and/or DNA. These call rate 258
thresholds may range from 93% to 98% and are often suggested 259
by the array manufacturer-based past performance of the array 260
(subject call rate = total number of genotype calls for an individual 261
subject/total number of variants attempted to be genotyped). 262
Poorly performing subjects should be removed prior to any down- 263
stream QC steps. Population stratification should be assessed using 264
a genome-wide procedure such as principal components analysis to 265
determine if there are slight variations in the genotype frequencies 266
between subpopulations within your dataset. This procedure can 267
detect any systematic differences that may be due to differences in 268

allele frequencies arising ancestry differences, but also due to exper- 269
imental/processing differences (e.g., plate effects). If significant 270
principal components (PCs) are detected it is suggested that these 271
PCs be adjusted for in the analysis. If the analysis assumes unrelated 272
subjects it is suggested that the dataset be examined for cryptic 273
relatedness using a procedure such as estimate pairwise identity-by- 274
decent (IBD). This pairwise measure that is often used is $\widehat{p}_i$. Again, 275
the threshold for identifying cryptically related subjects is arbitrary 276
one often chooses a threshold ranging from $0.125$ to $0.2$ and then 277
eliminates one of the two subjects in this cryptically related pair. 278
This can be done randomly, or one may want to eliminate the 279
subjects based on the subject call rate (eliminating the subject 280
with a fewer genotype calls) or if it is a case-control study and a 281
cryptically related pair is comprised of a case subject and a control 282
subject it may be beneficial to eliminate the control if cases are in 283
short supply. 284

To assess whether or not the QC steps that have been taken are 285
successful prior to conducting an interaction analysis, it would be 286
beneficial to conduct a genome-wide single-variant analysis and 287
examine QQ plots and/or estimate $\lambda$ from the data after 288
adjusting for PCs and other covariates. QQ plots can be generated 289
using an R script such as qqman.r (https://CRAN.R-project.org/ 290
package=qqman) and $\lambda$ estimated using PLINK. Deviations from 291
the expected line on the QQ plot are not expected except at the tail 292
(i.e., the true positives) and with deviations along much of the 293
expected line being an indication of residual population stratifica- 294
tion. $\lambda$ estimates greater than $1.05$ are routinely seen as indicators of 295
population stratification and additional PCs should be adjusted for 296
until the lambda value falls below this threshold. This assessment is 297
typically done in single-variant analyses prior to any interaction 298
analyses. 299

300

*3.1.3 Sequencing QC*    Specific workflows for alignment, variant calling, and variant 301
QA/QC and filtering are described in detail elsewhere (*see* Ref. 302
24 for a detailed step-by-step pipeline covering the major sequenc- 303
ing analysis tools). Briefly, a standard analysis pipeline would start 304
by aligning the FASTQ raw sequence reads to a reference genome 305
using the Burrows-Wheeler Aligner (BWA, [25]).Then converted 306
to BAM format, sorted, indexed, PCR duplicates marked and then 307
merged into one BAM file using SAMtools [26]. Finally, the align- 308 AU3
ments in the BAM file can be locally realigned around insertion/ 309
deletions, recalibrated and variants called using HaplotypeCaller in 310
the Genome Analysis Toolkit (GATK, [27, 28]). 311

Variant QC can utilize a variety of different metrics, but an 312
example of one approach is how we conducted our QC in our 313
whole-exome sequence analysis of a family segregating asthma [29]. 314
Variants were flagged (and not considered further in our analysis) if 315

Andrew T. DeWan

they met any of the following criteria: three or more variants detected within 10 bp; four or more alignments map to different locations equally well; coverage of less than five reads; quality score <50; low quality for a particular sequence depth (variant confidence/unfiltered depth <1.5); and strand bias (Phred-scaled $p$-values using Fisher's Exact Test >200). There are many variations that can be employed in your QC pipeline, but the most important aspect is achieving the highest quality set of variants to retain in your analysis.

*3.1.4 Linkage Disequilibrium Pruning*

Reducing the number of variants considered in the analysis, regardless of genotyping method, can be accomplished by linkage disequilibrium (LD) pruning. LD is an indication of the correlation, or non-independence, between variants and can be measured using $r^2$ or D'. A typical $r^2$ threshold used to prune variants is 0.8, but a lower threshold can be used to eliminate more variants at the risk of excluding some that are independently informative. This procedure is particularly important when conducting a case-only gene-gene interaction analysis (discussed in Subheading 3.2.12) as the variants must not be correlated with each other.

**3.2 Gene-Gene Interaction Analysis**

The basic analysis is straightforward, the assessment of a deviation from an additive or multiplicative model containing two or more variants. This is traditionally assessed through regression-based modeling but many different methods are available and several different methods will be discussed below to exemplify this approach. One could divide these approaches into explicit vs. non-explicit tests for interactions, with explicit tests determining if the null hypothesis of the sum on the additive/linear scale or the product on the multiplicative scale of the joint effects of the two variants is contributing to the outcome with the alternative hypothesis being that the joint effect of the two variants is greater or less than the expected. Non-explicit tests are able to determine if the grouping of variants is associated with the outcome, but not necessarily that the variants together deviate from an additive or multiplicative model. Different inheritance models can be imposed on the genetic variants in the interaction models (*see* **Note 1**) and the issue of outcome scale is also an important consideration when interpreting interactions (*see* **Note 2**).

Regression-based approaches are attractive not only for their ability to explicitly detect interactions in the data, but because of the ability to adjust for multiple covariates. In a genome-wide setting this is important so that significant principal components can be adjusted to remove population stratification, but other potential confounders can be adjusted as well, including age, sex, etc. depending on the outcome of interest.

*3.2.1 Logistic Regression*     For dichotomous outcomes, a logistic regression model is one of   362
the most straightforward approaches for testing an interaction,   363
either explicitly or non-explicitly. In the following logistic regres-   364
sion model we are modeling the probability of our dichotomous   365
outcome ($p$) using an intercept ($\beta_0$), two genetic variants (SNP1   366
and SNP2, for example) and their respective main effects ($\beta_1$ and   367
$\beta_2$) and their interaction effect ($\beta_3$):   368

$$\log\left(\frac{p}{1-p}\right) = \beta_0 + \beta_1 \text{SNP1} + \beta_2 \text{SNP2} + \beta_3 \text{SNP1} \times \text{SNP2}$$

In this model, we can test for the significance of the interaction   369
effect of the two variants by testing the null hypothesis that $\beta_3 = 0$.   370
This is an explicit test for interaction in that we are testing for an   371
interaction after adjusting for the main effects (i.e., independent   372
effects) of the two genetic variants. This model will produce an   373
effect estimate ($\beta_3$) which can easily be converted to an OR by   374
$\exp\beta_3$. It should be noted that the main effects should not be   375
interpreted as the main effect of the variants (or the variant and   376
the environmental factor) since they are adjusted for the interaction   377
term in the model.   378

Alternatively, we can test for the significance of the joint effects   379
of the two variants by using a 2 degree of freedom likelihood ratio   380
test of the full model against a model in which there is neither an   381
interaction term nor main effect term for one of the SNPs:   382

$$\log\left(\frac{p}{1-p}\right) = \beta_0 + \beta_1 \text{SNP1}$$

While this model contains an interaction term this is a   383
non-explicit test for interaction as we are testing the significance   384
of the main effect of one of the SNPs and interaction term com-   385
bined. This may be seen as a less desirable test but this can be more   386
powerful for detecting significant signals [30]. However, this joint   387
test can be significant if either of the main effects is strong enough   388
or if both main effects are strong enough without a significant   389
interaction term ($\beta_3$). As will be discussed later, one may want to   390
screen combinations of variants for joint effects and then follow up   391
for explicit tests of an interaction only on the limited set of variants   392
demonstrating significant joint effects.   393

When reporting the results of regression-based analyses, it is   394
important to report not only the significance of the interaction ($p$-   395
value), but also the parameter estimate of the interaction as well as   396
those of the main effects (*see* **Note 3**).   397
   398

*3.2.2 Linear Regression*     For continuous outcomes, linear regression can be used to test for   399
interactions in much the same way as logistic regression for dichot-   400
omous outcomes. $\beta$s are not converted to ORs, but are rather   401

interpreted directly with the sign of $\beta_3$, an indication of the direc- 402
tion of the interaction, and the $p$-value for the likelihood ratio test 403
of the full model versus one without the interaction term an indi- 404
cation of the significance of the interaction term. A two degree of 405
freedom joint test can also be constructed. 406

407

*3.2.3 Nonparametric Analyses*

The most popular nonparametric method for interaction analysis 408
are the multidimensional reduction methods. These are based on 409
the idea that across a contingency table of two or more genetic 410
variant genotype combinations or genotype and discrete environ- 411
mental factor combinations, each square in the contingency table 412
can be divided into increasing risk and decreasing risk in the context 413
of the phenotype. These collapsed contingency tables can then be 414
tested for association with the phenotype. This is typically done by 415
splitting the data into a training set (to build the collapsed contin- 416
gency table) and a testing set to determine the classification error. 417
The data are split multiple times and the model is assessed by the 418
classification error and prediction error estimated from these mul- 419
tiple cross-validations [21]. Due to the fact that this is a nonpara- 420
metric method, a $p$-value is obtained for the set of factors included 421
in the model, but an effect estimate is not. The significance of the 422
model is interpreted as the significance of the group of factors 423
(genetic variants or genetic variant(s) and environmental factors) 424
but cannot be interpreted as direct evidence of an interaction. 425

The model-based MDR (MB-MDR) methodology is an exten- 426
sion of the classic MDR framework that allows for the direct testing 427
of interactions through the use of the Wald statistic on high and 428
low-risk genotype categories [22, 31]. The significance of the 429
interaction is then tested using permutation testing of the maxi- 430
mum Wald statistic and can be considered an explicit test of an 431
interaction. 432

433

*3.2.4 Multiple Test Correction*

One of the biggest challenges in conducting genome-wide analysis 434
in general, and more specifically interaction analysis, particularly 435
gene-gene interactions, is the multiple testing problem. Given that 436
hundreds of thousands or millions of variants are considered in 437
most genome-wide studies, the chance of detecting a false positive 438
is immense. This is even more apparent when multiple environmen- 439
tal factors are considered and most problematic when all pairwise 440
variant interactions are considered in an exhaustive gene-gene 441
interaction search. Traditionally, genome-wide studies have 442
adopted a Bonferroni correction to adjust for the number of var- 443
iants analyzed, but this can be overly conservative, especially for 444
dense maps of markers in high LD. However, other procedures 445
such as the False Discovery Rate (FDR; [32]), LD-based variant 446
counting to account only for the number of independent tests [33] 447
and permutation testing [34, 35] have been adopted. Permutation 448

testing, while attractive because it accounts for the number of 449
independent tests by maintaining the LD structure in the permuted 450
data, is computationally intensive, even more so when all pairwise 451
variant combinations need to be considered in thousands of per- 452
muted datasets. Therefore, the FDR or LD-based methods are 453
more attractive alternatives. 454
455

*3.2.5   Genome-Wide
(Exhaustive) Approaches*

A typical genome-wide association study will genotype between 456
100,000 and 1,000,000 singe-nucleotide polymorphisms. While 457
exponential advances in computational power mean that the com- 458
puting power to run the analysis on the pairwise tests required for a 459
panel of 500,000 SNPs ($2.5 \times 10^{11}$ tests) is not insurmountable 460
and with parallel computing relatively quick, the major hurdle is 461
power. While large multi-center consortium studies mean that the 462
number of subjects for many subjects has greatly increased, with 463
studies reaching 50,000 or more subjects depending on the fre- 464
quency of the disease or phenotype, these may still not be powerful 465
enough to detect significant interactions on the genome-wide scale. 466
In a study of genome-wide data obtained from the Resource for 467
Genetic Epidemiology Research on Adult Health and Aging 468
(GERA), using 45,171 subjects for ten phenotypes and conducting 469
an exhaustive search for interactions, we failed to identify any 470
genome-wide significant interactions, suggesting that we were 471
underpowered to detect interactions with apparently weak effect 472
sizes [36]. 473

As demonstrated in Table 1, for a genome-wide gene-gene 474
interaction study of 200,000 markers using an exhaustive approach, 475
requires than 30,000 cases and 30,000 controls to achieve 476
genome-wide significance for two loci each with an MAF of 0.2 and 477
an interaction OR of 1.2. All power calculations presented in this 478
chapter were performed in Quanto [37]. The power calculations 479
assumed a population risk of the disease of 0.1 and a log-additive 480
mode of inheritance and a main effect of each variant of 1.2. For a 481
case-control study with 5000 cases and 5000 controls, a reasonably 482
sized case-control GWAS, and 200,000 markers genotyped, the 483
minimal detectable OR is just over 1.6 (Table 2). For an interaction 484
OR of 1.5, the number of markers considered would need to be 485
reduced to around 2000 to detect a significant interaction among 486
all pairs of interactions (4,000,000) and at an interaction OR of 487
1.2, only one interaction can be considered as only a nominally 488
associated interaction is detectable (Table 3). 489

This lack of power can be daunting, but as detailed in the 490
following sections there are several ways to increase power by 491
reducing the number of tests through filtering or using biologically 492
informed combinations of variants. 493
494

Andrew T. DeWan

t.1 **Table 1**
**Power of interaction analysis**

| | Number of cases and controls | | | | | |
|---|---|---|---|---|---|---|
| Interaction type | 5000 | 10,000 | 20,000 | 30,000 | 40,000 | 50,000 |
| G-G | 0.000 | 0.003 | 0.164 | 0.654 | 0.940 | 0.995 |
| G-E | 0.024 | 0.032 | 0.317 | 0.717 | 0.928 | 0.988 |

t.1 **Table 2**
**Power of case-control and case-only gene-gene interaction analyses**

| intOR | Case-control | Case-only |
|---|---|---|
| 1.1 | 0.000 | 0.000 |
| 1.2 | 0.000 | 0.013 |
| 1.3 | 0.004 | 0.547 |
| 1.4 | 0.076 | 0.992 |
| 1.5 | 0.388 | 1.000 |
| 1.6 | 0.782 | 1.000 |
| 1.7 | 0.961 | 1.000 |
| 1.8 | 0.996 | 1.000 |
| 1.9 | 1.000 | 1.000 |
| 2.0 | 1.000 | 1.000 |

*3.2.6 Data Reduction Approaches*

In order to overcome the multiple testing burden of a genome-wide screen for all gene-gene interactions, it may be advantageous to focus on a subset of SNPs that may have a one or more properties that may make the interaction analysis more likely to detect a significant interaction. If we reduce our set of 200,000 SNPs down to 2000 we can now detect an interaction OR of 1.5 using a set of 5000 cases and 5000 controls (Table 3). Below I outline a series of methods to reduce the set of SNPs considered.

Filtering by Allele Frequency

The simplest method to increase power to detect interactions is to impose a filter by minor allele frequency (MAF). Depending on the sample size and subsequent power estimates, it may be advantageous to filter out all SNPs with MAFs less than the power to detect a reasonable interaction effect size (e.g., all SNPs with an MAF <0.2).

**Table 3**                                                                 t.1
**Power of gene-gene interaction analysis by interaction OR and number of markers**

| # Markers | intOR = 1.2 | intOR = 1.4 | intOR = 1.5 | |
|---|---|---|---|---|
| 200,000 | 0.000 | 0.076 | 0.388 | t.3 |
| 100,000 | 0.000 | 0.107 | 0.465 | t.4 |
| 20,000 | 0.000 | 0.221 | 0.649 | t.5 |
| 10,000 | 0.001 | 0.289 | 0.724 | t.6 |
| 2000 | 0.004 | 0.488 | 0.869 | t.7 |
| 1000 | 0.008 | 0.584 | 0.914 | t.8 |
| 200 | 0.037 | 0.793 | 0.975 | t.9 |
| 100 | 0.066 | 0.864 | 0.988 | t.10 |
| 20 | 0.219 | 0.966 | 0.999 | t.11 |
| 10 | 0.337 | 0.986 | 1.000 | t.12 |
| 5 | 0.488 | 0.995 | 1.000 | t.13 |
| 2 | 0.713 | 0.999 | 1.000 | t.14 |
| 1 | 0.865 | 1.000 | 1.000 | t.15 |

The column header row is labeled t.2.

**Filtering by Marginal Effects**

While a significant marginal and/or main effect is not required to 511 detect a significant interaction, there are few reports of significant 512 interactions without also having detectable main effects. Therefore, 513 it may be advantageous to filter based on marginal effects (i.e., 514 single-variant effect) and only include those variants that have, for 515 example, a single-variant association $p$-value of $<0.05$ (nominal 516 significance threshold). This should reduce the number of SNPs 517 to approximately 5% of the starting number (for our 500,000 SNPs 518 this would result in $6.25 \times 10^8$ tests). This approach may be too 519 stringent, as it requires both variants be nominally significant. An 520 alternative approach would be to select a set of SNPs reaching a 521 predefined significance threshold (e.g., $p < 10^{-4}$) and testing this 522 set of SNPs against all other SNPs. In our example with 500,000 523 SNPs we would select ~50 SNPs with $p < 10^{-4}$ and test for an 524 interaction with the remaining 499,450 for a total of $2.5 \times 10^{-7}$ 525 SNPs. This balances the requirement that there be some marginal 526 effect of one of the SNPs with being able to detect interactions with 527 SNPs showing no marginal effects but having a significance effect 528 on disease only in the presence of a second SNP. 529

530

**Candidate Gene Approaches**

Filtering based on prior evidence of a gene's involvement in a 531 particular disease is another approach to reducing the search space 532 for interactions. Candidate genes could be selected by systematically 533

reviewing the literature or for some diseases databases of candidate genes exist based on association studies, linkage analyses, and/or expression studies. For example, there are databases for preterm birth [38], preeclampsia [39], and non-syndromic hearing loss [40] for which one can obtain lists of candidate genes. By annotating variants to their respective genes using programs such as Annovar [12] only those SNPs annotated to the set of candidate genes can be selected for inclusion in the interaction analysis. As for the marginal effects, it may be advantageous to consider those SNPs within candidate genes against all other SNPs genotyped. In this way, it is possible to detect interactions with SNPs in novel genes not previously identified to be associated with the disease of interest.

**Filtering by Function**

Filtering by the effects of specific variants is another approach to reducing the search space. Again, variants can be annotated, but now the selection may be made based on being a coding variant or a splicing variant that is more likely to be functional. Other biological information could also be utilized such as examining interactions between SNPs in genes known to be involved in protein-protein interactions with the rationale being that SNPs in these biologically interacting genes are more likely to also show evidence of a statistical interaction. Similarly, one could examine interactions between variants within transcription factors and those within their binding sites.

### 3.2.7 Multistage Approaches

In order to maintain a genome-wide approach but overcome the hurdle of the immense multiple testing problem, it may be beneficial to employ a multistage approach. This is possible in situations in which a study has a large sample size, but still not adequately powered to detect significant interactions genome-wide. In a study of asthma, we first screened all pairwise interactions ($9.1 \times 10^{10}$) in a small subset of the data and then carried through all interaction with a suggestive significance ($p < 10^{-5}$) to a follow-up stage of independent subjects and then attempted to replicate the top SNPs in a third set of independent subjects [41]. While this approach did not identify any genome-wide significant interactions, it did identify a candidate interaction between SNPs in two regions of the genome. The major advantage of this type of approach is that it allows for an unbiased examination of the SNPs without relying on previously reported biological and/or association data.

### 3.2.8 Gene-Based Interaction Tests

Another approach to reduce the multiple testing burden and combine information across multiple variants is to conduct a gene-based test of interaction. By considering each gene as a unit in the interaction rather than each individual variant, the total number of interactions considered is significantly reduced, thus increasing power. One such approach combines interaction $p$-values across all

combinations of genetic variants in two genes into a single 581
gene-gene interaction *p*-value that also accounts for linkage 582
disequilibrium [42]. 583

It should be appreciated that these gene-based tests can be 584
applied to the setting of rare variant which are increasingly being 585
studied. Extensions to rare variant tests to incorporate gene-gene 586
and gene-environment interactions have been developed and 587
include SKAT [43] which can handle gene-gene interactions and 588
iSKAT [44] and rareGE [45] for gene-environment interactions. 589
590

*3.2.9 Replication*

Regardless of the analysis approach taken to identify gene-gene 591
interactions, the gold standard is to conduct a replication analysis 592
using an independent dataset. While fairly common for single- 593
variant association studies, this is less routinely followed for gene- 594
gene interaction analyses, as demonstrated in our systematic review 595
of asthma gene-gene interactions where only 15.2% of interactions 596
were attempted to be replicated [46]. The challenge is often iden- 597
tifying an appropriate replication dataset in which both variants 598
were genotyped. Through the use of imputation the variant to be 599
replicated could be imputed if they were not directly in the inde- 600
pendent dataset. 601

A challenge for any replication of a genetic effect is the direc- 602
tionality of the effect. Differences due to the populations selected 603
that can alter the minor allele frequencies and linkage disequilib- 604
rium structure can result in differences in both the magnitude and 605
direction of effect [47]. This is amplified when looking at two or 606
more loci as the probability of subtle differences can result in 607
differences in the direction of effect when looking at interactions. 608
609

*3.2.10 Meta-Analysis*

This data analysis technique is commonly used to pool data across 610
multiple studies and increase evidence for an association with a 611
genetic variant. Meta-analysis is an attractive analytical technique 612
because it can be used to increase power by substantially increasing 613
the total sample size. While relatively straightforward for single- 614
variant analyses using either fixed or random effects models [48] or 615
by combining *p*-values [49], this is not always the case for gene- 616
gene and gene-environment interactions. Differences in the analyt- 617
ical strategy and the way in which the results are presented may 618
make meta-analyses more challenging. The *p*-value approach does 619
not take into account the effect estimates, but this may be a better 620
first approach as it can be used on a much broader set of analytical 621
strategies such as MDR, regression-based approaches, and Random 622
Forrest methods, for example. Meta-analysis methods that account 623
for the effect estimate are more attractive because they can account 624
for heterogeneity in both the effect estimates and between 625
populations. 626
627

Andrew T. DeWan

*3.2.11 Case-Only Approach*

The case-only study design to detect interactions was first described for gene-environment interactions [50, 51]. This design is premised on the idea that the variant is independent of the environmental exposure in the population. In the presence of a gene-environment interaction, there would be an association between the variant and the environmental exposure among cases only and this can be most easily tested using the regression-based approaches described above (linear or logistic depending on the environmental exposure being investigated); however, it should be noted that when using the case-only design the main effects of the variant and the environmental exposure cannot be determined as this is only possible using a case-control design.

The case-only design is more powerful than the case-control design (Table 4; 5000 cases and 5000 controls for case-control and 5000 cases for case-only, main effects ORs of 1.2 for both the environmental factor and genetic variant, a frequency of the environmental exposure of 0.25 in the population and MAF of 0.2, population risk of disease of 0.1 and 200,000 genetic markers). The major caveat is that there is an assumption that there is independence between the genetic variant and the environmental factor, i.e., that there is no association between these factors among the source population. Violation of this assumption can lead to biased estimates of the OR and there is a recommendation that a case-only gene-environmental interaction analysis only proceed for environmental factors with population-specific data [52]. This independence may be difficult to establish in population-based data, so independence could be tested for among controls (e.g., when a

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

**Table 4**
**Power of case-control and case-only gene-environment interaction analyses**

| intOR | Case-control | Case-only |
|---|---|---|
| 1.1 | 0.000 | 0.001 |
| 1.2 | 0.002 | 0.063 |
| 1.3 | 0.038 | 0.564 |
| 1.4 | 0.211 | 0.962 |
| 1.5 | 0.541 | 1.000 |
| 1.6 | 0.829 | 1.000 |
| 1.7 | 0.959 | 1.000 |
| 1.8 | 0.994 | 1.000 |
| 1.9 | 0.999 | 1.000 |
| 2.0 | 1.000 | 1.000 |

case-only analysis is being conducted among a case-control study) 655
or replication of the finding among an independent case-control 656
study. This is of particular importance when the environmental 657
factor could be under genetic influence such as BMI. The substan- 658
tial power gain of the case-only study makes this analysis strategy 659
relatively attractive, but this is somewhat diminished by the strong 660
independence assumption and the need to test this assumption in a 661
set of population-based samples. 662

The case-only approach can also be utilized to detect gene-gene 663
interactions [53]. In the presence of an interaction between two 664
genetic variants there will be an association between these genetic 665
variants among cases. This is most easily tested using a logistic 666
regression model using one of the genetic markers as the dependent 667
variable. As with case-only gene-environment interaction analyses, 668
only the interaction can be examined and not the main effects of the 669
genetic variants. Independence of the genetic variants in the source 670
population is assumed, but this is much easier to achieve given the 671
large amounts of genome-wide genotype data available for many 672
diverse populations (e.g., HapMap [54]) or the use of LD pruning 673
among controls, if available (*see* Subheading 3.1.4). 674

675

*3.2.12 Three or More Variants*

Gene-gene interactions are not limited to pairwise interactions. 676
Higher order interactions involving three or more variants have 677
been reported (e.g., renin-angiotensin system SNPs and hyperten- 678
sion [55]) and can be modeled in the regression framework. The 679
problem is that as higher order interactions are considered the data 680
become more and more sparse and the power to detect a significant 681
interaction decreases. One approach to detecting higher order 682
interactions would be to use a nonparametric, non-explicit method 683
such as MDR to screen for potential higher order combinations of 684
variants and then subsequently explicitly test for the interaction on 685
a subset of the best performing combinations using a regression- 686
based method. 687

688

*3.3 Gene-Environment Interaction Analysis*

At the core, the analysis of gene-environment interactions does not 689
fundamentally differ from gene-gene interactions. The analysis can 690
be conducted using the regression-based or nonparametric 691
approaches described above with the same caveats for testing for 692
explicit vs. non-explicit interactions in the dataset. Joint, or non-- 693
explicit, tests may be used to screen the data to detect potential 694
interactions, followed by explicit tests to determine if an interaction 695
exists between the identified variant and the environmental factor. 696
In the logistic regression framework, one of the genetic variant 697
variables is replaced by the environmental exposure of interest and 698
the significance of the interaction $\beta$ is tested in the model. Using the 699
non-explicit MDR approach, the environmental variable is entered 700
into the algorithm with the caveat that the environmental factor 701
must be categorical in order for the reduction algorithm to work. 702

The power to detect gene-environment interactions suffers from the same lack of power on the genome-wide scale as gene-gene interactions. More than 30,000 cases and controls are required to detect an interaction OR of 1.2 (Table 1) under similar parameters as the gene-gene interactions analysis (main effects of OR = 1.2 for gene and environmental factors, MAF of 0.2, environmental prevalence of 0.25, population risk of disease of 0.1 and 200,000 genetic markers tested). Using 5000 cases and 5000 controls, the smallest detectable OR is just under 1.6 (Table 4). The same approaches to increase power by reducing the search space can be applied to the search for gene-environment interaction analyses, through all of the data reduction techniques described for gene-gene interactions.

It cannot be stressed enough that replication is the key to describing true positive gene-environment interactions, as was discussed for gene-gene interactions. As mentioned previously for gene-gene interaction replication, if the exact genetic variants are not directly genotyped in the independent dataset, imputation can be used to estimate the genotypes to be replicated. While this will work for the genetic variants in gene-environment interactions to be replicated, this is not the case if there is not comparable environmental data in the independent dataset. This can make the identification and selection of an appropriate dataset for replication more challenging for gene-environment interactions, but makes the replication of the findings no less important.

### 3.3.1 Modifiable Environmental Factors

One attractive aspect of identifying and describing gene-environment interactions for a complex disease is that this gives us the possibility of potentially modifying one contributing factor to a disease. At this point in time, inherited genetic variants are not modifiable, but if, for example, a significant gene-environment interaction were identified for obesity with a genetic variant and high saturated fat diet, individuals carrying the risk variant could be more strongly encouraged to reduce their saturated fat intake. Despite the general benefit we could all gain from reducing our saturated fat intake, this may be more effective if it were targeted to individuals based on their genetic profile and increased risk for obesity when both factors are present (beyond the additive main effects).

### 3.3.2 Gene-Gene-Environment Interactions

We should not think of gene-gene and gene-environment interactions as being mutually exclusive. As with higher order gene-gene interactions, interactions involving multiple genetic variants and an environmental factor are possible to model. The same issues with power due to sparsity of data and the number of unique combinations of factors apply, but given sufficient sample size these types of interaction models can be tested. It may be more important in this setting to attempt a screening step using a non-explicit approach and then apply an explicit test only to the set of interactions that surpass an initial significance threshold.

*3.4 Conclusions*

I have outlined several methods for conducting interaction analysis 753
to detect both gene-gene and gene-environment interactions. 754
However, it should be clear that there is no optimal method to 755
detect either type of interaction. The method(s) chosen are often 756
dictated based on the type of data you have available (e.g., case-757
control, case-only), the number of markers you have genotyped, 758
and the number of subjects you have included. I strongly recom-759
mend reporting the results of explicit tests for interaction as this will 760
greatly improve ability of other groups to attempt to replicate your 761
results and meta-analyze where appropriate, but the use of 762
non-explicit tests for interaction can be extremely useful to initially 763
screen large numbers of interactions and when sample sizes are 764
limited. We must continue to invest time and resources into identi-765
fying interactions in genome-wide data as this will help us to fill in 766
the missing heritability gap and better understand the genetic 767
architecture of complex traits. 768

769

## 4  Notes

770

1. Inheritance models: It should be noted that as with single-771
variant approaches, inheritance models can be imposed on the 772
variants that include additive, multiplicative, dominant, reces-773
sive, and overdominant. These inheritance models can be 774
imposed on the variants in the interaction model independently 775
and can be considered in a combinatorial fashion for variant 776
1 and variant 2 (e.g., additive × additive, additive × multiplica-777
tive, additive × dominant, etc.). However, if this is done, care 778
must be taken to account for this additional multiple testing. 779
These models can be problematic when they are misspecified as 780
they can reduce power which makes nonparametric approaches 781
attractive. 782

2. Scale: The scale on which the outcome is measured or evaluated 783
can influence whether or not an interaction exists between two 784
variants or between a variant and an environmental factor. It 785
needs to be kept in mind that there may be a monotone trans-786
formation of the outcome that could remove the interaction. 787
For example, on the odds ratio scale, an interaction may exist 788
between two SNPs (coded dichotomously as in a dominant 789
model) in which the OR for having dominant alleles at both 790
SNPs is greater than the sum of the ORs for having one domi-791
nant allele at each SNP. However, on the log(OR) scale this 792
interaction is removed and termed a removable interaction. If 793
there is no monotone transformation that can remove this inter-794
action it is termed essential. The method to detect these types of 795
interactions is described in a paper by Wu et al. [56]. 796

Andrew T. DeWan

3. Reporting of interaction results: As was previously outlined in 797
our paper [46], there are several recommendations for how 798
results of interaction analyses are reported in order to increase 799
the interpretability and replicability of the interaction. Effect 800
estimates should be provided so that both the strength and 801
direction of the interaction can be assessed. If a regression- 802
based approach is used, parameter estimates of the main effects 803
and the interaction term should be provided. If a nonparametric 804
approach is used, such as MDR, effect estimates are not pro- 805
duced, however, counts of cases and controls for the contin- 806
gency table of genotype combinations should be provided. This 807
will allow for a better assessment of the interaction and the 808
possibility to incorporate the data into a meta-analysis. 809

## References 810

1. Niel C, Sinoquet C, Dina C et al (2015) A 812
survey about methods dedicated to epistasis 813
detection. Front Genet 6:285 814

2. Ritchie MD (2015) Finding the epistasis nee- 815
dles in the genome-wide haystack. Methods 816
Mol Biol 1253:19–33 817

3. Gusareva ES, Van Steen K (2014) Practical 818
aspects of genome-wide association interaction 819
analysis. Hum Genet 133(11):1343–1358 820

4. Tiret L (2002) Gene-environment interaction: 821
a central concept in multifactorial diseases. 822
Proc Nutr Soc 61(4):457–463 823

5. Ottman R (1990) An epidemiologic approach 824
to gene-environment interaction. Genet Epi- 825
demiol 7(3):177–185 826

6. Manolio TA, Collins FS, Cox NJ et al (2009) 827
Finding the missing heritability of complex dis- 828
eases. Nature 461(7265):747–753 829

7. Bateson W (1909) Mendel's principles of 830
heredity. Cambridge University Press, 831
Cambridge 832

8. Cordell HJ (2002) Epistasis: what it means, 833
what it doesn't mean, and statistical methods 834
to detect it in humans. Hum Mol Genet 11 835
(20):2463–2468 836

9. Moore JH (2005) A global view of epistasis. 837
Nat Genet 37(1):13–14 838

10. Ma J, Thabane L, Beyene J et al (2016) Power 839
analysis for population-based longitudinal 840
studies investigating gene-environment inter- 841
actions in chronic diseases: a simulation study. 842
PLoS One 11(2):e0149940 843

11. Dunham I, Kundaje A, Aldred SF et al (2012) 844
An integrated encyclopedia of DNA elements 845
in the human genome. Nature 489 846
(7414):57–74 847

12. Wang K, Li M, Hakonarson H (2010) ANNO- 848
VAR: functional annotation of genetic variants 849
from high-throughput sequencing data. 850
Nucleic Acids Res 38(16):e164 851

13. Bush WS, Dudek SM, Ritchie MD (2009) Bio- 852
filter: a knowledge-integration system for the 853
multi-locus analysis of genome-wide associa- 854
tion studies. Pac Symp Biocomput:368–379 855 AU4

14. Price AL, Patterson NJ, Plenge RM et al 856
(2006) Principal components analysis corrects 857
for stratification in genome-wide association 858
studies. Nat Genet 38(8):904–909 859

15. Howie BN, Donnelly P, Marchini J (2009) A 860
flexible and accurate genotype imputation 861
method for the next generation of genome- 862
wide association studies. PLoS Genet 5(6): 863
e1000529 864

16. Purcell S, Neale B, Todd-Brown K et al (2007) 865
PLINK: a tool set for whole-genome associa- 866
tion and population-based linkage analyses. 867
Am J Hum Genet 81(3):559–575 868

17. Ueki M, Cordell HJ (2012) Improved statistics 869
for genome-wide interaction analysis. PLoS 870
Genet 8(4):e1002625 871

18. Wu X, Dong H, Luo L et al (2010) A novel 872
statistic for genome-wide interaction analysis. 873
PLoS Genet 6(9):e1001131 874

19. Wan X, Yang C, Yang Q et al (2010) BOOST: a 875
fast approach to detecting gene-gene interac- 876
tions in genome-wide case-control studies. Am 877
J Hum Genet 87(3):325–340 878

20. Hahn LW, Ritchie MD, Moore JH (2003) 879
Multifactor dimensionality reduction software 880
for detecting gene-gene and gene-environment 881
interactions. Bioinformatics 19(3):376–382 882

21. Ritchie MD, Hahn LW, Roodi N et al (2001) 883
Multifactor-dimensionality reduction reveals 884
high-order interactions among estrogen- 885
metabolism genes in sporadic breast cancer. 886
Am J Hum Genet 69(1):138–147 887

22. Calle ML, Urrea V, Malats N et al (2010) mbmdr: an R package for exploring gene-gene interactions associated with binary or quantitative traits. Bioinformatics 26 (17):2198–2199

23. Gui J, Moore JH, Williams SM et al (2013) A simple and computationally efficient approach to multifactor dimensionality reduction analysis of gene-gene interactions for quantitative traits. PLoS One 8(6):e66545

24. Van der Auwera GA, Carneiro MO, Hartl C et al (2013) From FASTQ data to high confidence variant calls: the genome analysis toolkit best practices pipeline. Curr Protoc Bioinformatics 43:11.10 1–11.1033

25. Li H, Durbin R (2010) Fast and accurate long-read alignment with Burrows-Wheeler transform. Bioinformatics 26(5):589–595

26. Li H, Handsaker B, Wysoker A et al (2009) The sequence alignment/map format and SAMtools. Bioinformatics 25(16):2078–2079

27. McKenna A, Hanna M, Banks E et al (2010) The genome analysis toolkit: a MapReduce framework for analyzing next-generation DNA sequencing data. Genome Res 20 (9):1297–1303

28. DePristo MA, Banks E, Poplin R et al (2011) A framework for variation discovery and genotyping using next-generation DNA sequencing data. Nat Genet 43(5):491–498

29. Dewan AT, Egan KB, Hellenbrand K et al (2012) Whole-exome sequencing of a pedigree segregating asthma. BMC Med Genet 13 (1):95

30. Marchini J, Donnelly P, Cardon LR (2005) Genome-wide strategies for detecting multiple loci that influence complex diseases. Nat Genet 37(4):413–417

31. Calle ML, Urrea V, Vellalta G, Malats N, Steen KV (2008) Improving strategies for detecting genetic patterns of disease susceptibility in association studies. Stat Med 27 (30):6532–6546

32. Benjamini Y, Hochberg Y (1995) Controlling the false discovery rate: a practical and powerful approach to multiple testing. J R Stat Soc B 57 (1):289–300

33. Nyholt DR (2004) A simple correction for multiple testing for single-nucleotide polymorphisms in linkage disequilibrium with each other. Am J Hum Genet 74(4):765–769

34. North BV, Curtis D, Sham PC (2002) A note on the calculation of empirical P values from Monte Carlo procedures. Am J Hum Genet 71 (2):439–441

35. North BV, Curtis D, Sham PC (2003) A note on calculation of empirical P values from Monte Carlo procedure. Am J Hum Genet 72 (2):498–499

36. Murk W, DeWan AT (2016) Exhaustive genome-wide search for SNP-SNP interactions across 10 human diseases. G3 (Bethesda) 6 (7):2043–2050

37. Gauderma WJ, Morrison JM, QUANTO 1.1: A computer program for power and sample size calculations for genetic-epidemiology studies. http://hydra.usc.edu/gxe2006

38. Uzun A, Sharma S, Padbury J (2012) A bioinformatics approach to preterm birth. Am J Reprod Immunol 67(4):273–277

39. Uzun A, Triche EW, Schuster J et al (2016) dbPEC: a comprehensive literature-based database for preeclampsia related genes and phenotypes. Database (Oxford). https://doi.org/10.1093/database/baw006. pii:baw006

40. Shearer AE, Eppsteiner RW, Booth KT et al (2014) Utilizing ethnic-specific differences in minor allele frequency to recategorize reported pathogenic deafness variants. Am J Hum Genet 95(4):445–453

41. Murk W, DeWan AT (2016) Genome-wide search identifies a gene-gene interaction between 20p13 and 2q14 in asthma. BMC Genet 17(1):102

42. Ma L, Clark AG, Keinan A (2013) Gene-based testing of interactions in association studies of quantitative traits. PLoS Genet 9(2):e1003321

43. Wu MC, Lee S, Cai T et al (2011) Rare-variant association testing for sequencing data with the sequence kernel association test. Am J Hum Genet 89(1):82–93

44. Lin X, Lee S, Wu MC et al (2016) Test for rare variants by environment interactions in sequencing association studies. Biometrics 72 (1):156–164

45. Chen H, Meigs JB, Dupuis J (2014) Incorporating gene-environment interaction in testing for association with rare genetic variants. Hum Hered 78(2):81–90

46. Murk W, Bracken MB, DeWan AT (2015) Confronting the missing epistasis problem: on the reproducibility of gene-gene interactions. Hum Genet 134(8):837–849

47. Greene CS, Penrod NM, Williams SM et al (2009) Failure to replicate a genetic association may provide important clues about genetic architecture. PLoS One 4(6):e5639

48. Fleiss JL (1993) The statistical basis of meta-analysis. Stat Methods Med Res 2(2):121–145

49. Fisher RA (1948) Combining independent tests of significance. Am Stat 2:30

50. Piegorsch WW, Weinberg CR, Taylor JA (1994) Non-hierarchical logistic models and

Andrew T. DeWan

1001 case-only designs for assessing susceptibility in
1002 population-based case-control studies. Stat
1003 Med 13(2):153–162

1004 51. Begg CB, Zhang ZF (1994) Statistical analysis
1005 of molecular epidemiology studies employing
1006 case-series. Cancer Epidemiol Biomark Prev 3
1007 (2):173–175

1008 52. Hodgson ME, Olshan AF, North KE et al
1009 (2012) The case-only independence assump-
1010 tion: associations between genetic polymorph-
1011 isms and smoking among controls in two
1012 population-based studies. Int J Mol Epidemiol
1013 Genet 3(4):333–360

1014 53. Yang Q, Khoury MJ, Sun F et al (1999) Case-
1015 only design to measure gene-gene interaction.
1016 Epidemiology 10(2):167–170

1017 54. The International HapMap Consortium
1018 (2003) The international HapMap project.
1019 Nature 426:789–796

1020 55. Yang CH, Lin YD, Wu SJ et al (2015) High
1021 order gene-gene interactions in eight single
1022 nucleotide polymorphisms of renin-
1023 angiotensin system genes for hypertension
1024 association study. Biomed Res Int
1025 2015:454091

1026 56. Wu C, Zhang H, Liu X et al (2009) Detecting
1027 essential and removable interactions in
1028 genome-wide association studies. Stat Inter-
1029 face 2(2):161–170

# Author Queries

Chapter No.: 7      394545_1_En

| Query Refs. | Details Required | Author's response |
|---|---|---|
| AU1 | Please check whether the affiliation and correspondence details are presented correctly. | |
| AU2 | Please check the hierarchy of the section headings and confirm if correct. | |
| AU3 | Please check the sentence starting: "Then converted to BAM format, sorted, indexed, PCR duplicates marked and then…" for completeness. | |
| AU4 | Please provide volume number for Ref. [13]. | |